# BANGALORE INSTITUTE OF TECHNOLOGY

K.R. Road, V.V.Puram, Bengaluru-560 004

## DEPARTMENT OF COMPUTER SCIENCE & ENGG

## SYSTEM SOFTWARE AND COMPILER DESIGN

## NOTES

## SUBJECT CODE: 15CS63

### By

**Mrs. Hemavathi. P**
**Assistant Professor**
**Department of CSE**

## SYSTEM SOFTWARE AND COMPILER DESIGN
### [As per Choice Based Credit System (CBCS) scheme]
### (Effective from the academic year 2016 -2017)
### SEMESTER – VI

| Subject Code | 15CS63 | IA Marks | 20 |
|---|---|---|---|
| Number of Lecture Hours/Week | 4 | Exam Marks | 80 |
| Total Number of Lecture Hours | 50 | Exam Hours | 03 |

### CREDITS – 04

**Course objectives:** This course will enable students to

- Define System Software such as Assemblers, Loaders, Linkers and Macroprocessors
- Familiarize with source file, object file and executable file structures and libraries
- Describe the front-end and back-end phases of compiler and their importance to students

| Module – 1 | Teaching Hours |
|---|---|
| Introduction to System Software, Machine Architecture of SIC and SIC/XE. **Assemblers:** Basic assembler functions, machine dependent assembler features, machine independent assembler features, assembler design options. **Macroprocessors:** Basic macro processor functions, **Text book 1: Chapter 1: 1.1,1.2,1.3.1,1.3.2, Chapter2 : 2.1-2.4,Chapter4: 4.1.1,4.1.2** | **10 Hours** |
| **Module – 2** | |
| **Loaders and Linkers:** Basic Loader Functions**,** Machine Dependent Loader Features, Machine Independent Loader Features, Loader Design Options, Implementation Examples. **Text book 1 : Chapter 3 ,3.1 -3.5** | **10 Hours** |
| **Module – 3** | |
| **Introduction:** Language Processors, The structure of a compiler, The evaluation of programming languages, The science of building compiler, Applications of compiler technology, Programming language basics **Lexical Analysis:** The role of lexical analyzer, Input buffering, Specifications of token, recognition of tokens, lexical analyzer generator, Finite automate. **Text book 2:Chapter 1  1.1-1.6   Chapter 3    3.1 – 3.6** | **10 Hours** |
| **Module – 4** | |
| Syntax Analysis: Introduction, Role Of Parsers, Context Free Grammars, Writing a grammar, Top Down Parsers, Bottom-Up Parsers, Operator-Precedence Parsing **Text book 2: Chapter 4   4.1 4.2 4.3 4.4 4.5 4.6      Text book 1 : 5.1.3** | **10 Hours** |
| **Module – 5** | |
| Syntax Directed Translation, Intermediate code generation,  Code generation **Text book 2:  Chapter 5.1, 5.2, 5.3, 6.1, 6.2, 8.1, 8.2** | **10 Hours** |

**Course outcomes:** The students should be able to:

- Explain system software such as assemblers, loaders, linkers and macroprocessors
- Design and develop lexical analyzers, parsers and code generators
- Utilize lex and yacc tools for implementing different concepts of system software

**Question paper pattern:**

The question paper will have TEN questions.

There will be TWO questions from each module.

Each question will have questions covering all the topics under a module.

The students will have to answer FIVE full questions, selecting ONE full question from each module.

**Text Books:**

1. System Software by Leland. L. Beck, D Manjula, 3$^{rd}$ edition, 2012
2. Compilers-Principles, Techniques and Tools by Alfred V Aho, Monica S. Lam, Ravi Sethi, Jeffrey D. Ullman. Pearson, 2$^{nd}$ edition, 2007

**Reference Books:**

1. Systems programming – Srimanta Pal , Oxford university press, 2016
2. System programming and Compiler Design, K C Louden, Cengage Learning
3. System software and operating system by D. M. Dhamdhere TMG
4. Compiler Design, K Muneeswaran, Oxford University Press 2013.

## DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING

### COURSE OBJECTIVES AND OUTCOMES-2015-19

Course Title : System Software and Compiler Design     Course Code : : 15CS63

No. of Lecture Hrs./Week : 04     Exam Hours : 03

Total No. of Lecture Hrs. : 52     Exam Marks : 80

### Prerequisites

1. Microprocessors and Microcontrollers(15CS44)

2. Automata Theory and Computability (15CS54)

### Course Learning Objectives

This course will help students to achieve the following objectives:

1. To understand the concepts of System software, Application Software and different hypothetical machine architectures.

2. Familiarize with source file, symbol table creation (pass-1), object file creation (pass-2), loaders and linkers.

3. To know the fundamental concepts of translators.

4. To identify the methods and strategies for parsing techniques.

5. Devise and perform syntax-directed translation schemes for compiler.

6. Devise intermediate code generation schemes and analyze the optimized code generated after the synthesis phase.

### Course Outcomes

At the end of the course students should be able to:

1. Apply the knowledge of System Software such as Assemblers, Loaders, Linkers and Macro processors to build an application.

2. Understand the basic principles of compiler in high level programming language.

3. Analyze and design the analysis phase using different techniques.

4. Build the system software by associating synthesis phase with analysis phase for better optimization and performance.

MODULE-1

TEXTBOOK: System Software by Leland.L. Beck, D.Manjula, 3<sup>rd</sup> Edition, 2012

CHAPTER 1: **Introduction to System Software and Machine Architecture**

CHAPTER 2: **Assemblers**

CHAPTER 4: **Macro Processors**

# CHAPTER 1

## Introduction to System Software and Machine Architecture

The term "software" refers to the set of electronic program instructions or data a computer processor reads in order to perform a task.

"Hardware" refers to the physical components that you can see and touch, such as the computer hard drive, mouse and keyboard.



Fig: Relationship b/w system software, Appln software and Hardware

Defn: "System software" is a set of programs that are dedicated to manage the computer itself, such as operating system, file management utilities.

Application software are a set of productivity programs or end-user programs to perform their specific tasks.

# Difference between System software and Application software

| System software | Application software |
|---|---|
| 1. System software is a set of programs that are dedicated to manage the computer itself (mem. mgmt, process mgmt, protection security | Application software is a set of computer programs designed to permit the user to perform a group of functions, tasks or activities. |
| 2. Is written in a low-level language ie assembly language | Is written in a high-level language like C, C++, Java, .net, VB etc |
| 3. Starts running when the s/m is turned on and runs till the system is shut down | Runs as and when the user requests |
| 4. A system is unable to run without system software | Appln. software is even not required to run the system ∴ it is user specific |
| 5. System software is general purpose | Appln software is specific purpose software |
| 6. Ex: operating system, assembler, compiler, loader or linker, text editor, debugger, macro processors, | Ex: web browser, word processing, spreadsheet, database, Adobe creative suite, Audio master suite, games |
| 7. not machine dependent (machine architecture) | not machine dependent |

1a. System software and machine architecture

. Machine dependency of system software

→ System programs are intended to support the operation and use of the computer.

→ machine architecture differs in :

. machine code,

. instruction formats

. Addressing mode

. Registers

. machine independency of system software

→ General design and logic is basically the same:

. code optimization

. subprogram linking

## 1.3. Simplified Instructional Computer

As we know different systems have different features and different features are difficult to study one by one. so to avoid this problem we study Simplified Instructional Computer.

SIC is a hypothetical computer system introduced in system software. Due to the fact that most modern microprocessors include complex functions for the purpose of efficiency, it is very difficult to learn systems programming using a real-world system. The SIC solves this problem by abstracting away these complex behaviours in favour of an architecture that is clear and usable for those wanting to learn systems programming

. SIC comes in two versions

→ standard model

→ XE version ( Extra Equipment or Extra Expensive)

. The two versions has been designed to be upword compatible.

1.3 1    SIC machine architecture

Every machine architecture includes

a) memory

b) Registers

c) Data formats

d) Instruction formats

e) Addressing modes

f) Instruction set

g) Input and output


a) Memory

→ memory consists of 8 bit bytes

→ Any 3 consecutive bytes form a word (24 bits)

→ All addresses on SIC are byte addresses

→ words are addressed by the location of their lowest numbered byte

→ Total of 32,768 ($2^{15}$) bytes in the computer memory

∴ 15-line address bus


b) Registers

→ Five registers, all of which have special uses

→ Each register is 24 bits in length

→ table shows the mnemonic, number and use of each register.

| Mnemonic | Number | Use |
|---|---|---|
| A Accumulator | 0 | Used for arithmetic operations |
| X Index Register | 1 | Used for addressing |
| L Linkage Register | 2 | JSUB — Jump to subroutine Instruction stores return address |
| PC Program counter | 8 | Contain the address of the next instruction to be fetched for execution |
| SW Status word | 9 | Contains variety of information including a condition code in comp instruction |

c) Data formats

→ Integers are stored as 24-bit binary number.
$$( 0 - 2^{24} - 1 \Rightarrow 0 \text{ to } 16G - 1)$$

→ Negative value are represented as 2's complement
$$E_1: \quad -2h \text{ is represented as (8 bit representation)}$$

$$2h = 0 0 0 1 1 0 0 0$$

1's complement $\quad 1 1 1 0 0 1 1 1$

$$+ \quad \underline{\phantom{xxxxxx} 1}$$

$$1 1 1 0 1 0 0 0 \Rightarrow 232$$

→ characters are stored using their 8-bit ASCII codes

→ There is no floating-point hardware on the standard version of SIC

$$E_1:- \quad 5 = 0000\ 0000\ 0000\ 0000\ 0000\ 0101$$
$$-5 = 1111\ 1111\ 1111\ 1111\ 1111\ 1011$$
$$A = 0100\ 0001\ (65)$$

d) Instruction formats

→ All machine instructions on the standard version of SIC are have 24-bit format

| 8 | 1 | 15 |
|---|---|---|
| Opcode | x | Address |

x → indicates <u>indexed</u> <u>addressing</u> <u>mode</u>

e) Addressing modes

→ Two addressing modes based on **x** bit

• Direct Addressing

• Indexed addressing

| mode | Indication | Target address (TA) |
|---|---|---|
| Direct | x=0 | TA = address |
| Indexed | x=1 | TA = address + (x) <br> parenthesis indicate the content of a register or memory location |

→ Ex:- LDA TEN ; LDA = 00 (opcode)

| 8 | 1 | 15 |
|---|---|---|
| 0000 0000 | 0 | 001 0000 0000 0000 |
| opcode | x | address |

Target address / Effective address = 1000 ie contents of the address 1000 is loaded to accumulator

→ Indexed addressing mode

Ex: STCH   BUFFER , X     ;   opcode for STCH = 54

BUFFER = 1000

| 0101 | 0100 | 1 | 001 | 0000 | 0000 | 0000 |

  5      4     x        1   0    0    0

opcode                    address (BUFFER)

TA = address +(x)

= 1000 + content of the index register X

is the Accumulator content, the character is loaded to the effective address.


1) Instruction set

(i) load and store : LDA, LDX, STA, STX

(ii) Integer Arithmetic operations : ADD, SUB, MUL, DIV

• Arithmetic operation involve register A and a word in memory with the result being left in the register

Ex: ADD   WORD  ;  A ← A + WORD

is adds register A contents with WORD and result is stored in register A

(iii) Comparsion operations : COMP

• compares the value in register A with memory and sets a condition code (cc) accordingly

Ex: COMP  WORD ; compares A's contents with WORD and sets cc as < = >

(iv) Conditional jump instructions: JLT, JEQ, JGT
   - these instructions will tst the setting of CC and jump accordingly

(v) Subroutine linkage instructions ; JSUB, RSUB
   . JSUB — jumps to the subroutine by placing the return address in register L. (Fn. call)

   • RSUB — returns by jumping to the address contained in register L (Fn. return to the caller)

Eg: void main()
{

   Add (1,2);

   }
   Fn. call
   └► Add (int x, int y)
   { return (x+y);

   Fn. Return}



stack for main() → when it calls Add()

1000
base addr

stack frame for add()
Return address

bp for988  | 1000 |
add 992

bp→1000
for main

once call returns ie
Unwinding of stack happens it has to go back to caller routing base pointer of add point to main so 1000 is stored.

stack frame for main

1000
bp

(g) Input and output

→ Input and output is performed by transferring 1 byte data at a time to or from the rightmost 8 bits of register A.

→ Each device is assigned a unique 8-bit code

→ There are three I/o instruction which specifies the device code as on operand

(i) **TD (Test Device)** : checks whether the addressed device is ready to send or receive a byte of data. CC (condition code in sw register) is set according ( < = )

* < → device is ready to send/receive

= → device is not ready.

→ A program has to wait until the device is ready , then execute a Read Data (RD) and write data (WP) instructions.

→ This sequence should be continued for each byte of data . ( I/o).

→ RD : Transfers a byte of data from i/p device into rightmost byte of register A ( RD INDEV STA DATA )

→ WD : Byte of data is loaded into rightmost byte of reg A and then written to output device ( LDA DATA & WD OUTDEV )

Ex: 1) SIC instructions for data movement operations
(no memory-memory move instructions)

```
        LDA     FIVE        ; Load constant 5 into register A
        STA     ALPHA       ; Store in Alpha
        LDCH    CHAR2       , load character '2' into reg A
        STCH    C1          ; store in character variable C1
          .
          .
          .

ALPHA   RESW    1           ; one-word variable

FIVE    WORD    5           ; one-word constant

CHAR2   BYTE    C '2'       ; one-byte constant

C1      RESB    1           ; one-byte variable
```

⇓ some can be written as

```
        LDX     FIVE
        STX     ALPHA

            or

        LDL     FIVE
        STL     ALPHA
```

A — Accumulator
X — indexed register
L — Linkage register

2) SIC instructions for arithmetic operations

$$Delta = BETA = ALPHA + INCR - 1$$

$$DELTA = GAMMA + INCR - 1$$

| | | |
|---|---|---|
| LDA | ALPHA | ; loads Alpha into register A |
| ADD | INCR | ; $A \leftarrow (A) + (INCR)$ |
| SUB | ONE | ; $A \leftarrow (A) - 1$ |
| STA | BETA | ; $BETA \leftarrow (A)$ |
| LDA | GAMMA | ; $A \leftarrow (GAMMA)$ |
| ADD | INCR | ; $A \leftarrow (A) + (INCR)$ |
| SUB | ONE | ; $A \leftarrow (A) - 1$ |
| STA | DELTA | $DELTA \leftarrow (A)$ |
| ⋮ | | |

| | | |
|---|---|---|
| ONS | WORD | 1 |
| ALPHA | RESW | 1 |
| BETA | RESW | 1 |
| GAMMA | RESW | 1 |
| DELTA | RESW | 1 |
| INCR | RESW | 1 |

3) SIC instructions for looping and indexed operations
(program to copy 11-byte character string to another string)

```
// LDX  ZERO ;  initialize index register
         j=0;
. for (i=0; si[i] != '\0'; i++)
        s2[j++] = si[i];
    s2[j] = '\0';
```

```
        LDX    ZERO ;  initialize index register to 0

LOOP    LDCH   STR1, X ;  copies the first character q str1 to reg. A
                        ( TA :(address) → content q first byte
                                                q str1)

        STCH   STR2, X ;  store the first character into STR2

        TIX    ELEVEN ;  Add 1 to index, compare to 11
                          X = 0+1 = 1 ;  1 ⟷ 11  cc will
                                be set as <

        JLT    LOOP ;  repeat if index is < 11
                                    (x)
        .
        :

STR1    BYTE     C ' HELLO WORLD'
                            space

STR2    RESB     11

ZERO    WORD     0

ELEVEN  WORD     11
```

h) Program to add 2 arrays of 100 words each and store it in another array. Each word is 3 bytes.

$G = A + B$ ;

100 words : $3 \times 100 = 300$ bytes

```
ADDLOOP   LDX    INDEX ;  initialize index value X = 0

          LDA    ALPHA, X ;   A ← (ALPHA)

          ADD    BETA , X ;   A ← (A) + (BETA)  at 0th byte (index)

          STA    GAMMA, X ;   G ← (A) at 0th byte (index value)

          LDA    INDEX    ;   A ← 0

          ADD    THREE    ;   A ← (A) + 3 = 3 → m

          STA    INDEX    ;   INDEX = 3

          COMP   K300     ;   A ⟷ k300  ie  3 ⟷ 300  cc = <

          JLT    ADDLOOP  ;   repeat loop, now x = 3rd byte
```

```
k300     WORD    300        THREE    WORD   3
INDEX    RESW    1          BETA     RESW   100
ALPHA    RESW    100        GAMMA    RESW   100
```

5) To read one byte of data from input device and copies it
to device os

```
INLOOP    TD      INDEV       ; Test input device
          JEQ     INLOOP      ; cc := then loop until device ready
          RD      INDEV       ; once ready, read a byte into reg A
          STCH    DATA        ; store it in data (memory)
          .
          .
OUTLOOP   TD      OUTDEV      ; Test output device
          JEB     OUTLOOP     ; cc := then loop until device ready
          LDCH    DATA        ; load data byte into reg A
          WD      OUTDEV      ; write one byte to output device
          .
          .
INDEV     BYTE    X 'F1'
OUTDEV    BYTE    X '05'
DATA      RESB    1
```

6) Subroutine call to read a 100-byte record from an input
device into memory.

```
          JWB     READ        ; call Read subroutine wherein
                              ; it stores the return address in
          .                   ; linkage register
          .
READ      LDX     ZERO        ; X ← 0
RLOOP     TD      INDEV       ; Test input device
          JEQ     RLOOP       ; cc:=, loop until device is ready
          RD      INDEV       ; read one byte into reg A
          STCH    RECORD,X    ; store it into RECORD at Xth addr
          TIX     k100        ; X = (X)+1 & 1↔100 compare
          JLT     RLOOP       ; cc:< then loop back
          RSUB                ; Exit from subroutine; it
          .                   ; returns to the address stored in
                              ; linkage register
```

```
        :
        :

INDEX     BYTE      X 'F'

RECORD    RESB      100

ZERO      WORD      0

KICC      WORD      100
```

**1.3.2 : SIC/XE machine Architecture**

→ SIC/XE : Simple Instructional Computer with Extra Equipment

a) memory

b) Registers

c) Data formats

d) Instruction formats

e) Addressing modes

f) Instruction set

g) Input and output

a) **memory**

→ memory consists of 8 bit bytes

→ 3 consecutive bytes form a word (24 bits)

→ All addresses are byte addresses

→ words are addressed by the location of their lowest numbered byte

→ Total of 1 MB ($2^{20}$ bytes) in the memory. (20 bit address bus) which leads to change in instruction formats and addressing modes.

b) **Registers**

→ There are 9 registers

→ Each register is 24 bits in length except Floating Point reg

→ The registers are A, X, L, B, S, T, F, PC & SW

| Mnemonic | Number | Uses |
|---|---|---|
| A<br>Accumulator | 0 | Used for arithmetic operations |
| X<br>Index Register | 1 | Used for addressing (indexed) |
| L<br>Linkage register | 2 | used to store the return address for JSUB instruction |
| B<br>Base register | 3 | used for addressing |
| S<br>General Register | 4 | General working register – no special use |
| T<br>General Register | 5 | General working register - no special use |
| F<br>Floating Point Accumulator | 6 | ~~General working register~~ Floating point accumulator (48 bit) |
| PC<br>Program Counter | 8 | Contains the address of the next instruction to be fetched for execution |
| SW<br>status word | 9 | Contains a variety of information including a condition code (cc) |

c) Data formats:

→ Integers are stored as 24-bit binary numbers

→ negative values are represented as 2's complement (∮ 1's complement +1)

→ characters are stored using their 8-bit ASCII codes

→ There is a 48 bit floating point data type

| 1 | 11 | 36 |
|---|---|---|
| S | EXPONENT | FRACTION |

ns bit

→ The fraction is interpreted as a value between 0 and 1

→ The assumed binary point is immediately before the higher order bit

→ For normalized floating point numbers, the higher order bit of the fraction must be 1

→ The exponent is interpreted as an unsigned binary number between 0 and 8047 $(0 - (2^h - 1))$

→ If the expression has value e, fraction f and the absolute value of number is represented is

$$f \ast 2^{(e - 1024)}$$

→ The sign of floating point number is indicated by s ( s = 0 (+ve) and 1 (-ve))

Ex:  5 = 0000 0000 0000 0000 0000 0101

   -5 = 1111 1111 1111 1111 1111 1011

   AB = 0100 0001 (65)

Ex: 4.89 representation

As we know from computer organization it is
represented as $\pm m B^{\pm E}$
       ↓ └→ Base(2)
    Fraction (mantissa)

1) Represent 4 in binary form → 100

2) Convert 0.89 into binary form until it repeats
or until we get 36 bits which represents
the fraction part

   100.111 0001111010111000001010 0011101011
     000010100

3) normalization has to be done but not always
∵ They have specified that binary point
is immediately before the higher order bit

i.e ↖ 100. 111000 · · · ·

0.100111000111101011100001010001111010 X $2^3$ → Exponent
       fraction

Note: for normalized floating point number it will
be as 1.00111000111111 · · · · X $2^2$

i.e $\pm * 2^{e \pm 1024} = 0.10011 \cdots$ X $2^{3+1024}$
      $= 0.100111 \cdots$ X $2^{1027}$

$(1027)_2 \to 10000000011$

Right column:

0.89X2 →
0.78X2 → 1 i.e 1.78
0.56X2 → 1
0.12X2 → 1
0.24X2 → 0
0.48X2 → 0
0.96X2 → 0
0.92X2 → 1
0.84X2 → 1
0.68X2 → 1
0.36X2 → 1
0.72X2 → 0
0.44X2 → 1
0.88X2 → 0
0.76X2 → 1
0.52X2 → 1
0.04X2 → 1
0.08X2 → 0
0.16X2 → 0
0.32X2 → 0
0.64X2 → 0
0.28X2 → 1
0.56X2 → 0
0.12X2 → 1
0.24X2 → 0
0.48X2 → 0
0.96X2 → continue as above

0111101011100001010
0

| 11 bit | 36 bit | |
|---|---|---|
| 0 | 10000000011 | 100111000111101011100010100011101010 |
| S | exponent | fraction |

9) Represent $\overset{s=1(negative)}{\underset{-0.000489}{\nearrow}}$ in binary format

Given $-0.000489$

↦ Represent $0$ as binary $0$

→ Represent $0.000489$ as binary

$= 0.0000000000100000000011000000$
$11110000000011111110.$
$\underbrace{\quad}_{fraction}$

$= .10000000000110000001111100000001$
$\underbrace{1111110 \times 2^{-10}}$

⇒ $\frac{1}{f} \times 2^{e+1024} = .10000\cdots \times 2^{-10+1024}$

$= .10000\cdots \times 2^{\overset{1014 \to E}{}}$    $(1014) = 01111110110$

$\underbrace{}_{fraction}$

| | .000489 × 2 |
|---|---|
| | 0.000978 × 2 → 0 |
| | 0.001956 × 2 → 0 |
| | 0.003912 × 2 → 0 |
| | 0.007824 × 2 → 0 |
| | 0.015648 × 2 → 0 |
| | 0.031296 × 2 → 0 |
| | 0.062592 × 2 → 0 |
| | 0.125184 × 2 → 0 |
| | 0.250368 × 2 → 0 |
| | 0.500736 × 2 → 0 |
| | 0.001472 × 2 — 01 |
| | 0.002944 × 2 → 0 |

| 11 | | 36 | |
|---|---|---|---|
| 1 | 01111110110 | 10000000000110000001111000000011111 | |
| s | Exponent | fraction | 10 |

d) Instruction formats.

→ Since the memory used by SIC/XE may be $2^{20}$ bytes, the instruction format of SIC is not enough.

→ There are two possible options
  (i) Either use some form of relative addressing
  (ii) Extend the address field to 20 bits

→ if $e = 0$, then format 3
→ if $e = 4$, then format 4

1) Format 1 (1 byte)

```
      8 bit
  ┌──────────┐
  │  opcode  │
  └──────────┘
```

2) Format 2 (2 bytes)

```
     8        4       4    bits
  ┌────────┬──────┬──────┐
  │ opcode │  r₁  │  r₂  │
  └────────┴──────┴──────┘
```

3) Format 3 (3 bytes)

```
    6      1 1 1 1 1 1      12        bits
  ┌──────┬─┬─┬─┬─┬─┬─┬──────────┐
  │opcode│n│i│x│b│p│e│   disp   │
  └──────┴─┴─┴─┴─┴─┴─┴──────────┘
```

Note: e = 0

4) Format 4 (4 bytes)

```
    6      1 1 1 1 1 1        20         bits
  ┌──────┬─┬─┬─┬─┬─┬─┬──────────────┐
  │opcode│n│i│x│b│p│e│   address    │
  └──────┴─┴─┴─┴─┴─┴─┴──────────────┘
```

Ex: RSUB (return to subroutine) → 4C

⟹ it returns to the address stored in linkage register

```
  ┌───────────┐
  │ 0100 1100 │
  └───────────┘
      4    C    ⟶ object code
```

Ex: COMPR A,S (compare the contents of registers A & S)

opcode g  COMPR = A0

```
     8              4       4
  ┌───────────┬──────────┬──────┐
  │ 1010 0000 │0000      │0100  │
  └───────────┴──────────┴──────┘
      A         0      0    4 → obj.code
```

Ex: LDA #3 (load 3 to A)

```
     6     1 1 1 1 1 1         12
  ┌────────┬─┬─┬─┬─┬─┬─┬───────────────────┐
  │000000  │0│1│0│0│0│0│0000 0000 0011      │
  └────────┴─┴─┴─┴─┴─┴─┴───────────────────┘
  opcode    n i x b p e   0    0    3
```

0 1 0 0 3 ⟶ object code

Ex: +JSUB RDREC (Jump to address 1036)

opcode JSUB—48

```
     6      1 1 1 1 1 1          20
  ┌────────┬─────────┬────────────────────────┐
  │010010  │11 0001  │0000 0001 0000 0011 0110 │
  └────────┴─────────┴────────────────────────┘
  opcode   n i x b p e         addrew
     4  B      1     0    1    0   3  6
```

object code is 4B101036

e) Addressing modes

| | MODE | INDICATION | TARGET ADDRESS CALCULATION |
|---|---|---|---|
| 1. | Base Relative | $b=1, p=0$ | $TA = (B) + \text{displacement}$ <br> $(0 \leq disp \leq 4095)$ |
| 2. | Program Counter Relative | $b=0, p=1$ | $TA = (P) + \text{displacement}$ <br> $(-2048 \leq disp \leq 2047)$ |
| 3. | Direct Addressing | $b=0, p=0$ <br> (for format 3) | $TA = \text{displacement}$ |
| | | $b=0, p=0$ <br> (for format 4) | $TA = \text{Address field}$ |
| 4. | Base Relative Indexed addressing | $b=1, p=0$ <br> $x=1$ | $TA = (B) + (x) + \text{displacement}$ |
| 5. | Program Relative Indexed addressing | $b=0, p=1$ <br> $x=1$ | $TA = (P) + (x) + \text{displacement}$ |
| 6. | Immediate addressing | $i=1, n=0$ | Target address itself used as <br> $TA = \text{operand value}$ <br> (no memory reference) |
| 7 | Indirect addressing | $i=0, n=1$ | $TA = \text{displacement value}$ |
| 8. | Simple addressing | $i=0, n=0$ <br> OR <br> $i=1, n=1$ | $TA = \text{location of operand}$ |

Note:
→ Format 3 :
  ↳ in Base relative, disp is interpreted as 12 bit unsigned integer (1)
  ↳ in Pc relative, disp is interpreted as 12-bit signed integer & negative numbers if TA is complete (2)

Special symbols indication

1) #  :  Immediate address

2) @  :  Indirect address

3) +  :  Format 4

4) *  :  The current value of PC

5) c' ' :  character string

6) op m, x :  x - denotes the index register

7) base  :  Base - register


*) Instruction set :

Note:  Immediate addressing ( i=1, n=0) $\rightarrow$ Target address
 itself is used  as the operand value (no memory
 reference is performed )

$\rightarrow$  Indirect addressing (i=0, n=1); $\rightarrow$ the word at the
 location given by the target address is fetched and
 the value contained in this word is then taken
 as the address of the operand value.

$\rightarrow$  Indexing  cannot be used with immediate or
 indirect addressing mode.

** $\rightarrow$  we cant set both b=1 & p=1 which is invalid
 instruction set .

Examples of SIC/XE instructions and addressing modes

(B) = 006000
(PC) = 003000
(X) = 000090

Machine instruction

| Hex | opcode | n | i | x | b | p | e | disp/address | Target address | Value loaded into Reg A | Mode |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 032600 | 0000 00 | 1 | 1 | 0 | 0 | 1 | 0 | 0110 0000 0000 | 3600 | 103000 | Program Counter Relative. TA = (PC) + displacement = 003000 + 600 = 3600 |
| 03C300 | 0000 00 | 1 | 1 | 1 | 1 | 0 | 0 | 0011 0000 0000 | 6390 | 00C303 | Base relative Indexed. TA = (B)+(X)+disp = 006000 + 000090 + 300 = 6390 |
| 022030 | 0000 00 | 1 | 0 | 0 | 0 | 1 | 0 | 0000 0011 0000 | 3030 (contains 3600) | 103000 | Indirect + Program relative. TA = (PC) + disp = 003000 + 030 = 3030 |
| 010030 | 0000 00 | 0 | 1 | 0 | 0 | 0 | 0 | 0000 0011 0000 | 30 | 000030 | Immediate addressing. TA is used as operand value |
| 003600 | 0000 00 | 0 | 0 | 0 | 0 | 1 | 1 | 0110 0000 0000 | 3600 | 103000 | PC relative. TA = (PC)+disp = 003000 + 600 = 3600 |
| 0310C303 | 0000 00 | 1 | 1 | 0 | 0 | 0 | 1 | 1100 0011 0000 0000 0011 | C303 | 003030 | Simple addressing. TA = location of operand |

| Address | Contents |
|---|---|
| | ⋮ |
| 3030 | 003600, |
| | ⋮ |
| 3600 | 10 3000 |
| | ⋮ |
| 6390 | 00 C303 |
| | ⋮ |
| C303 | 0 03030 |
| | ⋮ |

(B) = 006000

(PC) = 003000

(x) = 000090

Fig: contents of registers
B, PC and x & memory
locations

1> Instruction set

* → Load and store instruction: LDA, LDX, STA, STX, LDB, STB

* → Integer and Floating point arithmetic operations:
    ADD, SUB, MUL, DIV, ADDF, SUBF, MULF, DIVF

* → Register move instructions (Rmo) ⟹ register-to-register
    operations such as ADDR, SUBR, MULR, DIVR

* → A special supervisor call (svc) instruction is provided.
    Executing this instruction generates an interrupt that
    can be used for communication with the operating system.

    → Comparision instruction: COMP, COMPR, COMPF

    → Conditional jump instruction: JLT, JEQ, JGT

    → Subroutine linkage instruction: JSUB, RSUB

2> Input and output

    → Input and output is performed by transferring 1 byte of
    data at a time to or from, the rightmost 8 bits of
    register A.

    → Each device is assigned a unique 8-bit code

    → Three I/O instructions which specifies the device code
    as an operand

        (i) TD (Test Device) → Tests whether the addressed
        device is ready to send or receive a byte of data
        and sets the condition code (CC)
            < : Device is ready to send/receive
            = : Device is not ready

→ Test continues until the device is ready.

→ Once ready, either <u>RD (Read Data)</u> : Transfer data from input device or keyboard into rightmost byte of register A and stored in buffer if required (RD INDEV & STA DATA)

→ WD (Write Data) : a byte of data is loaded into the rightmost byte of register A and then written to the addressed device (LDA DATA & WD OUTDEV)

*→ There are I/O channels that can be used to perform input and output while the CPU is executing other instructions. This allows overlap of computing and I/O, resulting in more efficient system operation.

    ↳ SIO ↔ start I/O

    ↳ TIO ↔ Test I/O

    ↳ HIO ↔ Halt I/O

1) Data movement operations

```
        LDA     #5          ; loads value 5 into register A
        STA     ALPHA       ; store in alpha : A ← (A)+(ALPHA)
                                                5+ (ALPHA)
        LDA     #90         ; load ascii code for 'z' into A
        STCH    C1          : store in character variable C1
        :
ALPHA   RESW    1           ; one word variable
C1      RESB    1           ; one byte variable
```

2) Arithmetic operations ( Beta = alpha + incr - 1)

```
        LDS     INCR        ; load value of incr to AS
        LDA     ALPHA       ; load value of alpha to A
        ADDR    S, A        ; A ← (A)+(S)
                source destination
        SUB     #1          ; A ← (A)-1
        STA     BETA        ; BETA ← (A)
        :
```

4) Looping and indexed operations

GAMMA = ALPHA + BETA    where ALPHA and BETA
                        are arrays of 100 words each.
                        (100x3 = 300 bytes)

```
        LDS     #3          ; S = 3
        LDT     #300        ; T = 300
        LDX     #0          ; X = 0 → which specifies index value
ADDLOOP LDA     ALPHA, X    ; A ← (ALPHA) at the specified address of
                                                        index register
        ADD     BETA, X     ; A ← (A)+(BETA)
        STA     GAMMA, X    ; GAMMA ← (A) at specified index value (0)
        ADDR    S, X        ; X ← (X)+(S) = 0+3 = 3 (index register
                                                    address is 3)
        COMPR   X, T        ; (X) ↔ (T) compared is  3 < 300, CC: <
        JLT     ADDLOOP     ; repeat loop till 300 = 300
```

```
              :

ALPHA     RESW    100

BETA      RESW    100

GAMMA     RESW    100


5)  To read one byte of data from input device F1 and copies it

    to device 05

              INLOOP    TD     INDEV
                        JGR    INLOOP
                        RD     INDEV
                        STCH   DATA
                          :
              OUTLOOP   TD     OUTDEV
                        JEG    OUTLOOP
                        LDCH   DATA
                        WD     OUTDEV
                          :
              INDEV     BYTE   X 'F1'
              OUTDEV    BYTE   X '05'
              DATA      RESB   1


6)  subroutine call to read 100-byte record from an input device

    into memory.

              JSUB    READ
                :

    READ      LDX    #0
              LDT    #100
    RLOOP     TD     INDEV
              JEQ    RLOOP
              RD     INDEV
```

```
        STCH    RECORD,X
        TIXR    T
        JLT     RLOOP
        RSUB
          .
          .
INDEV       BYTE    X 'F1'
RECORD      RESB    100
```

Write a SIC and SIC/XE program to copy 'SYSTEM SOFTWARE' to another string

a) SIC program

```
            LDX     ZERO

LOOP        LDCH    STR1,X

            STCH    STR2,X

            TIX     FIFTEEN

            JLT     LOOP
              .
              .
STR1        BYTE    C 'SYSTEM SOFTWARE'

STR2        RESB    15

ZERO        WORD    0

FIFTEEN     WORD    15
```

b) SIC/XE program

```
            LDX     #0
            LDT     #15

LOOP        LDCH    STR1,X

            STCH    STR2,X

            TIXR    T

            JLT     LOOP
              .
              .
STR1        BYTE    C 'SYSTEM SOFTWARE'

STR2        RESB    15
```

Exercises 1.3

1 Write a sequence of instructions for SIC to set ALPHA equal to the product of BETA and GAMMA. Assume ALPHA, BETA and GAMMA are 1 word ( ALPHA = BETA * GAMMA)

```
        LDA     BETA
        MUL     GAMMA
        STA     ALPHA
         :

ALPHA   RESW    1
BETA    RESW    1
GAMMA   RESW    1
```

2. Write a sequence of instructions for SIC/XE to set ALPHA equal to 4 * BETA - 9. ALPHA, BETA and GAMMA are 1 word. Use immediate addressing for the constants ( A = 4 * B - 9 )

```
        LDA     BETA
        LDS     #4
        MULR    S, A
        SUB     #9
        STA     ALPHA
         :
ALPHA   RESW    1
```

3   Write sic instruction to swap the values of ALPHA and
    BETA.

```
        LDA     ALPHA
        STA     GAMMA
        LDA     BETA
        STA     ALPHA
        LDA     GAMMA
        STA     BETA
        :
ALPHA   RESW 1
BETA    RESW 1
GAMMA   RESW 1
```

4   Write a sequence of instructions for sic to set ALPHA
    equal to the integer portion of BETA ÷ GAMMA. ALPHA,
    BETA, GAMMA are 1 word each

```
        LDA     BETA
        DIV     GAMMA
        STA     ALPHA
        :
ALPHA   RESW 1
BETA    RESW 1
GAMMA   RESW 1
```

5. Write a sequence of instructions for src/xe to divide BETA by GAMMA, setting ALPHA to the integer portion of the quotient and DELTA to the remainder. Use register-to-register instructions to make the calculation as efficient as possible.

| | | | Ex: | B = 5 |
|---|---|---|---|---|
| LDA | BETA | ; A = 5 | | G = 2 |
| DIVF | GAMMA | | | |
| LDS | GAMMA | ; S = 2 | | |
| DIVR | S, A | ; A = A/S = 5/2 = 2 | | |
| STA | ALPHA | ; A = 2 | | |
| MULR | S, A | ; A = A * S = 2 * 2 = 4 | | |
| LDS | BETA | ; S = 5 | | |
| SUBR | A, S | ; S = S - A = 5 - 4 = 1 | | |
| STS | DELTA | ; DELTA = 1 | | |
| : | | | | |
| ALPHA | RESW 1 | | | |
| BETA | RESW 1 | | | |
| GAMMA | RESW 1 | | | |
| DELTA | RESW 1 | | | |

note :

// To find the remainder

Quotient = Dividend / Divisor

Remainder = Dividend - (Quotient * Divisor)

Ex:- Dividend = 10, Divisor = 3, $Q = 10/3 = 3$ ; $R = 10 - (3 * 3) = 10 - 9$
$= 1$

Dividend = 15, Divisor = 3, $Q = 15/3 = 5$ ; $R = 15 - (5 * 3) = 15 - 15 = 0$

6. Write a sequence of instructions for SIC/XE to divide BETA by GAMMA, setting ALPHA to the value of the quotient, rounded to the nearest integer. Use register-to-register instructions to make the calculation as efficient as possible

```
              LDF     BETA
              DIVF    GAMMA
              FIX
              STA     ALPHA
               :
      ALPHA   RESW    1
      BETA    RESW    1
      GAMMA   RESW    1
```

7. Write a sequence of instructions for SIC to clear a 20-byte string to all blanks

```
              LDX     ZERO
      LOOP    LDCH    BLANK
              STCH    STR1 , X
              TIX     TWENTY      ; ADD 1 to index and compare
              JLT     LOOP        ; LOOP y index<100 with 20 & set
               :                                          CC ; < = >
      STR1    RESW    20
      BLANK   BYTE    C ' '
      ZERO    WORD    0
      TWENTY  WORD    20
```

8. Write a sequence of instructions for SIC/XE to clear a 20-byte string to all blanks. Use immediate addressing and register-to-register instructions to make this process as efficient as possible.

```
            LDT    #20
            LDX    #0
CLOOP       LDCH   #0
            STCH   STR1, X
            TIXR   T
            JLT    CLOOP
              :
STR1        RESW   20
```

9. Suppose that ALPHA is an array of 100 words. Write a sequence of SIC instruction to set all 100 elements of the array to 0.

```
            LDA    ZERO                          :
            STA    INDEX           INDEX   RESW   1
LOOP        LDX    INDEX           ALPHA   RECW   100
            LDA    ZERO            ZERD    WORD   0
            STA    ALPHA, X        K300    WORD   100
            LDA    INDEX           THREE   WORD   3
            ADD    THREE
            STA    INDEX
            COMP   K300
            TIX    TWENTY
            TIX    LOOP
```

10  ALPHA is an array of 100 words. write a sequence of instructions for SIC/XE to set all 100 elements of the array to 0. Use immediate addressing and register-to-register instructions to make the process as efficient as possible

```
            LDS     #3
            LDT     #300
            LDX     #0
    LOOP    LDA     #0
            STA     ALPHA, X
            ADDR    S, X
            COMPR   X, T
            JLT     LOOP
            .
            .
    ALPHA   RESW    100
```

11  ALPHA is an array of 100 words. write a sequence of SIC/XE instructions to arrange the 100 words in ascending order and store the result in an array BETA of 100 words.

```
            LDS     #3
            LDT     #300
            LDX     #0
    LOOP    LDA     ALPHA, X
            MUL     #4
```

12. ALPHA and BETA are the two arrays of 100 words. Another array of GAMMA elements are obtained by multiplying the corresponding ALPHA element by 4 and adding the corresponding BETA elements. write the sic|xe instructions for the same.

```
            LDS     #3
            LDT     #300
            LDX     #0
    LOOP    LDA     ALPHA, X
            MUL     #4
            ADD     BETA, X
            STA     GAMMA, X
            ADDR    S, X      ; X ← X+3
            COMPR   X, T      ; X ≥ 300
            JLT     LOOP
            :
    ALPHA   RESW    100
    BETA    RESW    100
    GAMMA   RESW    100
```

13. ALPHA is an array of 100 words. Write a sequence of
    SIC/XE instructions to find the maximum element in the
    array and store results in MAX.

```
              LDS     #3
              LDT     #300
              LDX     #0
     LOOP     LDA     ALPHA,X
              COMP    MAX
              JLT     NOMAX
              STA     MAX
     NOMAX    ADDR    S,X
              COMPR   X,T
              JLT     LOOP
                :
     ALPHA    RESW    100
     MAX      WORD    -32768
```

note: COMP MAX ⟹ indicates Accumulator value is
      compared with MAX and set the CC (condition code)
      ie CC ← < = > of (A) ? (MAX). Based on CC value
      check the condition ie JLT, JGT, JEQ

  → COMPR X,T ⟹ Register value are compared
       X: (X) ? (T) and CC is set ie < = > and
       Jump instruction is called

Explanation

$$ALPHA = \{10, 20, 30, 40, \cdots 32768, \cdots \}$$

$$\begin{array}{cc} \uparrow & \uparrow \\ 0 & 3 \end{array}$$

index

Each value is 1 word = 3 bytes

100 words = 3 x 100 = 300 bytes

index has to be incremented by 3.

ie initially $x = 0, 3, 6, 9 \cdots 300$

According to code : $S = 3, T = 100, X = 0$

1st iteration   LOOP :   Accumulator $(A) = 10$ at $0^{th}$ position

COMP MAX ; 10 ? 32768 sets CC : <

JLT NOMAX

NOMAX :   ADDR S, X   ; $X \leftarrow (X) + S = 0 + 3 \rightarrow$ increment by 3 (next element)

element ic 20

COMPR X, T   ; $X \leftarrow (X) ? ⊗ T$

$\leftarrow 300 ? 300$   CC : 0 <

JLT LOOP    $\Downarrow$

To check whether the array index has come to an end.

iteration

2)   LOOP :   $A \leftarrow 20$ which is at $3^{rd}$ position

LDA ALPHA 3 $\rightarrow$ value at $3^{rd}$ position

COMP MAX ; 20 ? 32768 — CC : <

JLT NOMAX

:

Continue

14. A RECORD contains a 100-byte record. Write a subroutine for sic that will write this record onto device 05.

```
              JSUB      WRREC
              .
              .
WRREC    LDX     ZERO          ; Initialize index register = 0
LOOP     TD      OUTPUT        ; Test output device
         JEQ     LOLOOP        ; Loop if device is busy
         LDCH    RECORD, X     ; load One byte to Accumulator
         WD      OUTPUT        ; write one byte to device
         TIX     LENGTH        ; Add 1 to index & compare to 100
         JLT     LOOP          ; Loop if index is < 100
         RSUB                  ; Exit from subroutine
         .
ZERO     WORD    0
LENGTH   WORD    1
OUTPUT   BYTE    X '05'
RECORD   RESB    100
```

Note: To read and write the data between the device, the device has to be ready to perform. This is done by using TD (Test device) instruction; status of the device is tested and CC is set to either { < (ready) = (not ready) } if ready then RD is executed. ⇒ reads 1 byte of data from device into rightmost byte of register A. If the input device is character-oriented (keyboard), the value placed in reg. A is the ASCII code for the character that was read. Ill'y WD → off device the byte is loaded into the rightmost byte of register A and

15. write a subroutine for SIC/XE to write a RECORD of 100 byte onto output device 05.

```
        JSUB    WRREC    ; Jump to subroutine

WRREC   LDX     #0
        LDT     #100

LOOP    TD      OUTPUT
        JEQ     LOOP
        LDCH    RECORD, X
        WD      OUTPUT
        TIXR    T
        JLT     LOOP
        RSUB
        :
OUTPUT  BYTE    X '05'
RECORD  RESB    100
```

16. write a subroutine for SIC that will read a record into a buffer. The record may be any length from 1 to 100 byte. The end of the record is marked with a "null" character (ASCII code 00). The subroutine should place the length of the record read into a variable named LENGTH.

```
        JSUB    RDREC
        :
RDREC   LDX     ZERO
RLOOP   TD      INDEV
```

```
                JEQ     RLOOP
                RD      INDEV
                COMP    NULL
                JEG     EXIT
                STCH    BUFFER ,X
                TIX     KIOO
                JLT     RLOOP
        EXIT    STX     LENGTH
                RSUB
                  .
                  .
                  .
        ZERO    WORD    'O
        NULL    WORD    O
        KIOO    WORD    1
        INDEV   BYTE    X 'FI'
        LENGTH  RESW    1
        BUFFER  RESB    100
```

17) SIC|XE :

```
                        JSUB    RDREC
                          .
                RDREC   LDX     #O
                        LDT     #100
                        LDS     #O

                RLOOP   TD      INDEV
                        JEB     RLOOP
                        RD      INDEV
                        COMPR.  A, S
                        JEQ     EXIT
```

```
                        STCH    BUFFER ,X
                        TIXR    T
                        JLT     RLOOP
                EXIT    STX     LENGTH
                        RSUB
                          :
                INDEV   BYTE    X 'FI'
                LENGTH  RESW    1
                BUFFER  RESB    100
```

11    To sort an array of 10 words in an ascending order

```
OUTER       LDX    INDEX ;

            LDS    ARR1, X ;

            LDX    #0

INNER       LDT    ARR1, X

            COMR   S,T

            JLT    LOOP

            JEQ    LOOP

            RMO    S,A

            RMO    T,S

            RMO    A,T

            RMO    X,A

            LDX    INDEX

            STS    ARR1,X

            RMO    A,X

            STT    ARR1,X

LOOP        RMO    X,A

            ADD    #3

            COMP   LENGTH

            RMO    A,X

            JLT    INNER

            LDA    INDEX

            ADD    #3

            COMP   LENGTH

            STA    INDEX

            JLT    OUTER
            :
            :
INDEX       WORD   0

ARR1        RESW   10

LENGTH      WORD   30
```

## OP

1 Write a SIC program to copy string 'SYSTEM SOFTWARE' to another string.

```
            LDX     ZERO        ; Initialize X to zero

   MOVECH   LDCH    STR1,X      ; X specifies indexing

            STCH    STR2,X

            TIX     FIFTEEN      ; increment X and compare
                                   with 15

            JLT     MOVECH


   STR1     BYTE    C 'SYSTEM SOFTWARE'

   STR2     RESB    15

   ZERO     WORD    0

   FIFTEEN  WORD    15
```

# Comparison Chart of SIC and SIC/XE machine

| Specification | | SIC | SIC/XE |
|---|---|---|---|
| **Memory** | | - **Word size:** 3 bytes (24 bits)<br>- **Total size:** 32,768 bytes (2^15). Thus any memory address will need at most 15 bits to be referenced ('almost' four hex characters). | - **Word size:** 3 bytes (24 bits)<br>- **Total size:** 32,768 bytes (2^15). Thus any memory address will need at most 15 bits to be referenced ('almost' four hex characters). |
| **Register** | | - **Total Registers:** 5<br>- **Accumulator (A):** Used for most of the operations (number 0)<br>- **Index (X):** Used for indexed addressing (number 1)<br>- **Linkage (L):** Stores return addresses for JSUB (number 2)<br>- **Program Counter (PC):** Address for next instruction (number 8)<br>- **Status Word (SW):** Information and condition codes (number 9). | - **Total Registers:** 9 , same 5 from SIC plus 4 additional ones.<br>- **Base (B):** Used for base-relative addressing (number 3)<br>- **General (S and T):** General use (numbers 4 and 5 resp.)<br>- **Floating Point Accumulator (F):** Used for floating point arithmetic, 48 bits long (number 6) |
| **Instruction Formats** | | - Only one instruction format of 24 bits (3 bytes / 1 word)<br>- Opcode: first 8 bits, direct translation from the Operation Code Table<br>- Flag (x): next bit indicates address mode (0 direct - 1 indexed)<br>- Address: next 15 bits, indicate address of operand according to address mode. | - Four instruction formats<br>- **Format 1 (1 byte):** contains only operation code (straight from table)<br>- **Format 2 (2 bytes):** first eight bits for operation code, next four for register 1 and following four for register 2.<br>- The numbers for the registers go according to the numbers indicated at the registers section (ie, register T is replaced by hex 5).<br>- If the operation uses only one register the last hex digit becomes \0" (ie, TIXR T becomes B850)<br>- **Format 3 (3 bytes):** First 6 bits contain operation code, next 6 bits contain flags, last 12 bits contain displacement for the address of the operand.<br>- Operation code uses only 6 bits, thus the second hex digit will be a affected by the values of the first two flags (n and i)<br>- The flags, in order, are: n, i, x, b, p, and e. Its functionality is explained in the next section.<br>- The last flag e indicates the instruction format (0 for 3 and 1 for 4)<br>- **Format 4 (4 bytes):** same as format 3 with an extra 2 hex digits (8 bits) for addresses that require more than 12 bits to be represented |

| Specification | SIC | SIC/XE |
|---|---|---|
| **Addressing Modes** | • Only two possible addressing modes<br>• Direct (x = 0): operand address goes as it is Indexed (x = 1): value to be added to the register x to obtain real address of the operand. | • five possible addressing modes plus combinations (see page 11 for examples)<br>• Direct (x, b, and p all set to 0): operand address goes as it is. While in general that value is 1, if set to 0 for format 3 we can assume that the rest of the flags (x, b, p, and e) are used as a part of the address of the operand, to make the format compatible to the SIC format<br>• Relative (either b or p equal to 1 and the other one to 0): the address of the operand should be added to the current value stored at the B register (if b = 1) or to the value stored at the PC register (if p = 1)<br>• Immediate (i = 1, n = 0): The operand value is already enclosed on the instruction (ie. lies on the last 12/20 bits of the instruction)<br>• Indirect (i = 0, n = 1): The operand value points to an address that holds the address for the operand value<br>• Indexed (x = 1): value to be added to the value stored at the register x to obtain real address of the operand. This can be combined with any of the previous modes except immediate. |
| **Assembler Considerations** | • Operation code gets translated directly from table (no need to check other bits)<br>• x bit dependent on the addressing mode of the operand. If indexed the code will have to indicate it with \,x" after the operand name (ie. BUFFER,X)<br>• The last 3 hex digits of the address will remain the same, the first hex digit (leftmost) will change if the address is indexed (first bit becomes one, thus the hex digit increases by 8). Ie. if the address of the operand is 124A and the addressing is indexed, the object code will indicate 924A. | • Operation code gets translated directly from table. While the first hex digit remains the same, the second one can change according to the values of the n and i flags. Thus, we can add 1, 2 or 3 to the operation code.<br>• Direct addressing is mainly used in extended format (format 4) and is indicated with a \+" before the operand (an indication that the format is 4, which will also make the e flag to be 1).<br>• Relative: for Base relative, the instruction BASE will precede the current instruction.<br>• Any other format, except immediate, will be considered Program Counter relative. If the displacement with respect to the PC does not fit into the 12 bits, the assembler should try to compute the displacement with respect to the Base register. If neither case works, the instruction should be extended to format 4, where the addressing mode becomes direct. |

| Specification | SIC | SIC/XE |
|---|---|---|
| | | • **Immediate addressing** will be indicated by the use of \#" before the operand name/value (ie. #1)<br><br>• **Indirect addressing** will be indicated by adding the prex \@" to the operand name (ie. @RETADR)<br><br>• **Indexed addressing** will be indicated the same way as it was for the SIC machine, \,X" after the operand name (ie. BUFFER,X)<br><br>• Hex digits for the address are not affected by the content of the flags, since the first two flags affect the second digit of the operation code, and the following four make up its own hex digit. |

# CHAPTER 2
# **Assemblers**

Chapter 2 : <u>Assemblers</u>



Assembler does two functions

1) It converts the mnemonic operation codes into their machine language equivalent

2) Converts symbolic labels into their machine address

The <u>design of assembler can be</u> of

1. Convert mnemonic operation codes to their machine language equivalent. Ex: Translate STL to 14

2. Convert symbolic operands to their equivalent machine addresses. Ex: Translate RETADR to 1033

3. Build the machine instructions in the proper format.

4. Convert the data constents specified in the source program into their internal machine representation. Ex: Translate 'EOF' to 454F46

5. Write the object program and the assembly listing

Different datastructures for assemblers

1 Operation code Table (OPTAB)

2. Symbol Table (SYMTAB)

3 Location Counter Variable (LOCCTR)

I OPERATION CODE TABLE (OPTAB)

a) Contents
  - Mnemonic operation codes
  - machine language equivalent
  - Instruction format and length

b) During pass-1
  - Validates op codes
  - Find instruction length to increase location counter value (LOCCTR)

c) During pass-2
  - Determines the instruction format (3 or 4)
  - Translates the operation codes to their machine language equivalents

d) Implementation
  - static hash table, easy for searching

| Mnemonic name | op code | Format |
|---|---|---|
| ADD m | 18 | 3/4 |
| . | | |
| . | | |
| . | | |

# I. SYMBOL TABLE (SYMTAB)

a) contents

- label name
- label address
- Flags to indicate error conditions
- Data type or length

b) During pass-1

- store label name and assigned address (from LOCCTR) in SYMTAB

c) During pass-2

- symbols used as operands are looked up in SYMTAB

d) Implementation

- A dynamic hash table for efficient insertion and retrieval
- should perform well with non-random keys (LOOP1, LOOP2...)

| Label name | value/addr | Flags | Length |
|------------|------------|-------|--------|
| CLOOP | 0003 | | |
| | | | |
| | | | |

## ii) LOCATION COUNTER VARIABLE (LOCCTR)

- Variable accumulated for address assignment
  ie LOCCTR gives the address of the associated labels

- LOCCTR is initialized to be the beginning address specified in the "START" statement

- After each source stmt is processed during pass-1, the instruction length or data area is added to LOCCTR

→ The functionality of assembler looks like this



note: During pass-1, the address of labels is not known ∵ it is defined later ie called forward reference. To resolve this coc go for pass-2.
    Ei. JEQ RETADR

## 8.1 Assembler Directives

1) **START** — specifies name and starting address for the program

2) **END** — Indicates the end of the source program and optionally specify the first executable instruction in the program

3) **BYTE** — Generate character or hexadecimal constant, occupying as many bytes as needed to represent the constant

4) **WORD** — Generate one-word integer constant

5) **RESB** — Reserve the indicated number of bytes for a data area

6) **RESW** — Reserve the indicated number of words for a data area.

7) **LTORG** — create a literal pool that contains all of the literal operands used since the previous LTORG or the beginning of the program

8) **EQU** — Establishes symbolic name that can be used for improved readability inplace of numeric values and also used to define mnemonic names for registers.

9) ORG — Used to indirectly assign values to symbols

10) USE — Indicates which portion of the source program belong to the various blocks and also indicates a continuation of a previously begun block

11) BASE — Indicates that the base register will contain the address of operand

12) NOBASE — Indicates that the contents of the base register can no longer be relied upon for addressing.

2.1.1 **A simple sic assembler**

The usual (general) format to represent the assembly language program for sic machine with generated assembly code :

| LINE NO | LOCATION | SOURCE STATEMENT | OBJECT CODE |
|---------|----------|------------------|-------------|

| LABEL | OPCODE | OPERAND | COMMENTS |
|-------|--------|---------|----------|

where

→ LABEL : An identifier and optional. labels are used to reduce reliance upon programmers remembering where data or code is located. The length of label differs between assemblers.

Ex:- FIRST STL #4096.

→ OPCODE : Is a machine code instruction. It may require additional information like operands (optional)

Ex:- COMP ZERO ; with operand

OR

RSUB ; without operands

→ OPERAND : Is an additional data or information that the opcode requires. Operands are used to specify constants, labels, immediate data, data contained in another register, an address etc

Advantages and Disadvantages of assembly language

Advantages : → Reduced Errors

→ Faster Translation time

→ changes could be made easier and faster

Disadvantages : → many instructions are required to achieve small tasks

→ Source program tend to be large and difficult to follow

→ Programs are machine dependent, thus the complete program has to be rewritten if the hardware is changed

→ The programmer has to have the complete knowledge of the process architecture and instruction set.

| Mnemonic | Format | Opcode | Effect | Notes |
|---|---|---|---|---|
| ADD m | 3/4 | 18 | A ← (A) + (m..m+2) | X |
| ADDF m | 3/4 | 58 | F ← (F) + (m..m+5) | X F |
| ADDR r1,r2 | 2 | 90 | r2 ← (r2) + (r1) | X |
| AND m | 3/4 | 40 | A ← (A) & (m..m+2) | |
| CLEAR r1 | 2 | B4 | r1 ← 0 | X |
| COMP m | 3/4 | 28 | (A) : (m..m+2) | C |
| COMPF m | 3/4 | 88 | (F) : (m..m+5) | X F C |
| COMPR r1,r2 | 2 | A0 | (r1) : (r2) | X C |
| DIV m | 3/4 | 24 | A ← (A) / (m..m+2) | X |
| DIVF m | 3/4 | 64 | F ← (F) / (m..m+5) | X F |
| DIVR r1,r2 | 2 | 9C | r2 ← (r2) / (r1) | X |
| FIX | 1 | C4 | A ← (F) [convert to integer] | X F |
| FLOAT | 1 | C0 | F ← (A) [convert to floating] | X F |
| HIO | 1 | F4 | Halt I/O channel number (A) | P X |
| J m | 3/4 | 3C | PC ← m | |
| JEQ m | 3/4 | 30 | PC ← m if CC set to = | |
| JGT m | 3/4 | 34 | PC ← m if CC set to > | |
| JLT m | 3/4 | 38 | PC ← m if CC set to < | |
| JSUB m | 3/4 | 48 | L ← (PC); PC ← m | |
| LDA m | 3/4 | 00 | A ← (m..m+2) | |
| LDB m | 3/4 | 68 | B ← (m..m+2) | |
| LDCH m | 3/4 | 50 | A [rightmost byte] ← (m) | X |
| LDF m | 3/4 | 70 | F ← (m..m+5) | X F |
| LDL m | 3/4 | 08 | L ← (m..m+2) | |
| LDS m | 3/4 | 6C | S ← (m..m+2) | X |
| LDT m | 3/4 | 74 | T ← (m..m+2) | X |
| LDX m | 3/4 | 04 | X ← (m..m+2) | |
| LPS m | 3/4 | D0 | Load processor status from information beginning at address m (see Section 6.2.1) | P X |
| MUL m | 3/4 | 20 | A ← (A) * (m..m+2) | |

| Mnemonic | Format | Opcode | Effect | Notes |
|---|---|---|---|---|
| MULF m | 3/4 | 60 | F ← (F) * (m..m+5) | X F |
| MULR r1,r2 | 2 | 98 | r2 ← (r2) * (r1) | X |
| NORM | 1 | C8 | F ← (F) [normalized] | X F |
| OR m | 3/4 | 44 | A ← (A) | (m..m+2) | |
| RD m | 3/4 | D8 | A [rightmost byte] ← data from device specified by (m) | P |
| RMO r1,r2 | 2 | AC | r2 ← (r1) | X |
| RSUB | 3/4 | 4C | PC ← (L) | |
| SHIFTL r1,n | 2 | A4 | r1 ← (r1); left circular shift n bits. [In assembled instruction, r2 = n-1] | X |
| SHIFTR r1,n | 2 | A8 | r1 ← (r1); right shift n bits, with vacated bit positions set equal to leftmost bit of (r1). [In assembled instruction, r2 = n-1] | X |
| SIO | 1 | F0 | Start I/O channel number (A); address of channel program is given by (S) | P X |
| SSK m | 3/4 | EC | Protection key for address m ← (A) (see Section 6.2.4) | P X |
| STA m | 3/4 | 0C | m..m+2 ← (A) | X |
| STB m | 3/4 | 78 | m..m+2 ← (B) | X |
| STCH m | 3/4 | 54 | m ← (A) [rightmost byte] | X |
| STF m | 3/4 | 80 | m..m+5 ← (F) | X F |
| STI m | 3/4 | D4 | Interval timer value ← (m..m+2) (see Section 6.2.1) | P X |
| STL m | 3/4 | 14 | m..m+2 ← (L) | |
| STS m | 3/4 | 7C | m..m+2 ← (S) | X |
| STSW m | 3/4 | E8 | m..m+2 ← (SW) | P |
| STT m | 3/4 | 84 | m..m+2 ← (T) | X |
| STX m | 3/4 | 10 | m..m+2 ← (X) | X |
| SUB m | 3/4 | 1C | A ← (A) - (m..m+2) | |
| SUBF m | 3/4 | 5C | F ← (F) - (m..m+5) | X F |

*Appendix A: SIC/XE Instruction Set and Addressing Modes*

| Mnemonic | Format | Opcode | Effect | Notes |
|---|---|---|---|---|
| SUBR r1,r2 | 2 | 94 | r2 ← (r2) - (r1) | X |
| SVC n | 2 | B0 | Generate SVC interrupt. [In assembled instruction, r1 = n] | X |
| TD m | 3/4 | E0 | Test device specified by (m) | P C |
| TIO | 1 | F8 | Test I/O channel number (A) | P X C |
| TIX m | 3/4 | 2C | X ← (X) + 1; (X): (m..m+2) | C |
| TIXR r1 | 2 | B8 | X ← (X) + 1; (X): (r1) | X C |
| WD m | 3/4 | DC | Device specified by (m) ← (A) [rightmost byte] | P |

range); m indicates a memory address or a constant value larger than 4095. Further information can be found in Section 1.3.2.

The letters in the Notes column have the following meanings:

4  Format 4 instruction

D  Direct-addressing instruction

A  Assembler selects either program-counter relative or base-relative mode

S  Compatible with instruction format for standard SIC machine. Operand value can be between 0 and 32,767 (see Section 1.3.2 for details).

| Addressing type | Flag bits n i x b p e | Assembler language notation | Calculation of target address TA | Operand | Notes |
|---|---|---|---|---|---|
| Simple | 110000 | op c | disp | (TA) | D |
| | 110001 | +op m | addr | (TA) | 4 D |
| | 110010 | op m | (PC) + disp | (TA) | A |
| | 110100 | op m | (B) + disp | ((TA) | A |
| | 111000 | op c,X | disp + (X) | (TA) | D |
| | 111001 | +op m,X | addr + (X) | (TA) | 4 D |
| | 111010 | op m,X | (PC) + disp + (X) | (TA) | A |
| | 111100 | op m,X | (B) + disp + (X) | (TA) | A |
| | 000--- | op m | b/p/e/disp | ((TA) | D S |
| | 001--- | op m,X | b/p/e/disp + (X) | ((TA) | D S |
| Indirect | 100000 | op @c | disp | ((TA)) | D |
| | 100001 | +op @m | addr | ((TA)) | 4 D |
| | 100010 | op @m | (PC) + disp | ((TA)) | A |
| | 100100 | op @m | (B) + disp | ((TA)) | A |
| Immediate | 010000 | op #c | disp | TA | D |
| | 010001 | +op #m | addr | TA | 4 D |
| | 010010 | op #m | (PC) + disp | TA | A |
| | 010100 | op #m | (B) + disp | TA | A |

## Instruction Formats

**Format 1 (1 byte):**

| op |
|---|
| 8 |

**Format 2 (2 bytes):**

| op | r1 | r2 |
|---|---|---|
| 8 | 4 | 4 |

**Format 3 (3 bytes):**

| op | n | i | x | b | p | e | disp |
|---|---|---|---|---|---|---|---|
| 6 | 1 | 1 | 1 | 1 | 1 | 1 | 12 |

**Format 4 (4 bytes):**

| op | n | i | x | b | p | e | address |
|---|---|---|---|---|---|---|---|
| 6 | 1 | 1 | 1 | 1 | 1 | 1 | 20 |

## Addressing Modes

The following addressing modes apply to Format 3 and 4 instructions. Combinations of addressing bits not included in this table are treated as errors by the machine. In the description of assembler language notation, c indicates a constant between 0 and 4095 (or a memory address known to be in this

| Line | | Source statement | | |
|------|--------|--------|-----------|----------------------------------|
| 5    | COPY   | START  | 1000      | COPY FILE FROM INPUT TO OUTPUT |
| 10   | FIRST  | STL    | RETADR    | SAVE RETURN ADDRESS |
| 15   | CLOOP  | JSUB   | RDREC     | READ INPUT RECORD |
| 20   |        | LDA    | LENGTH    | TEST FOR EOF (LENGTH = 0) |
| 25   |        | COMP   | ZERO      | |
| 30   |        | JEQ    | ENDFIL    | EXIT IF EOF FOUND |
| 35   |        | JSUB   | WRREC     | WRITE OUTPUT RECORD |
| 40   |        | J      | CLOOP     | LOOP |
| 45   | ENDFIL | LDA    | EOF       | INSERT END OF FILE MARKER |
| 50   |        | STA    | BUFFER    | |
| 55   |        | LDA    | THREE     | SET LENGTH = 3 |
| 60   |        | STA    | LENGTH    | |
| 65   |        | JSUB   | WRREC     | WRITE EOF |
| 70   |        | LDL    | RETADR    | GET RETURN ADDRESS |
| 75   |        | RSUB   |           | RETURN TO CALLER |
| 80   | EOF    | BYTE   | C'EOF'    | |
| 85   | THREE  | WORD   | 3         | |
| 90   | ZERO   | WORD   | 0         | |
| 95   | RETADR | RESW   | 1         | |
| 100  | LENGTH | RESW   | 1         | LENGTH OF RECORD |
| 105  | BUFFER | RESB   | 4096      | 4096-BYTE BUFFER AREA |
| 110  | .      |        |           | |
| 115  | .      |        | SUBROUTINE TO READ RECORD INTO BUFFER | |
| 120  | .      |        |           | |
| 125  | RDREC  | LDX    | ZERO      | CLEAR LOOP COUNTER |
| 130  |        | LDA    | ZERO      | CLEAR A TO ZERO |
| 135  | RLOOP  | TD     | INPUT     | TEST INPUT DEVICE |
| 140  |        | JEQ    | RLOOP     | LOOP UNTIL READY |
| 145  |        | RD     | INPUT     | READ CHARACTER INTO REGISTER A |
| 150  |        | COMP   | ZERO      | TEST FOR END OF RECORD (X'00') |
| 155  |        | JEQ    | EXIT      | EXIT LOOP IF EOR |
| 160  |        | STCH   | BUFFER,X  | STORE CHARACTER IN BUFFER |
| 165  |        | TIX    | MAXLEN    | LOOP UNLESS MAX LENGTH |
| 170  |        | JLT    | RLOOP     | HAS BEEN REACHED |
| 175  | EXIT   | STX    | LENGTH    | SAVE RECORD LENGTH |
| 180  |        | RSUB   |           | RETURN TO CALLER |
| 185  | INPUT  | BYTE   | X'F1'     | CODE FOR INPUT DEVICE |
| 190  | MAXLEN | WORD   | 4096      | |
| 195  | .      |        |           | |
| 200  | .      |        | SUBROUTINE TO WRITE RECORD FROM BUFFER | |
| 205  | .      |        |           | |
| 210  | WRREC  | LDX    | ZERO      | CLEAR LOOP COUNTER |
| 215  | WLOOP  | TD     | OUTPUT    | TEST OUTPUT DEVICE |
| 220  |        | JEQ    | WLOOP     | LOOP UNTIL READY |
| 225  |        | LDCH   | BUFFER,X  | GET CHARACTER FROM BUFFER |
| 230  |        | WD     | OUTPUT    | WRITE CHARACTER |
| 235  |        | TIX    | LENGTH    | LOOP UNTIL ALL CHARACTERS |
| 240  |        | JLT    | WLOOP     | HAVE BEEN WRITTEN |
| 245  |        | RSUB   |           | RETURN TO CALLER |
| 250  | OUTPUT | BYTE   | X'05'     | CODE FOR OUTPUT DEVICE |
| 255  |        | END    | FIRST     | |

Figure 2.1 Example of a SIC assembler language program.

*ASCII code*
*A : 65*
*Q : 77*

| Line | Loc | Length | LABEL | OPCODE | OPERAND | Object code |
|------|-----|--------|-------|--------|---------|-------------|
| 5 | 1000 | | COPY | START | 1000 | |
| 10 | 1000 | 3 | FIRST | STL | RETADR | 141033 |
| 15 | 1003 | 3 | CLOOP | JSUB | RDREC | 482039 |
| 20 | 1006 | 3 | | LDA | LENGTH | 001036 |
| 25 | 1009 | 3 | | COMP | ZERO | 281030 |
| 30 | 100C | 3 | | JEQ | ENDFIL | 301015 |
| 35 | 100F | 3 | | JSUB | WRREC | 482061 |
| 40 | 1012 | 3 | | J | CLOOP | 3C1003 |
| 45 | 1015 | 3 | ENDFIL | LDA | EOF | 00102A |
| 50 | 1018 | 3 | | STA | BUFFER | 0C1039 |
| 55 | 101B | 3 | | LDA | THREE | 00102D |
| 60 | 101E | 3 | | STA | LENGTH | 0C1036 |
| 65 | 1021 | 3 | | JSUB | WRREC | 482061 |
| 70 | 1024 | 3 | | LDL | RETADR | 081033 |
| 75 | 1027 | 3 | | RSUB | | 4C0000 |
| 80 | 102A | 3 | EOF | BYTE | C'EOF' | 454F46 |
| 85 | 102D | 3 | THREE | WORD | 3    *3 Byte* | 000003 |
| 90 | 1030 | 3 | ZERO | WORD | 0 | 000000 |
| 95 | 1033 | 3 | RETADR | RESW | 1   *→ 3 bytes* | |
| 100 | 1036 | 3 | LENGTH | RESW | 1 | |
| 105 | 1039 | *1000* | BUFFER | RESB | 4096   *( 1000 in hexadecimal )* | |
| 110 | | | | | | |
| 115 | | | | . | | |
| 120 | | | . | | SUBROUTINE TO READ RECORD INTO BUFFER | |
| 125 | 2039 | 3 | RDREC | LDX | ZERO | 041030 |
| 130 | 203C | 3 | | LDA | ZERO | 001030 |
| 135 | 203F | 3 | RLOOP | TD | INPUT | E0205D |
| 140 | 2042 | 3 | | JEQ | RLOOP | 30203F |
| 145 | 2045 | 3 | | RD | INPUT | D8205D |
| 150 | 2048 | 3 | | COMP | ZERO | 281030 |
| 155 | 204B | 3 | | JEQ | EXIT | 302057 |
| 160 | 204E | 3 | | STCH | BUFFER,X | 549039 |
| 165 | 2051 | 3 | | TIX | MAXLEN | 2C205E |
| 170 | 2054 | 3 | | JLT | RLOOP | 38203F |
| 175 | 2057 | 3 | EXIT | STX | LENGTH | 101036 |
| 180 | 205A | 3 | | RSUB | | 4C0000 |
| 185 | 205D | 1 | INPUT | BYTE | X'F1' | F1 |
| 190 | 205E | 3 | MAXLEN | WORD | 4096 | 001000 |
| 195 | | | | . | *← needs 3 bytes* | |
| 200 | | | . | | SUBROUTINE TO WRITE RECORD FROM BUFFER | |
| 205 | | | . | | | |
| 210 | 2061 | 3 | WRREC | LDX | ZERO | 041030 |
| 215 | 2064 | 3 | WLOOP | TD | OUTPUT | E02079 |
| 220 | 2067 | 3 | | JEQ | WLOOP | 302064 |
| 225 | 206A | 3 | | LDCH | BUFFER,X | 509039 |
| 230 | 206D | 3 | | WD | OUTPUT | DC2079 |
| 235 | 2070 | 3 | | TIX | LENGTH | 2C1036 |
| 240 | 2073 | 3 | | JLT | WLOOP | 382064 |
| 245 | 2076 | 3 | | RSUB | | 4C0000 |
| 250 | 2079 | 1 | OUTPUT | BYTE | X'05' | 05  *1 Byte* |
| 255 | 207A | | | END | FIRST | |

**Figure 2.2**   Program from Fig. 2.1 with object code.

The following program contains a main routine that reads records from an input device (code: F1) and copies them to output device (code: 05).

main function calls subroutine RDREC to read a record into a buffer and subroutine WRREC to write record from the buffer to output device.

Each subroutine transfers one record one character at a time because only I/O instructions available are RD and WD.

Since the I/O rates of two devices (disk and a printing terminal) may be different, a buffer is used. The end of each record is marked with a null character is 00 (in hexadecimal). If a record is longer than length of buffer (4096) bytes) then only the first 4096 bytes are copied. The end of file to be copied is indicated by a zero length record.

The program indicates EOF (End of File) on output device when the zero length record (ie end of file) is detected. The program terminates by executing the RSUB instruction since it was called by JSUB instruction.

Procedure to generate object code and object Program
( Intermediate File).

note : we have assumed that the program starts at
address 1000.

1) First and foremost write the LOCCTR addresses

→ START 1000
→ Add 3 bytes for each instruction. (∵ instruction
format for SIC m/c is 24 bits. i.e 3 bytes)

→ BYTE C 'EOF' : count the length of constant
and add those many bytes

→ RESW 2000 : then it should be 2000×3 bytes =
6000B = 1770 (H) added to
previous address

→ RESW 1 : add just 3 bytes

→ RESB 2000 : convert 2000 to hexadecimal (700)
ie 7D0 bytes and add

RESB 4096 : 4096 → 1000$_{(H)}$ is added to
previous value

→ WORD 3 or WORD 0 → 3 bytes added.

2) Start creating the object code.
→ Convert mnemonic operation codes to their machine
language equivalent. Ex: STL to 14
→ Convert symbolic operands to their equivalent machine
address Ex: RETADR to 1033 (forward reference)

→ Build machine instruction in proper format

    a) Direct addressing : $X = 0$ : TA = address

    b) Indexed addressing : $X = 1$ : TA = address + (X)

        → indicated by symbol 'X

        Ex:    STCH    BUFFER, X : → line no. 160

→ Convert the data constants into their machine representation. Ex :- EOF to $h54F46$ (Line no 80)

    $(A = 65, a = 97)$

        $(41)_{16}$      $(61)_{16}$

3) Write the object program ( Intermediate File)

    → object program contains three types of records

    a) Header Record      b) Text Record      c) End Record.

a) <u>Header Record</u> : Contains program name, starting address and length of program.

| column 1 | H |
|---|---|
| col. 2-7 | Program name |
| col. 8-13 | starting address of object Program (Hexadecimal) |
| col. 14-19 | Length of object program in byte (Hexadecimal) |

        → name of program

Ex :-    5   COPY   START   1000

                 → starting address

    :
    :

    255    207A    END

Length of program = last address - starting address

              = 207A - 1000 = 107A

∴ H^COPY^001000^00107A (Header Record)

b) Text Record :

Text record contains the translated instructions (machine code) and data of the program together with an indication of addresses where there are to be loaded.

| col. 1 | T |
|---|---|
| col. 2-7 | starting address for object code in this record (Hexadecimal) |
| col 8-9 | Length of object code in this record in bytes (hexadecimal) |
| col 10-69 | object code represented in hexadecimal (2 columns per byte of object code) |

↓ note·

60 columns
⇒ 10 words ⇒ 30 bytes ⇒ $(1E)_{16}$

length of object code

Ex:-  10    10000   FIRST   STL   RETADDR   141033
                                          3 bytes each } 10 words
      ⋮
      55    101B            LDA   THREE     001020

Text record →

T∧001000∧1E∧141033∧ · · · · · · ·   · · ·   ∧001020

↳ marker for separation

c) End Record :

End record marks the end of the object program and specifies the address in the program where execution is to begin. If no operand is specified then the address of the first executable instruction is used.

| Col 1 | E |
|-------|---|
| Col 2-7 | Address of first Executable instruction in object program (hexadecimal) |

Ex:- 10     1000     FIRST     STL     RETADR     141033
          ⋮
        255               END     FIRST

End record → E∧001000

Let us start for the given program in Fig. 2.11

Given opcodes

| | | | | |
|---|---|---|---|---|
| STL - 14 | J - 3C | LDX - 04 | JLT - 38 | F - 46 |
| JSUB - 48 | STA - 0C | TD - E0 | LDCH - 50 | |
| LDA - 00 | STX - 10 | RD - D8 | WD - DC | |
| COMP - 28 | LDL - 08 | STCH - 54 | E - 45 | |
| JEQ - 30 | RSUB - 4C | TIX - 2C | 0 - 4F | |

①    start incrementing LOCAR

Initially it is <u>1000</u>.

→ start adding 3 bytes each time from line no. 5 to 105

→ 105   1039   BUFFER   RESB   <u>4096</u>
                          ⇃ convert to Hexadecimal ie
                              $(4096)_{16} = 1000$

∴ add 1000 bytes to 1039 = 2039

∴ line m. 125 starts at $2039^{th}$ address continue till line no. 185.

→ 185   205D   INPUT   BYTE   X 'F1'   $\underset{1 \text{ Byte}}{F1}$

∴ add only 1 byte to 205D ⟹ 205E at line no 190.

→ 190   205E   MAXLEN   <u>WORD</u>   4096   001000
                          ⇃ word ⇂ 3 bytes not 1000 byte

∴ 3 bytes added to 205E ⟹ 2061 at line No. 210

→ 210   2061   WRREC   LDX   ZERO   041030.
     continue the same till end.

→ 255   <u>207A</u>   END   FIRST

③ object code for each line.

→ Every line is direct addressing except line  oo 160

and 225

| Line | Loc CTR | label | opcode | operand | Object code |
|------|---------|-------|--------|---------|-------------|
| (i) | 5 | 1000 | COPY | START | 1000 |
| | 10 | 1000 | FIRST | STL | RETADR | 14 1033 |
| | ⋮ | | | | |
| | 95 | 1033 | RETADR | RESW | 1 |

mnemonic code

(ii) 40 1042

(ii) 75 1027 RSUB · · · HC0000

NO operand ·

mnemonic code

(iii) 160 2045 STCH BUFFER, X → indicates indexed addressing

54

address of buffer = 1039

| 8 | 1 | 15 |
|---|---|----|
| opcode | X | address |

| 0101 0100 | 1 | 001 0000 0011 1001 |
|-----------|---|---------------------|

1   0   3   9

⇓

5   4   9   0   3   9

∴ 160 2046 STCH BUFFER, X  549039

(iv) Same for line number 225

225     206A     LDCH     BUFFER, X

                     ↓            ↓

                    50         1039

| opcode | X | | Address | | |
|---|---|---|---|---|---|
| 0101 0000 | 1 | 001 | 0000 | 0011 | 1001 |

      5    0     9     0    3    9

∴   225    206A    LDCH    BUFFER, X    <u>509039</u>

① Object program for Fig 2.2

H＾COPY.  ＾001000＾00107A
T＾001000＾1E＾141033＾H82037＾001036＾281030＾301015＾H82061＾3C1003＾ ··· ＾001029
T＾00101E＾15＾0C1036＾H82061＾081033＾HC0000＾H5HFH6＾000003＾000000
T＾002039＾1E＾0H1030＾001030＾E0205D＾30203F＾D8205D＾281030＾30205h ··· ＾38203F
T＾002057＾1C＾101036＾HC0000＾F1001000＾0H1036＾E02077＾30206H＾ ··· ＾2C1036
T＾002073＾07＾38206H＾HC0000＾05

E＾001000.

② 

SYMBOL TABLE

| Symbol name | Value of the symbol |
|---|---|
| FIRST | 1000 |
| CLOOP | 1003 |
| ENDFIL | 1015 |
| EOF | 102A |
| THREE | 102D |
| ZERO | 1030 |
| RETADR | 1033 |
| LENGTH | 1036 |
| BUFFER | 1039 |
| RDREC | 2039 |
| RLOUP | 203F |

| Symbol name | Value |
|---|---|
| EXIT | 2057 |
| INPUT | 205D |
| MAXLEN | 205E |
| WRREC | 2061 |
| WLOOP | 206H |
| OUTPUT | 2079 |

loader loads into main memory

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0000 | | | | | | | | | | | | | | | | | |
| 0010 | | | | | | | | | | | | | | | | | |
| ⋮ | | | | | | | | | | | | | | | | | |
| 1000 | 14 | 10 | 33 | 48 | 20 | 39 | 00 | 10 | 36 | 28 | 10 | 30 | 30 | 10 | 15 | 48 | ① |
| 1010 | 20 | 61 | 3C | 10 | 08 | 00 | 10 | 2A | 0C | 10 | 39 | 00 | 10 | 2D | 0C | 10 | ② |
| 1020 | 36 | 48 | 20 | 61 | 08 | 10 | 33 | 4C | 00 | 00 | 45 | 4F | 46 | 00 | 00 | 08 | |
| 1030 | 00 | 00 | 00 | RETADR | | | LENGTH | | | | | | | | | | |
| 1040 | | | | | | | | | | | | | | | | | |
| 1050 | | | | | | | | | | | | | | | | | |
| ⋮ | | | | | | BUFFER | | | | | | | | | | | |
| 2030 | | | | | | | | | | 04 | 10 | 80 | 00 | 10 | 30 | E0 | ③ |
| 2040 | 20 | 5D | 30 | 20 | 3F | D8 | 20 | 5D | 28 | 10 | 30 | 30 | 20 | 51 | 54 | 90 | |
| 2050 | 39 | 2C | 20 | 5E | 38 | 20 | 3F | 10 | 10 | 36 | 4C | 00 | 00 | F1 | 00 | 10 | ④ |
| 2060 | 00 | 04 | 10 | 30 | E0 | 20 | 79 | 30 | 20 | 6h | 50 | 90 | 39 | DC | 20 | 79 | |
| 2070 | 2C | 10 | 36 | 38 | 20 | 6H | 4C | 00 | 00 | 05 | ⑤ | | | | | | |
| 2080 | | | | | | | | | | | | | | | | | |

# Functions of Pass-I and Pass-II

Pass 1 :

→ Assign addresses to all statements in the program

→ Save the values (addresses) assigned to all labels for use in Pass 2

→ Perform some processing of assembler directives (includes processing that affects address assignment, such as determining the length of data areas defined by BYTE, RESW etc)

Pass 2 :

→ Assemble instructions (Translating operation codes and looking up address)

→ Generate data values defined by BYTE, WORD etc

→ Perform processing of assembler directives not done during pass-1

→ write the object program and the assembly listing.

Pass 1:

```
begin
   read first input line
   if OPCODE = 'START' then
      begin
         save #[OPERAND] as starting address
         initialize LOCCTR to starting address
         write line to intermediate file
         read next input line
      end {if START}
   else
      initialize LOCCTR to 0
   while OPCODE ≠ 'END' do
      begin
         if this is not a comment line then
            begin
               if there is a symbol in the LABEL field then
                  begin
                     search SYMTAB for LABEL
                     if found then
                        set error flag (duplicate symbol)
                     else
                        insert (LABEL,LOCCTR) into SYMTAB
                  end {if symbol}
               search OPTAB for OPCODE
               if found then
                  add 3 {instruction length} to LOCCTR
               else if OPCODE = 'WORD' then
                  add 3 to LOCCTR
               else if OPCODE = 'RESW' then
                  add 3 * #[OPERAND] to LOCCTR
               else if OPCODE = 'RESB' then
                  add #[OPERAND] to LOCCTR
               else if OPCODE = 'BYTE' then
                  begin
                     find length of constant in bytes
                     add length to LOCCTR
                  end {if BYTE}
               else
                  set error flag (invalid operation code)
            end {if not a comment}
         write line to intermediate file
         read next input line
      end {while not END}
   write last line to intermediate file
   save (LOCCTR - starting address) as program length
end {Pass 1}
```

**Figure 2.4(a)**   Algorithm for Pass 1 of assembler.

Pass 2:

```
begin
   read first input line {from intermediate file}
   if OPCODE = 'START' then
       begin
           write listing line
           read next input line
       end {if START}
   write Header record to object program
   initialize first Text record
   while OPCODE ≠ 'END' do
       begin
           if this is not a comment line then
               begin
                   search OPTAB for OPCODE
                   if found then
                       begin
                           if there is a symbol in OPERAND field then
                               begin
                                   search SYMTAB for OPERAND
                                   if found then
                                       store symbol value as operand address
                                   else
                                       begin
                                           store 0 as operand address
                                           set error flag (undefined symbol)
                                       end
                               end {if symbol}
                           else
                               store 0 as operand address
                           assemble the object code instruction
                       end {if opcode found}
                   else if OPCODE = 'BYTE' or 'WORD' then
                       convert constant to object code
                   if object code will not fit into the current Text record then
                       begin
                           write Text record to object program
                           initialize new Text record
                       end
                   add object code to Text record
               end {if not comment}
           write listing line
           read next input line
       end {while not END}
   write last Text record to object program
   write End record to object program
   write last listing line
end {Pass 2}
```

Figure 2.4(b)   Algorithm for Pass 2 of assembler.

2.2. Machine Dependent Assembler Features

. Here we consider an example of SIC/XE machine

→ As we know already, SIC/XE has

a) Registers : A  X  L  B  S  T  F  PC  SW
   (0  1  2  3  4  5  6   8   9)

b) Data format : Integer : 3 bytes
                 Character : 1 byte
                 Float : 6 bytes

c) Instruction Formats :

Format 1 : 1 byte

$$\boxed{\overset{8}{\text{opcode}}}$$   Ex: FLOAT, FIX

Format 2 : 2 bytes

$$\boxed{\overset{8}{\text{opcode}} | \overset{4}{r_1} | \overset{4}{r_2}}$$   Ex: ADDR A,X

Format 3 : 3 bytes

$$\boxed{\overset{6}{op} | n | i | x | b | p | e | \overset{12}{disp}}$$   Ex:- STL RETADR

Format 4 : 4 bytes

$$\boxed{\overset{6}{op} | n | i | x | b | p | e | \overset{20}{address}}$$   Ex: +JSUB RDREC

→ we have 20 address lines ∴ we can have $2^{20}$ addresses



| nibble | 19 | 18 | 17 | 16 | 15 | 14 | 13 | 12 | 11 | | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 00000 | | | | | | | | | | | | | |
| 00001 | | | | | | | | | | | | | |
| 00002 | | | | | | | | | | | | | |
| ⋮ | | | | | | | | | | | | | |
| FFFFF | | | | | | | | | | | | | |

d) Addressing mode are determined based on 6 bits

n i x b p e

(i) →

| n | i | x | Addressing mode |
|---|---|---|---|
| 1 | 0 | | Indirect addressing |
| 0 | 1 | | Immediate |
| 1 | 1 | | not immediate, not indirect |
| 0 | 0 | | · Simple addressing |
| | | 1 | Indexed addressing |
| | | 0 | Direct addressing |

(ii)

| b | p | e | Addressing mode |
|---|---|---|---|
| 0 | 1 | | Program Counter Relative |
| 1 | 0 | | Base relative |
| 1 | 1 | | Invalid (can't be set) |
| 0 | 0 | | No PC relative, no base relative |
| | | 1 | Format 4 instruction |
| | | 0 | Format 3 instruction |

Different addressing mode notations

1) Indirect Addressing : @

2) Immediate Addressing : #

3) Extended Format : +

4) Indexed Addressing : operand, X

5) character string : C ' '

6) Base - Register : BASE

7) Current value of PC : *

→ The addressing priority are as follows

a) PC relative addressing : $-2048 \leq disp \leq 2047$
$$(FFFF F800 \leq disp \leq 7FF)$$

b) Base relative addressing : $0 \leq disp \leq 4095$
$$(0 \leq disp \leq FFF)$$

c) Extended Instruction Format :

note : <u>negative</u> <u>numbers</u> <u>are</u> represented <u>in</u> 2's compliment

<u>Procedure</u> to create object code for SIC/XE program

1) Write the LOCCTR addresses for each instruction in the program.

→ if operand field is
(i) memory address → Format 3 ⇒ Add 3 bytes
(ii) Register - Register → Format 2 ⇒ Add 2 bytes
(iii) + before operand → Format 4 ⇒ Add 4 bytes

→ if it is RESW 2000
. $2000 \times 3$ bytes $= (6000)_d = (1770)_H$ ⇒ Add these
many bytes to previous address.
. multiplication by 3 ∵ each word is 3 bytes

→ RESW 1 ⇒ Add just 3 bytes

→ RESB 2000 ⇒ Add 2000 bytes in hexadecimal
ie $(2000)_d = (7D0)_H$

$\longrightarrow$ RESB 4096

$\cdot (4096)_d = (1000)_H \Rightarrow$ Add 1000 bytes

$\longrightarrow$ BYTE C 'EOF' $\Rightarrow$ Count the length of constant
and add those many bytes

$\longrightarrow$ Enter the labels onto SYMTAB (pass 1)

2) Once we are done with LOCCTR calculation and then finding program length = End Address — start address

3) now start creating the object code (Pass 2) based on different addressing mode and set corresponding bits and calculate displacement

$\rightarrow$ For extended format, displacement = address

$\rightarrow$ For Reg-to-Reg instruction, write the opcode address followed by register numbers.

Ex: CLEAR X $\Rightarrow$ B410 (Format 2)
(1) $\rightarrow$ number of X register in the list

$\rightarrow$ For PC relative, disp = TA - PC

$\rightarrow$ For Base relative, disp = TA - (B)

| Line | | Source statement | | |
|------|--------|---------|----------|-------------------------------|
| 5 | COPY | START | 0 | COPY FILE FROM INPUT TO OUTPUT |
| 10 | FIRST | STL | RETADR | SAVE RETURN ADDRESS |
| 12 | | LDB | #LENGTH | ESTABLISH BASE REGISTER |
| 13 | | BASE | LENGTH | |
| 15 | CLOOP | +JSUB | RDREC | READ INPUT RECORD |
| 20 | | LDA | LENGTH | TEST FOR EOF (LENGTH = 0) |
| 25 | | COMP | #0 | |
| 30 | | JEQ | ENDFIL | EXIT IF EOF FOUND |
| 35 | | +JSUB | WRREC | WRITE OUTPUT RECORD |
| 40 | | J | CLOOP | LOOP |
| 45 | ENDFIL | LDA | EOF | INSERT END OF FILE MARKER |
| 50 | | STA | BUFFER | |
| 55 | | LDA | #3 | SET LENGTH = 3 |
| 60 | | STA | LENGTH | |
| 65 | | +JSUB | WRREC | WRITE EOF |
| 70 | | J | @RETADR | RETURN TO CALLER |
| 80 | EOF | BYTE | C'EOF' | |
| 95 | RETADR | RESW | 1 | |
| 100 | LENGTH | RESW | 1 | LENGTH OF RECORD |
| 105 | BUFFER | RESB | 4096 | 4096-BYTE BUFFER AREA |
| 110 | . | | | |
| 115 | . | | SUBROUTINE TO READ RECORD INTO BUFFER | |
| 120 | . | | | |
| 125 | RDREC | CLEAR | X | CLEAR LOOP COUNTER |
| 130 | | CLEAR | A | CLEAR A TO ZERO |
| 132 | | CLEAR | S | CLEAR S TO ZERO |
| 133 | | +LDT | #4096 | |
| 135 | RLOOP | TD | INPUT | TEST INPUT DEVICE |
| 140 | | JEQ | RLOOP | LOOP UNTIL READY |
| 145 | | RD | INPUT | READ CHARACTER INTO REGISTER A |
| 150 | | COMPR | A,S | TEST FOR END OF RECORD (X'00') |
| 155 | | JEQ | EXIT | EXIT LOOP IF EOR |
| 160 | | STCH | BUFFER,X | STORE CHARACTER IN BUFFER |
| 165 | | TIXR | T | LOOP UNLESS MAX LENGTH |
| 170 | | JLT | RLOOP | HAS BEEN REACHED |
| 175 | EXIT | STX | LENGTH | SAVE RECORD LENGTH |
| 180 | | RSUB | | RETURN TO CALLER |
| 185 | INPUT | BYTE | X'F1' | CODE FOR INPUT DEVICE |
| 195 | . | | | |
| 200 | . | | SUBROUTINE TO WRITE RECORD FROM BUFFER | |
| 205 | . | | | |
| 210 | WRREC | CLEAR | X | CLEAR LOOP COUNTER |
| 212 | | LDT | LENGTH | |
| 215 | WLOOP | TD | OUTPUT | TEST OUTPUT DEVICE |
| 220 | | JEQ | WLOOP | LOOP UNTIL READY |
| 225 | | LDCH | BUFFER,X | GET CHARACTER FROM BUFFER |
| 230 | | WD | OUTPUT | WRITE CHARACTER |
| 235 | | TIXR | T | LOOP UNTIL ALL CHARACTERS |
| 240 | | JLT | WLOOP | HAVE BEEN WRITTEN |
| 245 | | RSUB | | RETURN TO CALLER |
| 250 | OUTPUT | BYTE | X'05' | CODE FOR OUTPUT DEVICE |
| 255 | | END | FIRST | |

**Figure 2.5** Example of a SIC/XE program.

| Line | Loc | | Source statement | | | Object code | |
|------|-----|---|------------------|---|---|-------------|---|
| 5 | 0000 | | COPY | START | 0 | | |
| 10 | 0000 | FIRST | STL | RETADR | | 17202D | |
| 12 | 0003 | | LDB | #LENGTH | | 69202D | |
| 13 | | | BASE | LENGTH | | | |
| 15 | 0006 | CLOOP | +JSUB | RDREC | | 4B101036 | |
| 20 | 000A | | LDA | LENGTH | | 032026 | |
| 25 | 000D | | COMP | #0 | | 290000 | |
| 30 | 0010 | | JEQ | ENDFIL | | 332007 | |
| 35 | 0013 | | +JSUB | WRREC | | 4B10105D | |
| 40 | 0017 | | J | CLOOP | | 3F2FEC | |
| 45 | 001A | ENDFIL | LDA | EOF | | 032010 | |
| 50 | 001D | | STA | BUFFER | | 0F2016 | |
| 55 | 0020 | | LDA | #3 | | 010003 | |
| 60 | 0023 | | STA | LENGTH | | 0F200D | |
| 65 | 0026 | | +JSUB | WRREC | | 4B10105D | |
| 70 | 002A | | J | @RETADR | | 3E2003 | |
| 80 | 002D | EOF | BYTE | C'EOF' | | 454F46 | |
| 95 | 0030 | RETADR | RESW | 1 | | | |
| 100 | 0033 | LENGTH | RESW | 1 | | | |
| 105 | 0036 | BUFFER | RESB | 4096 | | | |
| 110 | | | . | | | | |
| 115 | | | . | SUBROUTINE TO READ RECORD INTO BUFFER | | | |
| 120 | | | . | | | | |
| 125 | 1036 | RDREC | CLEAR | X | | B410 | |
| 130 | 1038 | | CLEAR | A | | B400 | |
| 132 | 103A | | CLEAR | S | | B440 | |
| 133 | 103C | | +LDT | #4096 | | 75101000 | |
| 135 | 1040 | RLOOP | TD | INPUT | | E32019 | |
| 140 | 1043 | | JEQ | RLOOP | | 332FFA | |
| 145 | 1046 | | RD | INPUT | | DB2013 | |
| 150 | 1049 | | COMPR | A,S | | A004 | |
| 155 | 104B | | JEQ | EXIT | | 332008 | |
| 160 | 104E | | STCH | BUFFER,X | | 57C003 | |
| 165 | 1051 | | TIXR | T | | B850 | |
| 170 | 1053 | | JLT | RLOOP | | 3B2FEA | |
| 175 | 1056 | EXIT | STX | LENGTH | | 134000 | |
| 180 | 1059 | | RSUB | | | 4F0000 | |
| 185 | 105C | INPUT | BYTE | X'F1' | | F1 | |
| 195 | | | . | | | | |
| 200 | | | . | SUBROUTINE TO WRITE RECORD FROM BUFFER | | | |
| 205 | | | . | | | | |
| 210 | 105D | WRREC | CLEAR | X | | B410 | |
| 212 | 105F | | LDT | LENGTH | | 774000 | |
| 215 | 1062 | WLOOP | TD | OUTPUT | | E32011 | |
| 220 | 1065 | | JEQ | WLOOP | | 332FFA | |
| 225 | 1068 | | LDCH | BUFFER,X | | 53C003 | |
| 230 | 106B | | WD | OUTPUT | | DF2008 | |
| 235 | 106E | | TIXR | T | | B850 | |
| 240 | 1070 | | JLT | WLOOP | | 3B2FEF | |
| 245 | 1073 | | RSUB | | | 4F0000 | |
| 250 | 1076 | OUTPUT | BYTE | X'05' | | 05 | |
| 255 | 1077 | | END | FIRST | | | |

**Figure 2.6**  Program from Fig. 2.5 with object code.

Consider the example of figure 3.5

1) Add the length of each instruction and add it to LOCCTR and find the program length

Program length = End address - start address

$= 1077 - 0000 = 1077$

2) Create the symbol table

| Symbol Name | PC value |
|---|---|
| FIRST | 0000 |
| CLOOP | 0006 |
| ENDFIL | 001A |
| EOF | 002D |
| RETADR | 0030 |
| LENGTH | 0033 |
| BUFFER | 0036 |
| RDREC | 103C |
| RLOOP | 1040 |
| EXIT | 1056 |
| INPUT | 105C |
| WRREC | 105D |
| WLOOP | 1062 |
| OUTPUT | 1076 |

3) start creating object code (pass-2)

10    0000    FIRST    STL    RETADR    (Format 3 ∴ oprd is
                        ͜          ͜                memory address)
                        14        0030

By default assembler uses PC relative addressing

| opcode | n | i | x | b | p | e | disp |
|--------|---|---|---|---|---|---|------|
| 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | |
| | | | | | | | | 02D |

→    TA = PC + disp
↳    Displacement = TA - PC
                = RETADR - LOCCTR (location of next inst$^n$ to be
                                                        executed)
                = 0030 - 0003 = $\cancel{0}$02D
                                    ∵ format 3 displacement
↳ 02D is within range of $-2048 \leq disp \leq 2047$      is 12 bits
↳ opcode is 6 bits ( h + 2 ⟹ 1 nibble + 2-bits)
        ⟹ last 2 bits can be represented by 4 bits
        but always last 2 bits are "zero"

        Ex:- h ⟹ 0$\cancel{1}$0$\cancel{0}$
                    only 2 bits

        C ⟹ 11$\cancel{0}$$\cancel{0}$
                C

↳ not immediate, not indirect so set n=1, i=1
↳ not indexed    x=0, not base relative b=0 but it
        is pc relative p=1, not format 4 e=0
↳ write the whole instruction's object code ie

| 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 02D |
|---|---|---|---|---|---|---|---|-----|

nibble
representation    1      7        2    02D

        STL    RETADR    17202D

13    0003    $\underset{68}{LDB}$    $\underset{0033}{\#LENGTH}$

→ opcode for LDB = 68

→ it is imma calculate disp.

$$TA = PC + disp$$

$$disp = TA - PC = 0033 - 0006 = \emptyset 02D$$

→ PC relative ∴ operand is memory address

→ it is immediate so P = 1



|opcode| n | i | x | b | p | e | |
|---|---|---|---|---|---|---|---|
| 6 | 1 | 0 | 0 | 1 | 0 | 0 | 02D |

6     9     2 02D

LDB    #LENGTH ⟹ 692.02D

15    0006    $\underset{}{CLOOP}$    $\underset{48}{+JSUB}$    $\underset{0030}{RDREC}$    → format 4

→ disp = (operand) ∴ it is extended format (F4)

$= 001036$ (20 bits)

→ not immediate, not indirect n=1, i=1.

→ extended e=1

|opcode| n | i | x | b | p | e | Address |
|---|---|---|---|---|---|---|---|
| 4 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 001036 |

4     B         1   01036

CLOOP    + JSUB    RDREC ⟹ 4B101036

20    CODA      LEN      LENGTH   ⟶ Format 3
                     00       0033

→  Pc relative,   disp = TA - PC
         P=1                    = 0033 - 000D = 026

→  not immediate, not indirect, not indexed co

         n=1, i=1, x=0



         n  i  x  b  p  e
      0 | 0  0  1  1  0  0  1  0 |     026

    0  0  3         2      0 2 6

    LDA   LENGTH  ⟹  032026


25    000D      COMP      #0
                     28

→  immediate not Pc relative because operand is
   direct value but not memory address.
         ∴ displacement = operand = 000

→  Immediate addressing   n=0, i=1, b=0, p=0

         n  i  x  b  p  e
      2 | 1  0  0  1  0  0  0  0 |  000

    2      9        0 000

    COMP  #0  ⟹  290000

3C     CRLO     $\underset{3C}{JEG}$     $\underset{001A}{ENDFIL}$     $\Rightarrow$ Format 3

• Pc relative $\therefore$ disp = TA - PC

$$= 001A - 0013 = 007$$

within range

• not immediate, not indirect    n=1, i=1

$$\begin{array}{|c|c|c|c|c|c|c|c|c|c|} \hline & n & i & x & b & p & e & & \\ 3 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & 0 & 007 \\ \hline \end{array}$$

3 3        2 007

$\therefore$ JEQ    ENDFIL $\Rightarrow$ 332007

35 / 65    0013    $\underset{48}{+JSUB}$    $\underset{105D}{WRREC}$    $\Rightarrow$ Format 4

• Displacement = address of operand

$$= 0105D$$

~~within~~

$$\begin{array}{|c|c|c|c|c|c|c|c|c|c|} \hline & n & i & x & b & p & e & & \\ 4 & 1 & 0 & 1 & 1 & 0 & 0 & 0 & 1 & 0105D \\ \hline \end{array}$$

4 B        10105D

+ JSUB   WRREC $\Rightarrow$ 4B10105D

40    0017    $\underset{3C}{J}$   $\underset{0006}{CLOOP}$    $\Rightarrow$ Format 3

disp = TA - PC = 0006 - 001A

$$= -1A \; (\text{it takes 2's complement})$$

$$= FEC$$

- It is Pc relative $P=1$
- not immediate, not indirect $n=0, i=1$



object code : 3     F     2-FEC

H S    CO1A    BADFR:    LDA    EOF  $\Rightarrow$ Format 3 (Pc relative)
                          $\underbrace{\phantom{LDA}}_{00}$  $\underbrace{\phantom{EOF}}_{001D}$

$$disp = TA - PC = 002D - 001D = \cancel{0}010$$



0     3     2     010

LDA    EOF $\Rightarrow$ 032010

SU    CO1D              STA    BUFFER $\Rightarrow$ Format 3 (Pc relative)
                        $\underbrace{\phantom{STA}}_{0C}$  $\underbrace{\phantom{BUFFER}}_{0036}$

$$disp = TA - PC = 0036 - 0020 = \cancel{0}016$$



0     F     2.016

STA    BUFFER $\Rightarrow$ OF2016

                LDA    #3                      $\Rightarrow$ immediate address

$$disp = 003$$



LDA    #3 $\Rightarrow$ 010 003

65    0026    STA    @ORG 9TH    ⟹ Format 3 (pc relative)
                ‾OC‾        ‾0033‾

dsp = TA - PC = 0033 - 0026 = ∅00D

```
     n i x b p e
 0 |1|1|1|0|0|1|0|     00D
    0   F      200D   ⟹ OF200D
```

70    002D    ST    @PCINDR    ⟹ Format 3 - Indirect
              ‾3C‾      ‾0030‾

dsp = TA - PC = 0030 - 002D = ∅003

```
     n i x b p e
 3 |1|1|1|0|0|1|0|     003
    3   E      2   003   ⟹ 3E2003
```

80    002D    EOF    BYTE    C 'EOF'

⟹ Convert EOF to hexadecimal ascii value

    E — 45
    0 — 4F
    F — 46

125    1036    RDREC    CLEAR    X    A X L B S T F PC SW
                        ‾B4‾               0 1 2 3 4 5 6 8 9

                    ⟹ B410
                        ↳ this is not accumulator
        ⟹ only 2 bytes since it is register-to-register mode

130    1038    CLEAR    A    ⟹ B400
              ‾B4‾

182    103A    CLEAR    S    ⟹ B440
              ‾B4‾

130     1033     COMPR   A, S    $\Rightarrow$ A004
$$\underbrace{A0}$$

145     1057     TIXR   T    $\Rightarrow$ B850
$$\underbrace{B8}$$

93     1066     TIXR   T   $\Rightarrow$ B850
$$\underbrace{B8}$$

210     106A     LDREC   CLEAR   X   $\Rightarrow$ B410
$$\underbrace{B4}$$

135     103C     +LDT   #4096   $\Rightarrow$ Format 4 & Immediate addressing

$disp = (4096)_{16} = 01000$



| | n | i | x | b | p | e | |
|---|---|---|---|---|---|---|---|
| 7 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 01000 |

7   5    1 01000 $\Rightarrow$ 75101000

138     1040     RLOOP   TD   INPUT $\Rightarrow$ Format 3 + PC relative
$$\underbrace{E0} \qquad \underbrace{105C}$$

$disp = TA - PC = 105C - 1043 = 0019$



| | n | i | x | b | p | e | |
|---|---|---|---|---|---|---|---|
| E | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 019 |

E   3    2 019 $\Rightarrow$ E32019

line         1032        PRC        BUFEP        $\Rightarrow$ Format 3 + Pc relative
                          ‿            ‿
                          30          1040

$$dsp = TA - PC = 1040 - 1046 = -6$$



$$\Downarrow$$
2's complement
$$\Downarrow$$
FFA

3  | 0 0 1 1 0 0 1 0 | FFA

n i x b p e

3      3      2  FFA $\Rightarrow$ 332FFA

line        1043        JLSI        EXIT        $\Rightarrow$ Format 3 + Pc relative
                          ‿            ‿
                          30          1056

$$dsp = TA - PC = 1056 - 104E = 008$$



3  | 0 0 1 1 0 0 1 0 | 008

n i x b p e

3      3      2  008      $\Rightarrow$  332008

line        1049        RD        INPUT        $\Rightarrow$ Format 3 + Pc relative
                          ‿            ‿
                          D8          105C

$$dsp = TA - PC = 105C - 1049 = 013$$



D  | 1 0 1 1 0 0 1 0 | 013

n i x b p e

D      B      2  013  $\Rightarrow$  DB2013

160        1040        STCH        BUFFER, X  $\Rightarrow$  Indexed + Pc relative
*.*                    $\underbrace{\quad}_{SH}$   $\underbrace{\quad}_{0036}$

$dsp = TA - PC = 0036 - 1051 = -101B$

$= \underbrace{6FE5}_{(H123)_{10}}  > 2047$

$\therefore$ it is not pc relative, go for base relative

$dsp = \overset{TA}{BUFFER} - \overset{B}{B}$ (length is stored in base register at 0033)

$= 0036 - 0033 = \cancel{\emptyset}003$



| S | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 003 |

    5        7           C   003  $\Rightarrow$  57C003


170        1056        JLT         KLOUP  $\Rightarrow$  Format 3 + Pc relative
                       $\underbrace{\quad}_{38}$   $\underbrace{\quad}_{1040}$

$dsp = TA - PC = 1040 - 1056 = FEA$



          n i x b p e
| 3 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | FEA |

    3        B         2   FEA  $\Rightarrow$  3B2FEA

175    1056    STIT    STX    LENGTH    $\Rightarrow$ Format 3 + PC relative

$\underbrace{\qquad}_{10}$    $\underbrace{\qquad}_{0033}$

$disp = TA - PC = 0033 - 1059 = 6FDA > 2047$

$\therefore$ go for base relative mode

$disp = TA - (B) = 0033 - 0033 = \emptyset000$



$$\underbrace{\qquad}_{1} \quad \underbrace{\qquad}_{3} \quad \underbrace{\qquad}_{4000} \Rightarrow 134000$$

180    1059    RSUB    $\Rightarrow$ Format 3
$\underbrace{\qquad}_{4C}$

no displacement $\therefore$ no operand.



$$\underbrace{\qquad}_{4} \quad \underbrace{\qquad}_{F} \quad \underbrace{\qquad}_{0000} \Rightarrow 4F0000$$

185    105C    INPUT    BYTE    X'F1'

$\Rightarrow$ character string $\therefore$ store as it is

212    105F    LDT    LENGTH    $\Rightarrow$ Format 3
$\underbrace{\qquad}_{74}$    $\underbrace{\qquad}_{0033}$

$disp = TA - PC = 0033 - 1062 = \underset{\underbrace{4443 \, > \, 2047}}{EFD1}$

$\therefore$ go for base relative

disp = TA - (B) = 0033 - 0033 = 0000

| | n | i | x | b | p | e | |
|---|---|---|---|---|---|---|---|---|
| 7 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 000 |

7   7    4000 ⟹ 774000

---

013   1062   +LDBUF  #0  #0774 ⟹ Format 3 + PC relative
                     E0    1076

disp = TA - PC = 1076 - 1065 = 011

| | n | i | x | b | p | e | |
|---|---|---|---|---|---|---|---|
| C | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 011 |

E   3    2 011 ⟹ E32011

---

030   1065   JTFB  +LDBUF ⟹ Format 3 + PC relative
              30    1062

disp = TA - PC = 1062 - 1068 = FFA

| | n | i | x | b | p | e | |
|---|---|---|---|---|---|---|---|
| 3 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | FFA |

3   3    2 FFA ⟹ 332FFA

---

033   1068   LDCH  BUFFER, X ⟹ indexed
              50    0036

disp = TA - PC = 0036 - 106B = -1035 > 2047

∴ go for base relative mode

disp = TA - B = 0036 - 0033 = 0003

n i x b p e

5 | 0 0 1 1 1 0 0 | 003

5    3        C 003  ⟹ 53 C003

---

33C . 1069    100   OUTPUT   ⟹ F3 + Pc relative
                └PC┘  └1076┘

disp = TA - PC = 1076 - 106E = 0008

n i x b p e

D | 1 1 1 1 0 0 1 0 | 008

D    F        2 008  ⟹ DF2008

---

3?? 1070   JLT   LOODOP   ⟹ Format 3 + PC relative
                 └38┘  └1062┘

disp = TA - PC = 1062 - 1073 = FEF

n i x b p e

3 | 1 0 1 1 0 0 1 0 | FEF

3    B        2 FEF  ⟹ 3B2FEF

---

3?? 1073   RSUB    ⟹ Format 3

n i x b p e

4 | 1 1 1 0 0 0 1 0 | 000

4    F        0000  ⟹ 4F0000

---

3?. 1076   OUTPUT   BYTE   X 'OS'

⟹ character string  ∴ store as it is ⟹ OS

## 4) Object program

H␸COPY ␸000000␸001077

T␸000000␸1D␸172020␸692020␸4B101036␸032026␸290000␸332007␸4B10105D␸
                                                        3F2FEC␸032010

T␸00001D␸13␸0F2016␸010003␸0F202D␸4B10105D␸3E2003␸454F46

T␸001036␸1D␸B410␸B400␸B440␸75101000␸E32019␸332FFA␸DB2013␸A004␸
                          332FFA␸DB2013␸A004␸332008␸57C003␸B850

T␸001053␸1D␸3B2FEA␸134000␸4F0000␸F1␸B410␸774000␸E32011␸332FFA␸
                                        5BC003␸DF2008␸B850

T␸001070␸07␸3B2F6F␸4F0000␸05

E␸000000

## Loading into memory

| Address | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0000 | 17 | 20 | 2D | 69 | 20 | 20 | 4B | 10 | 10 | 36 | 03 | 20 | 26 | 29 | 00 | 00 | ① |
| 0010 | 33 | 20 | 07 | 4B | 10 | 10 | 5D | 3F | 2F | EC | 03 | 20 | 10 | 0F | 20 | 16 | ② |
| 0020 | 01 | 00 | 03 | 0F | 20 | DD | 4B | 10 | 10 | 5D | 3E | 20 | 08 | 45 | 4F | 46 | |
| 0030 | RETADR | | | LENGTH | | | | | | | | | | | | | |
| ⋮ | | | | | | | B | U | F | F | E | R | | | | | |
| 1030 | | | | | | | Bh | 10 | B4 | DD | Bh | h0 | 75 | 10 | 10 | 00 | |
| 1040 | 33 | 20 | 19 | 33 | 2F | FA | DB | 20 | 13 | A0 | 04 | 33 | 20 | 08 | 57 | C0 | ③ |
| 1050 | 03 | B8 | 50 | 3B | 2F | 6A | 13 | h0 | 00 | 4F | 00 | 00 | F1 | Bh | 10 | 77 | ④ |
| 1060 | h0 | 00 | 63 | 20 | 11 | 33 | 2F | FA | 53 | C0 | 03 | DF | 20 | 08 | B8 | 50 | |
| 1070 | 3B | 2F | EF | 4F | 00 | 00 | 05 | ⑤ | | | | | | | | | |

1. Generate the complete object program for the following assembly level program

   CLEAR - B4, LDA- 00, LDB - 68, ADD- 18, TIX- 2C,

   JLT - 38    STA - 0C

| PASS-I | LENGTH | LABEL | OPCODE | OPERAND | PASS-II |
|--------|--------|-------|--------|---------|---------|
| 0000 | | SUM | START | 0 | |
| 0000 | 2 | | CLEAR | X | B410 |
| 0002 | 3 | | LDA | #0 | 010000 |
| 0005 | H | | +LDB | #TOTAL | 69101789 |
| | | | BASE | TOTAL | |
| 0009 | 3 | LOOP | ADD | TABLE,X | 1BA00P |
| 000C | 3 | | TIX | COUNT | 2F2007 |
| 000F | 3 | | JLT | LOOP | 3F2FF7 |
| 0012 | 4 | | +STA | TOTAL | 0F101789 |
| 0016 | 3 | COUNT | RESW | 1 | |
| 0019 | 1770 | TABLE | RESW | 2000 (1770)H | |
| 1789 | 3 | TOTAL | RESW | 1 | |
| 178C | | | END | FIRST | |

RESW   2000 ⟹ 2000×3 = (6000)bytes = $(1770)_H$

∴  0019 + 1770 = 1789

Program Length = 178C - 0000 = 178C

1) 0000    CLEAR    X    (Register-to-Register)

  ⟹ directly opcode with register numbers

  ⟹ B410

2) 0002    LDA    #0    ⟹ Immediate Addressing

disp = 000

| | n | i | x | b | p | e | |
|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 000 |

  0    01    0    000    ⟹ 01 0000

3) 0005    + LDB    #TOTAL    ⟹ Extended & Immediate — Format 4 with PC relative    + immediate

disp = operand address
    = 1789

| | n | i | x | b | p | e | |
|---|---|---|---|---|---|---|---|
| 6 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 1789 |

  6    9    1    01789    ⟹ 69101789

4) 0009    LOOP    ADD    TABLE,X    ⟹ indexed with PC relative

TA = PC + disp
disp = TA - PC
    = 0019 - 000C = 00D

| | n | i | x | b | p | e | |
|---|---|---|---|---|---|---|---|
| 1 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 00D |

  1    B    A    00D    ⟹ 1BA00D

5) 000C    TIX    COUNT    ⟹ Format-3

disp = 0016 - 000F = 007

| | n | i | x | b | p | e | |
|---|---|---|---|---|---|---|---|
| 2 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 007 |

  2    F    2    007    ⟹ 2F2007

6) 000F     JLT    LOOP $\Longrightarrow$ Format -3   PC relative

$$disp = TA - PC$$
$$= 0009 - 0012 = FF7$$

$\underset{range}{\underline{within}} \longrightarrow -2048 \leq FF7 \leq 2048$

| 3 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | FF7 |
|---|---|---|---|---|---|---|---|---|-----|

n f a b p e

3     B     2    FF7 $\Longrightarrow$ 3F2FF7

7) 0012     + STA    TOTAL $\Longrightarrow$ Format 4

$$disp = operand \ value = 1789$$

| 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 0178 9 |
|---|---|---|---|---|---|---|---|---|--------|

n i x b p e

0    F     1 01789 $\Longrightarrow$ 0F101789

object program (pass-2)

H∧ Sum ∧ 0000 , ∧ 00178C

T∧ 000000 ∧ 16 ∧ BH10 ∧ 010000 ∧ 69101789 ∧ 1BA00D ∧ 2F2007 ∧ 3F2FF7 ∧ 0F101789

E∧ 000000

SYMTAB $\longrightarrow$
(pass-1)

| SYMBOL NAME | VALUE |
|-------------|-------|
| LOOP | 0009 |
| COUNT | 0016 |
| TABLE | 0019 |
| TOTAL | 1789 |

2. Generate the complete object program for the following assembly level program. Also indicate the contents of symbol table at the end. Assume standard SIC model and assume the following mle codes in HEX

$$LDA = 00 \qquad STA = 0C \qquad TIX = 2C \qquad JLT = 38$$
$$LDX = 04 \qquad ADD = 18 \qquad RSUB = 4C$$

| LOCCTR (PASS-1) | LENGTH | LABEL | OPCODE | OPERAND | OBJECT CODE |
|---|---|---|---|---|---|
| | | SUM | START | 4000 | |
| 4000 | 3 | FIRST | LDX | ZERO | 045788 |
| 4003 | 3 | | LDA | ZERO | 005788 |
| 4006 | 3 | LOOP | ADD | TABLE,X | 18C015 |
| 4009 | 3 | | TIX | COUNT | 2C5785 |
| 400C | 3 | | JLT | LOOP | 384006 |
| 400F | 3 | | STA | TOTAL | 0C578B |
| 4012 | 3 | | RSUB | | 4C0000 |
| 4015 | 1770 | TABLE | RESW | 2000 (1770)H | |
| 5785 | 3 | COUNT | RESW | 1 | |
| 5788 | 3 | ZERO | WORD | 0 | 000000 |
| 578B | 3 | TOTAL | RESW | 1 | |
| 578E | | | END | FIRST | |

Program length = END address - starting address
= 578E - 4000 = 178E

→ since it is SIC program, we have two
addressing mode { direct addressing (x=0)
indexed addressing (x=1)

• so directly put opcode with operand address

| opcode | x | address |
|--------|---|---------|
| 8 | 1 | 15 |

i) H000   FIRST   LDX   ZERO

| OH | 0 1 0 1 | 788 |
|----|---------|-----|

$\Rightarrow$ 0h5788

2) H006   ADD   TABLE, X   $\Rightarrow$ indexed addressing
              18      H015

| 18 | 1 φ 0 0 | 015 |
|----|---------|-----|

   18        c         015   $\Rightarrow$ 18C015

**SYMTAB**

| NAME | VALUE |
|-------|-------|
| FIRST | H000 |
| LOOP | H006 |
| TABLE | H015 |
| COUNT | 5785 |
| ZERO | 5788 |
| TOTAL | 578B |

object program

H∧ SUM   ∧⁰⁰H000 ∧00178E

T∧ 00H000 ∧ 15 ∧ 0H5788 ∧ 005788 ∧ 18C015 ∧ 2C5785 ∧ 38 H006 ∧ 0C5788 ∧ H0000

T∧ 005788 ∧ 03 ∧ 000000

E∧ 00H000

LOADING INTO MAIN MEMORY

|      | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | A  | B  | C  | D  | E  | F  |
|------|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 0000 |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 0010 |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |
| :    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |
| H000 | 0A | 57 | 88 | 00 | 57 | 88 | 18 | C0 | 15 | 2C | 57 | 85 | 38 | 40 | 06 | 0C |
| H010 | 57 | 8B | HC | 00 | 00 |    |    |    |    |    |    |    |    |    |    |    |
| H020 |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |
| .    |    |    |    |    |    | TABLE |  |    |    |    |    |    |    |    |    |    |
| .    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |
| .    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 5780 |    |    |    |    |    | COUNT |  | 00 | 00 | 00 | TOTAL |  |    |    |    |    |
| 5790 |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |

3. Generate the object code for each statement in the following sic/xe program and generate the object program for the same.

| LOCCTR | LENGTH | LABEL | OPCODE | OPERAND | OBJECT-CODE |
|--------|--------|-------|--------|---------|-------------|
|        |        | SUM   | START  | 0       |             |
| 0000   | 3      | FIRST | LDX    | #0      | 050000      |
| 0003   | 3      |       | LDA    | #0      | 010000      |
| 0006   | 4      |       | +LDB   | #TABLE2 | 69101790    |
|        |        |       | BASE   | TABLE2  |             |
| 000A   | 3      | LOOP  | ADD    | TABLE, X | 1BA013     |
| 000D   | 3      |       | ADD    | TABLE2, X | 1BC000    |
| 0010   | 3      |       | TIX    | COUNT   | 2F200A      |
| 0013   | 3      |       | JLT    | LOOP    | 3B2FFA      |
| 0016   | 4      |       | +STA   | TOTAL   | 0F102F00    |
| 001A   | 3      |       | RSUB   |         | 4F0000      |
| 001D   | 3      | COUNT | RESW   | 1       |             |
| 0020   | 1770   | TABLE | RESW   | 2000 (1770)₁ |        |
| 1790   | 1770   | TABLE2 | RESW  | 2000 (1770)₁₁ |       |
| 2F00   | 3      | TOTAL | RESW   | 1       |             |
| 2F03   |        |       | END    | FIRST   |             |

LDX = 04    LDB= 68    TIX= 2C    STA= 0C
LDA =00    ADD = 18    JLT = 38    RSUB=4C

→ keep assigning the length for each instruction based on

    (i) 1st operand is memory address — 3 byte

    (ii) 1st operand is register — 2 byte

    (iii) + (Extended format) — 4 bytes

→ Find the LOCCTR value ; Program length = 2F03 - 0000

                                                  = 2F03

→ Create SYMTAB

| Name | Value |
|------|-------|
| FIRST | 0000 |
| LOOP | 000A |
| COUNT | 001D |
| TABLE | 0020 |
| TABLE2 | 1790 |
| TOTAL | 2F00 |

→ object code for each instruction

1) 0000   FIRST   $\underset{04}{LDX}$  #0  → immediate addressing

disp = 000



   0    S      0    000 ⟹ 050000

2) 0003   $\underset{00}{LDA}$  #0  → immediate addressing

disp = 000



   0    1      0000 ⟹ 010000

3) 0006      +LDB $\underset{\cancel{68}}{}$    #TABLE2 $\longrightarrow$ Extended + immediate

TA= PC+disp

disp = TA = pc = $\cancel{1790}$ = 080A

disp = $\underset{Target}{\overset{}{}}$ address = $\underline{01790}$ $\underset{20bit}{}$



| | n | i | x | b | p | e | |
|---|---|---|---|---|---|---|---|
| 6 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 01790 |

6      9       101790 $\implies$ 69101790

4) 000A    LOOP    ADD$\underset{18}{}$    TABLE,X $\longrightarrow$ Indexed + Pc relative

TA = PC + disp

disp = TA - PC = 0020 - 000D = $\cancel{0}$013

| | n | i | x | b | p | e | |
|---|---|---|---|---|---|---|---|
| 1 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 013 |

1    B     A 013 $\implies$ 1BA013

5) 000D    ADD$\underset{18}{}$   TABLE2,X $\longrightarrow$ indexed + base relative

( ∵ TABLE2 is stored in base register)

Initially we can try for Pc - relative & checkout whether displacement is within the range.

disp= TA - PC

= 1790 - 0010 = $(1780)_H$ ⩾ $(6016)_d > (2047)_d$

∴ go for base relative

disp= TA - B (look for address of TABLE2 in SYMTAB)

= 1790 - 1790 = $\cancel{0}$000

| | n | i | x | b | p | e | |
|---|---|---|---|---|---|---|---|
| 1 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 000 |

1    B     C   000 $\implies$ 1BC000

6) 0010     $\underset{2C}{TIX}$    COUNT   → Pc relative

     disp : TA - PC = 001D - 0013 = ∅00A

| 2 | n i x b p e<br>1 1 1 1 0 0 1 0 | 00A |
|---|---|---|

   2      F     2 00A  ⇒  2F200A


7) 0013     $\underset{3\emptyset}{JLT}$    LOOP ⇒ Pc relative

     disp = TA - PC = 000A - 0016 = FF4h

| 3 | n i x b p e<br>1 0 1 1 0 0 1 0 | FF4 |
|---|---|---|

   3      B     2 FF4  ⇒  3B2FF4h


8) 0016    $+\underset{0C}{STA}$    TOTAL ⇒ Extended (Format 4)

     disp = address q TOTAL = 2F00

| 0 | n i x b p e<br>1 1 1 1 0 0 0 1 | 02F00 |
|---|---|---|

   0      F     1 02F00


9) RSUB 001A    RSUB<br>
     ⇒ no operand ∴ no displacement

| 4 | 1 1 1 1 0 0 0 0 | 000 |
|---|---|---|

   4      F     0 000  ⇒  4F0000

→ Object program

H ∧ sum ∧ 000000 ∧ 002F03

T ∧ 000000 ∧ 1D ∧ 050000 ∧ 010000 ∧ 6910179C ∧ 1BA013 ∧ 1BC000 ∧ 2F260A ∧ 3B2FFh
∧ 0F102F00 ∧ hF0000

E ∧ 000000

→ loader loads into memory

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0000 | 05 | 00 | 00 | 01 | 00 | 00 | 69 | 10 | 17 | 90 | 1B | A0 | 13 | 1B | C0 | 00 |
| 0010 | 2F | 20 | 0A | 3B | 2F | Fh | 0F | 10 | 2F | 00 | hF | 00 | 00 | COUNT | | |
| 0040 | | | | | | TABLE | | | | | | | | | | |
| 1790 | | | | | | TABLE 2 | | | | | | | | | | |
| 2F00 | TOTAL | | | | | | | | | | | | | | | |

A. Generate the machine code for the following sic/xe program

Given   JSUB = A0, LDA = 80, LDX = 60, STA = 5O, COMP = 9O,

RSUB = HC   JEQ = BD ,  J = B8

| LOCCTR | LENGTH | LABEL | OPCODE | OPERAND | OBJECT CODE |
|---|---|---|---|---|---|
|  |  | COPY | START | 1000 |  |
| 1000 | 4 | CLOOP | +JSUB | RDREC |  |
|  | 3 |  | LDA | LENGTH |  |
|  | 3 |  | COMP | ZERO |  |
|  | 3 |  | JEQ | EXIT |  |
|  | 3 |  | J | CLOOP |  |
|  | 3 | EXIT | STA | BUFFER |  |
|  | 3 |  | LDA | THREE |  |
|  | 3 |  | STA | TOTAL LENGTH) |  |
|  | 3 |  | RSUB |  |  |
|  |  | BUFFER | RESCO | 10U |  |
|  | 3 | EOF | BYTE | C EOF' |  |
|  | 3 | ZERU | WORD | 0 |  |
|  | 9 | THREE | WORD | 3 |  |
|  | 3 | LENGTH | RESW | 1 |  |
|  | 3 | TOTAL LENGTH | RESW | 1 |  |
|  | 3 | RDREC | LDX | ZERO |  |

## 2.2.2. Program Relocation

Absolute assembly program is one which executes properly, only if program is loaded from specified location.

Ex: All SIC programs are absolute assembly program

Consider the SIC program

| 5 | 1000 | COPY | START | 1000 |
|---|------|------|-------|------|
| 10 | 1000 | FIRST | STL | RETADDR | 141033 |
| 15 | 1003 | CLOOP | JSUB | RDREC | 482039 |
| | : | | | | |
| 55 | 101B | | LDA | THREE | 001020 |
| | : | | | | |
| 85 | 1020 | THREE | WORD | 3 | 000003 |

→ Here program is loaded at address 1000.

→ Line no. 55 specifies that the register A is to be loaded from memory address 1020 (object code).

→ Suppose we attempt to load and execute the program at address 2000 instead of address 1000, the address 1020 will not contain the value that we expected, as it might be part of some other user's program.

→ Obviously we need to make some change in the address portion of this instruction so we can load and execute the program at address 2000.

→ At the same time, there are statements like line no. 95 which generates a constant 3, that should remain the same regardless of where the program is loaded.

→ From the object code, we can't it is not possible to tell which values represent addresses and which represent constant data items.

→ This is all because the assembler doesnot know the actual location where the program will be loaded till load time. ∴ it cannot make the necessary changes required.

→ only parts of the program that require modification at load time are those that specify direct addresses.

This is achieved through relocatable program for SIC/XE program relocation

machine.)

2.2.2. Program Relocation

Program relocation is a process of modifying the addresses used in address sensitive instructions of a program such that program can execute correctly from allocated memory area. It is often needed to have more than one program at same time, sharing the memory and other resources of the machine. Because of this, it is necessary to load a program into memory whenever it is available. Hence relocation of the addresses in the program is required and this will be done during loading time. Assembler only indicates those instructions which need modification and this information is passed to loader.

The Assembler solves the relocation problem as follows:

→ keeping track of operand address relative to start of a program

→ Generating commands for loader which add the beginning address to operand relative address

The An object program that contains the information necessary to perform this kind of modification is called a "relocatable program". we can accomplish this with a modification record as follows

Modification Record

col. 1          m

col 2-7          starting location of the address field to be
                 modified, relative to the beginning of the program

col 8-9          length of the address field to be modified,
                 in half-bytes (hexadecimal).

→  The length is stored in half-bytes (rather than bytes)
   because the address field to be modified may not
   occupy an integral number of bytes.
              Ex: 20 bits = 5 half-bytes

→  The starting location is the location of the byte
   containing the leftmost bits of the address field to be
   modified. If this field occupies an odd number of
   half-bytes, it is assumed to begin in the middle of
   the first byte at the starting location.

Ex: SIC/XE program

| 5 | 0000 | COPY | START | | |
|---|---|---|---|---|---|
| 10 | 0000 | FIRST | STL | RETADR | 17202D |
| 12 | 0003 | | LDB | #LENGTH | 69202D |
| 13 | | | BASE | LENGTH | |
| 15 | 0006 | CLOOP | +JSUB | RDREC | HB101036 |
| | | | ⋮ | | |
| | | | +JSUB | WRREC | 4B1010SD |
| 35 | 0013 | | J | CLOOP | 3E2FEC |
| 40 | 0017 | | ⋮ | | |
| 65 | 0026 | | +JSUB | WRREC | 4B1010SD |
| | | | ⋮ | | |
| 100 | 0036 | BUFFER | RESB | 4096 | |
| | | | ⋮ | | B410 |
| | | RDREC | CLEAR | X | |
| 125 | 1036 | | | | |

program is loaded at address 0000



Fig: Example for program relocation

→ JSUB instruction at line 15 is loaded at address 0006

→ The address field contains 01036 (address of RDREC)

15 0006 +JSUB RDREC 4B101036

→ Suppose we want to load this program beginning at address 5000, as shown in fig (b), the address of instruction labeled RDREC will be 6036.

→ Likewise if we load at 7420 as in fig (c), then address of RDREC will be 4B108456.

→ It means, irrespective of the starting address loaded, RDREC is always 1036 bytes past the starting address of the program. This is the reason we initialized the location counter to 0. (ie relative to the starting address)

→ The modification record looks like

0006    CLOOP    +JSUB    RDREC    4B101036



note: 05 because it is 20 bits address ⇒ 05 half byte ⇓ 05 X 4 = 20 bits

m̲ 000007 ̲ 05

Actually at location 0007, first half byte is part of flag bits x, b, p, e. But length 05 tells loader to modify only last 5 half bytes. Hence instruction 4B1 remains unchanged.

Relocation for instruction of line 35 and 65

| 35 | 0013 | +JSUB | CORREC | 4B10105D |
| 65 | 0026 | +JSUB | CORREC | 4B10105D |

$m_A 0000144.05$ & $m_A 000022.05$

→ if we add 5000 then address should be

4B1/0#1036            4B1/01036
    5000          +    7420 — relocatable address
———————          ———————————
4B1/06036            4B108456
    7420
A

→ Some instructions like
    CLEAR S  }  doesnot need modification ∵ operand is
    LDA  #3  }  not a memory address.

→ 10  STL  RETADR : doesnot need modification ∵
   operand is specified using program-counter relative or
   base-relative addressing. Here the displacement is
   always 02D. Irrespective of location of program loaded,
   it is always 2D bytes away from the STL instruction.

→ The ten distance between LENGTH and BUFFER
   will always be 3 bytes.

The object program is rewritten as (Fig 2.6)

H^COPY   ^000000^001077
T^000000^1D^172027D^692027D^HB101036^032026^290000^·-^4B101050D^032010
T^00001D^13^0F2016^010003^0F2007D^HB101050D^3E2003^454F46
T^001036^1D^B410^B400^B440^75101000^·· ^B850
T^001053^1D^3B2FEA^132000^4F0000^·· ^B850
T^001070^07^3B2FEF^4F0000^05
M^000007^05
M^00001A^05
M^000027^05
E^000000

## 2.3 Machine Independent Assembler features

→ machine independent means some assembler features that are not closely related to machine architecture.

This section includes

2.3.1 → The implementation of literals within an assembler

2.3.2 → Two assembler directives EQU and ORG used to define the symbols

2.3.3 → Use of expression in assembler language statements

2.3.4 → Implementation of program blocks

2.3.5 → Implementation of control sections

## 2.3.1 Literals

→ Constant operand can be specified as a part of the instruction that uses it, instead of using a label which is defined as constant else where. Such an operand is called a literal because the value is stated 'literally' in the instruction

Ex:-

45    001A    ENDFIL    LDA    $\underbrace{EOF}_{Label}$    032010

⋮

8c    0020    EOF    BYTE    C'EOF'    454F46

⇓ can be written as

45    001A    ENDFIL    LDA    =C'EOF'    032010

⋮

LTORG
=C'EOF'                                    454F46

The object code generated for line 45, 215 and 230 in fig 2.6 and fig 2.10 are identical.

(i) 45    001A    ENDFIL    LDA = C 'EOF'    032010

| opcode | n | i | x | b | p | e | disp |
|---|---|---|---|---|---|---|---|
| 0000 00 | 1 | 1 | 0 | 0 | 1 | 0 | 0000 0001 0000 |

$$disp = opaddr - pc$$
$$= 0020 - 001D = 01D$$

TA = (pc) + disp

$$\Rightarrow 032010$$

(ii)  215    1062    WLOOP    TD = X '05'    E3 2011

TD = ED

$$disp = opaddress - pc$$
$$= 1076 - 1065 = 011$$

| 11 0 | 00 | 1 | 1 | 0 | 0 | 1 | 0 | 0000 | 0001 | 0001 |
|---|---|---|---|---|---|---|---|---|---|---|

E     3         2     0        1     1

(iii)  230    106B    WD    = X '05'    DF2008

$$1076 \quad * \quad = X '05'$$

WD = DC

$$disp = opaddress - pc = 1076 - 106E = 008$$

| | | D | i | x | b | p | e | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 1101 | 11 | 1 | 1 | 0 | 0 | 1 | 0 | 0000 | 0000 | 1000 |

D   C              2      0    0    8
    D   F

| Line | | Source statement | | |
|---|---|---|---|---|
| 5 | COPY | START | 0 | COPY FILE FROM INPUT TO OUTPUT |
| 10 | FIRST | STL | RETADR | SAVE RETURN ADDRESS |
| 13 | | LDB | #LENGTH | ESTABLISH BASE REGISTER |
| 14 | | BASE | LENGTH | |
| 15 | CLOOP | -JSUB | RDREC | READ INPUT RECORD |
| 20 | | LDA | LENGTH | TEST FOR EOF (LENGTH = 0) |
| 25 | | COMP | #0 | |
| 30 | | JEQ | ENDFIL | EXIT IF EOF FOUND |
| 35 | | -JSUB | WRREC | WRITE OUTPUT RECORD |
| 40 | | J | CLOOP | LOOP |
| 45 | ENDFIL | LDA | =C'EOF' | INSERT END OF FILE MARKER |
| 50 | | STA | BUFFER | |
| 55 | | LDA | #3 | SET LENGTH = 3 |
| 60 | | STA | LENGTH | |
| 65 | | +JSUB | WRREC | WRITE EOF |
| 70 | | J | @RETADR | RETURN TO CALLER |
| 93 | | LTORG | | |
| 95 | RETADR | RESW | 1 | |
| 100 | LENGTH | RESW | 1 | LENGTH OF RECORD |
| 105 | BUFFER | RESB | 4096 | 4096-BYTE BUFFER AREA |
| 106 | BUFEND | EQU | * | |
| 107 | MAXLEN | EQU | BUFEND-BUFFER | MAXIMUM RECORD LENGTH |
| 110 | . | | | |
| 115 | . | | SUBROUTINE TO READ RECORD INTO BUFFER | |
| 120 | . | | | |
| 125 | RDREC | CLEAR | X | CLEAR LOOP COUNTER |
| 130 | | CLEAR | A | CLEAR A TO ZERO |
| 132 | | CLEAR | S | CLEAR S TO ZERO |
| 133 | | +LDT | #MAXLEN | |
| 135 | RLOOP | TD | INPUT | TEST INPUT DEVICE |
| 140 | | JEQ | RLOOP | LOOP UNTIL READY |
| 145 | | RD | INPUT | READ CHARACTER INTO REGISTER A |
| 150 | | COMPR | A,S | TEST FOR END OF RECORD (X'00') |
| 155 | | JEQ | EXIT | EXIT LOOP IF EOR |
| 160 | | STCH | BUFFER,X | STORE CHARACTER IN BUFFER |
| 165 | | TIXR | T | LOOP UNLESS MAX LENGTH |
| 170 | | JLT | RLOOP | HAS BEEN REACHED |
| 175 | EXIT | STX | LENGTH | SAVE RECORD LENGTH |
| 180 | | RSUB | | RETURN TO CALLER |
| 185 | INPUT | BYTE | X'F1' | CODE FOR INPUT DEVICE |
| 195 | . | | | |
| 200 | . | | SUBROUTINE TO WRITE RECORD FROM BUFFER | |
| 205 | . | | | |
| 210 | WRREC | CLEAR | X | CLEAR LOOP COUNTER |
| 212 | | LDT | LENGTH | |
| 215 | WLOOP | TD | =X'05' | TEST OUTPUT DEVICE |
| 220 | | JEQ | WLOOP | LOOP UNTIL READY |
| 225 | | LDCH | BUFFER,X | GET CHARACTER FROM BUFFER |
| 230 | | WD | =X'05' | WRITE CHARACTER |
| 235 | | TIXR | T | LOOP UNTIL ALL CHARACTERS |
| 240 | | JLT | WLOOP | HAVE BEEN WRITTEN |
| 245 | | RSUB | | RETURN TO CALLER |
| 255 | | END | FIRST | |

**Figure 2.9** Program demonstrating additional assembler features.

| Line | Loc | Source statement | | | Object code |
|---|---|---|---|---|---|
| 5 | 0000 | COPY | START | 0 | |
| 10 | 0000 | FIRST | STL | RETADR | 17202D |
| 13 | 0003 | | LDB | #LENGTH | 69202D |
| 14 | | | BASE | LENGTH | |
| 15 | 0006 | CLOOP | +JSUB | RDREC | 4B101036 |
| 20 | 000A | | LDA | LENGTH | 032026 |
| 25 | 000D | | COMP | #0 | 290000 |
| 30 | 0010 | | JEQ | ENDFIL | 332007 |
| 35 | 0013 | | +JSUB | WRREC | 4B10105D |
| 40 | 0017 | | J | CLOOP | 3F2FEC |
| 45 | 001A | ENDFIL | LDA | =C'EOF' | 032010 |
| 50 | 001D | | STA | BUFFER | 0F2016 |
| 55 | 0020 | | LDA | #3 | 010003 |
| 60 | 0023 | | STA | LENGTH | 0F200D |
| 65 | 0026 | | +JSUB | WRREC | 4B10105D |
| 70 | 002A | | J | @RETADR | 3E2003 |
| 93 | | | LTORG | | |
| | 002D | * | =C'EOF' | | 454F46 |
| 95 | 0030 | RETADR | RESW | 1 | |
| 100 | 0033 | LENGTH | RESW | 1 | |
| 105 | 0036 | BUFFER | RESB | 4096 | |
| 106 | 1036 | BUFEND | EQU | * | |
| 107 | 1000 | MAXLEN | EQU | BUFEND-BUFFER | |
| 110 | | . | | | |
| 115 | | . | SUBROUTINE TO READ RECORD INTO BUFFER | | |
| 120 | | . | | | |
| 125 | 1036 | RDREC | CLEAR | X | B410 |
| 130 | 1039 | | CLEAR | A | B400 |
| 132 | 103A | | CLEAR | S | B440 |
| 133 | 103C | | +LDT | #MAXLEN | 75101000 |
| 135 | 1040 | RLOOP | TD | INPUT | E32019 |
| 140 | 1043 | | JEQ | RLOOP | 332FFA |
| 145 | 1046 | | RD | INPUT | DB2013 |
| 150 | 1049 | | COMPR | A,S | A004 |
| 155 | 104B | | JEQ | EXIT | 332008 |
| 160 | 104E | | STCH | BUFFER,X | 57C003 |
| 165 | 1051 | | TIXR | T | B850 |
| 170 | 1053 | | JLT | RLOOP | 3B2FEA |
| 175 | 1056 | EXIT | STX | LENGTH | 134000 |
| 180 | 1059 | | RSUB | | 4F0000 |
| 185 | 105C | INPUT | BYTE | X'F1' | F1 |
| 195 | | . | | | |
| 200 | | . | SUBROUTINE TO WRITE RECORD FROM BUFFER | | |
| 205 | | . | | | |
| 210 | 105D | WRREC | CLEAR | X | B410 |
| 212 | 105F | | LDT | LENGTH | 774000 |
| 215 | 1062 | WLOOP | TD | =X'05' | E32011 |
| 220 | 1065 | | JEQ | WLOOP | 332FFA |
| 225 | 1068 | | LDCH | BUFFER,X | 53C003 |
| 230 | 106B | | WD | =X'05' | DF2008 |
| 235 | 106E | | TIXR | T | B850 |
| 240 | 1070 | | JLT | WLOOP | 3B2FEF |
| 245 | 1073 | | RSUB | | 4F0000 |
| 255 | | | END | FIRST | |
| | 1076 | * | =X'05' | | 05 |

**Figure 2.10**   Program from Fig. 2.9 with object code.

Literal Pools:-

→ All the literal operands used in a program are gathered together into one or more literal pools.

→ Normally literals are placed into a pool at the end of the program, which shows the assigned address and the generated data values.

→ The drawback of keeping literal pool at the end of the program is use the literal operand is too far away from the instruction referring it and requires a large amount of storage reservation for the buffer too.

→ To avoid this we use an assembler directive LTORG (ORIGIN OF LITERALS) which instructs the assembler to assemble the current literals pool immediately

→ when the assembler encounters a LTORG statement, it creates a literal pool that contains all of the literal operands used since the previous LTORG (or the beginning of the program). ⇒ ie keep the literal operand close to the instruction.

→ Same literal may be used more than once in the program ie duplicate literals, but it stores only one copy of the specified data value.

Ex:- 215    1062    WLOOP    TD = X'05'
     230    106B              WD = X'05

→ Apart from one copy of data value, it stores only one data area with this value generated. Both instructions refer to the same address in the literal pool for their operand

→ There are two ways of recognising the duplicate literals

    (a) Compare the character strings defining them. Same literal name with different value

        Ex:  X 'os'

    (b) Compare the generated data value. This is better but increases the complexity of the assembler.

        Ex: = C'EDF' and = X'H5H5H6'

→ The problem of using character strings to recognize duplicate literals is, as we see '*' denotes a literal refers to the current value of program counter after line no. 93. There may be some literals that have the same name but different values.

    for example the statements

        BASE  *       ——①

        LDB  = *    ——②

      ① —— loads the beginning address of the program into register B. This value will be available later for base relative addressing.

x* ⟶ 1) causes a problem if we use of line no. 13

  ie  13   0003   LDB  = *    <u>692003</u>

  it specifies an operand with value 0003.

  55   0020   LDA  = *    010020

  ie  literal operands have identical names but they
  have different values and both must appear in the
  literal pool.

**x* ⟶** The same problem arises if a literal refers to any
  other item whose value changes between one point
  in the program and another.

⟶ The datastructure used to store literal operands is
  literal table <u>"LITTAB"</u>

⟶ literal Table ( LITTAB) : It is a hashtable using literal
                         name or value as the
                         key.

contains  ↳ literal name
          ↳ operand value
          ↳ operand length
          ↳ address assigned

| NAME | OPERAND VALUE | LENGTH | ADDRESS |
|------|---------------|--------|---------|
| = C 'EOF' | EOF | 03 | 002D |
| = X '05' | 05 | 01 | 1076 |

pass-1
⟶ Builds LITTAB with literal name, operand value and
  length, leaving the address unassigned

→ when LTORG statement is encountered, assign an address to each literal not yet assigned an address. Along with this, location counter is updated to reflect the number of bytes occupied by each literal

**Pass - 2**

→ Search LITTAB for each literal operand encountered to generate respective object code

→ Generate data values using BYTE or WORD statements

→ Generate modification record for literals that represent an address in the program.

**Difference between literal and an immediate operand**

| Literal (=) | Immediate operand (#) |
|---|---|
| 1. Literal is an assembler directive | Immediate is a machine recognizable data. |
| 2. The assembler generates the specified value as a constant at some other memory location. The address of this generated constant is used as target address for machine instruction | 2. Here value is assembled as part of machine instruction. |
| 3 Architectural support is required | Architectural support not required |
| 4. very slow since values are obtained from data memory | faster than literal ∵ data is within the instruction |
| 5. capable of storing large data | can't store larger data ∵ fullword is opcode, register |

55 0020 LDA #3 010003

## 2.3.2  Symbol - Defining statements

Most assemblers provide an assembler directive that allows the programmer to define symbols and specify their values.

The assembler directives are    a) EQU    b) ORG

### a)  EQU (Equate)

→  allows the programmer to define symbols (i.e enters it into SYMTAB) and assigns to it the specified value.

syntax:    Symbol    EQU    value ⟨ constant
                                          Expression
                                          previously defined symbol

→  The value may be

(i)  constant

(ii)  An expression involving constant

(iii)  previously defined symbols.

→  Uses of EQU

1) To establish symbolic names that can be used for improved readability in place of numeric values.

Ex:    +LDT    #4096 ;  load the value h096 into reg. T

⇓ replace with

        MAXLEN    EQU    H096
        +LDT    MAXLEN

→  when assembler encounters the EQU statement, it enters MAXLEN to SYMTAB with value H096.

→ During assembly of WDT instruction, the assembler searches the SYMTAB for the MAXLEN symbols and using its value as the operand in the instruction.

→ The advantage of doing so is if we want to change the value 4096 to some other value, we need to change in only one place ^MAXLEN instead of searching or scanning through the program for #4096 for the replacement (required change) ⇒ #define in c

2) To define mnemonic names for registers.

Ex:     A    EQU    0

        X    EQU    1

        L    EQU    2

→ The symbols A, X, L has to be entered into SYMTAB with their corresponding values 0,1,2.

→ instruction RMO A,X searches the SYMTAB for A and X and their values to assemble the instruction

3) To reflect the logical function of the registers

Ex:    BASE    EQU    R1

       COUNT   EQU    R2

       INDEX   EQU    R3

→ imply Register R1 is used as base register, R2 as program counter, R3 as index registers etc

→ forward reference is not allowed in EQU ic all terms in the value field must have been defined previously during pass-1

Ex:  ALPHA  RESW  1  } Allowed          BETA  EQU  ALPHA } not
     BETA   EQU   ALPHA                   ALPHA RESW  1   } allowed

b) <u>ORG (origin)</u>

→ Assembler directive used to indirectly assign values to symbols.

→ syntax :

ORG value

→ value can be

(i) constant

(ii) expression involving constant

(iii) Previously defined symbols

→ When ORG is encountered, the assembler resets its LOCCTR to the specified value.

→ Location counter is used to control assignment of storage in the object program. Hence altering its value would result in an incorrect assembly.
∴ the directive should be minimum used.

→ The ORG statement will affect the values of all labels defined until the next ORG.

→ If the previous value of LOCCTR is automatically remembered, then we can return to the normal use of LOCTR just by writing

<u>ORG</u>

→ Example: To define a symbol table with the following structure.

→ Symbol table with the given structure

| SYMBOL | VALUE | FLAGS |
|--------|-------|-------|
| STAB (100 entries) | | |
| | | |
| | | |
| | | |

6 bytes     3 bytes     2 byte
(1 word)

⌐> SYMBOL field contains user defined symbols.

⌐> VALUE field represents the value assigned to the symbol

⌐> FLAG field specifies symbol type and other information

⌐> The space for this table is reserved as

     STAB1    RESB     1100 ; $\underset{\text{entries}}{100} \times \underset{\text{each entry}}{11} = 1100$

⌐> we can access the label entries in two ways
(usage of EQU and ORG)

⌐> Using EQU ⟹

| SYMBOL | EQU | STAB |
|--------|-----|------|
| VALUE | EQU | STAB+6 |
| FLAGS | EQU | STAB+9 |

offset from STAB

(i) To fetch the value field,

     LDA     VALUE, X ; where X = 0, 11, 22, ... for each entry

       ⌐> index register

**\*** (ii) This method of definition simply defines the labels, it does not make the structure of the table as clear as it might be.

(iii) Therefore we make use of ORG.

↳ Using __ORG__ ⟹

```
            STAB      RESB    1100
                      ORG     STAB    ← set LOCCTR to STAB
            SYMBOL    RESB    ⎡ 6 ⎤
            VALUE     RESW    ⎢ 1 ⎥  ← Size of each
            FLAGS     RESB    ⎣ 2 ⎦     field
                      ORG     STAB + 1100  ← Restore
                                              LOCCTR
```

(i) The first ORG sets the location counter to the value of STAB

(ii) RESB statement defines SYMBOL to have the current value in LOCCTR

(iii) LOCCTR is then advanced, so the label on RESW statement assigns to VALUE to address (STAB + 6) and then advanced to assign to FLAGS to address (STAB + 9).

(iv) The last ORG statement sets LOCCTR back to its previous value, That which is the address of the next unassigned byte of memory after the table STAB.

→ Forward reference is not allowed in ORG ie all symbols used to specify the new location counter value have to be previously defined.

→ Example :

```
            ORG     ALPHA
    BYTE1   RESB    1
    BYTE2   RESB    1
```

```
BYTE3        RESB     1

             ORG

ALPHA        RESB     1
```

↳ cannot processed :. the assemble does not know what value has to be assigned to the location counter in response to the first ORG statement. The symbols BYTE1, BYTE2, BYTE3 are not assigned addresses during pass 1.

↳ It has to be written as

```
ALPHA    RESB      1
         ORG       ALPHA
BYTE1    RESB      1
BYTE2    RESB      1
BYTE3    RESB      1
         ORG
```

2.3.3. Expressions

→ The assembler allows the use of expressions as operand.

→ It calculates the expressions and produces a single operand address or value.

→ The expression consists of

(i) operators : + - * / ( division is usually defined to produce an integer value)

(ii) Individual terms : Constants, user-defined symbols, special terms like * ( current value of the location counter).

Ex:. MAXLEN   EQU   BUFEND - BUFFER

     STAB   RESB   (6+3+1) * MAX

     BUFEND   EQU   *

→ The values of terms can be absolute (independent of program location) such as constants or relative ( to the beginning of the program) such as Address labels, data areas, references to the location counter value.

→ Expressions are classified as

(i) Absolute Expressions } based on type of value
(ii) Relative Expressions } produced.

(i) **Absolute Expressions :**

→ Absolute means independent of program location and contains absolute terms like constants

→ It may also contain relative terms provided the relative terms occur in pair and the terms in each such pair have opposite signs.

→ It is not necessary that the paired terms be adjacent to each other in the expression however, all relative terms must be capable of being paired in this way.

→ None of the relative terms may enter into a multiplication or division operation.


(ii) **Relative Expressions :**

→ Relative means relative to the beginning of the program, such as labels on the instruction, data areas, references to the location counter value.

→ Here, all of the relative terms except one can be paired as in absolute expressions and the remaining unpaired relative term must have a positive sign.

→ no relative terms may enter into a multiplication or division operation.

note : If Either of absolute expression or relative expression do not meet the conditions, they are flagged as errors

→ A relative term or expression represents some value which is written as $(s+r)$

   • s = starting address of the program
   • r = value of term or expression relative to the starting address.

Ex: 1) MAXLEN EQU BUFEND − BUFFER

⤷ both BUFEND and BUFFER are relative terms representing an address within the program. The expression represents an absolute value.

(2) Illegal Expressions

   BUFEND + BUFFER   ; no opposite signs
   100 − BUFFER   ; both are not relative terms
   3 ∗ BUFFER   ; ∗ can't be used

→ Type of expression is determined by keeping track of symbol types in the program. This is done by adding a flag (R or A) in the SYMTAB for each symbol defined

Ex:

| Symbol | Type | Value |
|---|---|---|
| RETADR | R | 0030 |
| BUFFER | R | 0036 |
| BUFFND | R | 1036 |
| MAXLEN | A | 1000 |

} two symbols
} of fig 3.10

2.3.4   Program Blocks

→ Till now, we have seen that the program being assembled was treated as a single unit, eventhough it had subroutine, data areas etc resulting in a single block of object code.

→ within this object code (program) the generated machine instructions and data appeared in the same order as they were written in the source program.

→ But sometime it is required to logically rearrange the statements of the source program so that the 1)large buffer area can be moved to the end of object program, 2)no need of using extended instruction format, 3)the base register usage is not required, 4)the problem of placing literals in program has to be more flexible etc.

→ All these are achieved through some of the assembler features such as program blocks and control sections.

→ Program blocks : Allows the generated machine instructions and data to appear in the object program in a different order from the corresponding source statement.

or

Program blocks are segment of code that are rearranged within a single object program unit.

Assembler Directive : USE

syntax : USE BLOCKNAME

Fig. 2.11 shows the source program with program blocks

→ There are three blocks in the program

(1) Unnamed program block contains the executable instructions of the program

(b) CDATA program block contains all data areas that consists of few block of memory ie few coords or few in length

(iii) CBLKS program block contain all data areas that consists of larger blocks of memory.

→ At the beginning, statements are assumed to be part of the unnamed (default) block. If no USE statements are included, the entire program belongs to this single block.

→ USE on line 92 indicates the beginning of CDATA block

→ USE on line 103 indicates the beginning of CBLKS block

→ USE on line 123 resume the default block

→ Each program block may contain several separati segments of the source program but assembler will logically rearrange these segments to gather together the pieces of each block and assign address.

→ Program readability is better if data areas are placed in the source program close to the statements that reference them.

The assembler accomplishes this logical rearrangement of code by maintaining during pass-1 and pass 2

(i) Pass-1

Fig 3.13(b) shows the pass-1 of program blocks

→ A separate location counter for each program block is assigned and is assigned to ZERO when a program block begins

↳ Saving and Restoring the current value of LOCCTR when occurs whebe switching between blocks

↳ Each label is assigned an address relative to the start of the block.

↳ Stores the block name and number in the SYMTAB along with the assigned relative address of the label.

↳ At the end of pass-1, indicates the block length as the latest value of LOCCTR for each block.

↳ constructs a table which contains the starting address and length for all blocks.

↳ Assembler assigns to each block a starting address in the object program (beginning with relative location 0).

latest value of LC
↓

| Block name | Block Number | Address | Length | |
|---|---|---|---|---|
| default | 0 | 0000 | 0066 | 0063 + 0003 = 0066 |
| CDATA | 1 | 0066 | 000B | 0066 + 000B = 0071 |
| CBLKS | 2 | 0071 | 1000 | 1071 - 0071 = 1000 |

↳ Flag is also added in this table

Diagram showing memory blocks:
- 0000 to 0066: Default (Default block)
- 0066: ← CDATA starts loading
- 0066 to 0071: CDATA
- 0071: ← CBLKS starts loading
- 0071 to 1071: CBLKS

ii) Pass 2:

→ calculates the address for each symbol relative to the start of the object program (not the start of an individual program block) by adding

(i) The location of the symbol relative to the start of its block (from SYMTAB)

(ii) The starting address of this block.

## Example

| | | | | |
|---|---|---|---|---|
| 20 | 0006 | 0 | LDA | LENGTH |
| : | | | | |
| 92 | 0000 | 1 | | USE CDATA |
| 100 | 0003 | 1 | LENGTH RESW 1 |

→ The value of the operand LENGTH is 0003 relative to block 1 (CDATA)

∴ address = 0003 + 0066 = 0069 relative to program
  (TA)    when this instruction is executed

| Line | Loc/Block | | | Source statement | | Object code |
|---|---|---|---|---|---|---|
| 5 | 0000 | 0 | | COPY | START | 0 | |
| 10 | 0000 | 0 | FIRST | STL | RETADR | 172063 |
| 15 | 0003 | 0 | CLOOP | JSUB | RDREC | 4B2021 |
| 20 | 0006 | 0 | | LDA | LENGTH | 032060 |
| 25 | 0009 | 0 | | COMP | #0 | 290000 |
| 30 | 000C | 0 | | JEQ | ENDFIL | 332006 |
| 35 | 000F | 0 | | JSUB | WRREC | 4B203B |
| 40 | 0012 | 0 | | J | CLOOP | 3F2FEE |
| 45 | 0015 | 0 | ENDFIL | LDA | =C'EOF' | 032055 |
| 50 | 0018 | 0 | | STA | BUFFER | 0F2056 |
| 55 | 001B | 0 | | LDA | #3 | 010003 |
| 60 | 001E | 0 | | STA | LENGTH | 0F2048 |
| 65 | 0021 | 0 | | JSUB | WRREC | 4B2029 |
| 70 | 0024 | 0 | | J | @RETADR | 3E203F |
| 92 | 0000 | 1 | | USE | CDATA | |
| 95 | 0000 | 1 | RETADR | RESW | 1 | |
| 100 | 0003 | 1 | LENGTH | RESW | 1 | |
| 103 | 0000 | 2 | | USE | CBLKS | |
| 105 | 0000 | 2 | BUFFER | RESB | 4096 | |
| 106 | 1000 | 2 | BUFEND | EQU | * | |
| 107 | 1000 | | MAXLEN | EQU | BUFEND-BUFFER | |
| 110 | | | | . | | |
| 115 | | | | . | SUBROUTINE TO READ RECORD INTO BUFFER | |
| 120 | | | | . | | |
| 123 | 0027 | 0 | | USE | | |
| 125 | 0027 | 0 | RDREC | CLEAR | X | B410 |
| 130 | 0029 | 0 | | CLEAR | A | B400 |
| 132 | 002B | 0 | | CLEAR | S | B440 |
| 133 | 002D | 0 | | +LDT | #MAXLEN | 75101000 |
| 135 | 0031 | 0 | RLOOP | TD | INPUT | E32038 |
| 140 | 0034 | 0 | | JEQ | RLOOP | 332FFA |
| 145 | 0037 | 0 | | RD | INPUT | DB2032 |
| 150 | 003A | 0 | | COMPR | A,S | A004 |
| 155 | 003C | 0 | | JEQ | EXIT | 332008 |
| 160 | 003F | 0 | | STCH | BUFFER,X | 57A02F |
| 165 | 0042 | 0 | | TIXR | T | B850 |
| 170 | 0044 | 0 | | JLT | RLOOP | 3B2FEA |
| 175 | 0047 | 0 | EXIT | STX | LENGTH | 13201F |
| 180 | 004A | 0 | | RSUB | | 4F0000 |
| 183 | 0006 | 1 | | USE | CDATA | |
| 185 | 0006 | 2 | INPUT | BYTE | X'F1' | F1 |
| 195 | | | | . | | |
| 200 | | | | . | SUBROUTINE TO WRITE RECORD FROM BUFFER | |
| 205 | | | | . | | |
| 208 | 004D | 0 | | USE | | |
| 210 | 004D | 0 | WRREC | CLEAR | X | B410 |
| 212 | 004F | 0 | | LDT | LENGTH | 772017 |
| 215 | 0052 | 0 | WLOOP | TD | =X'05' | E3201B |
| 220 | 0055 | 0 | | JEQ | WLOOP | 332FFA |
| 225 | 0058 | 0 | | LDCH | BUFFER,X | 53A016 |
| 230 | 005B | 0 | | WD | =X'05' | DF2012 |
| 235 | 005E | 0 | | TIXR | T | B850 |
| 240 | 0060 | 0 | | JLT | WLOOP | 3B2FEF |
| 245 | 0063 | 0 | | RSUB | | 4F0000 |
| 252 | 0007 | 1 | | USE | CDATA | |
| 253 | | | | LTORG | | |
| | 0007 | 1 | | =C'EOF' | | 454F46 |
| | 000A | 1 | | =X'05' | | 05 |
| 255 | | | | END | FIRST | |

**Figure 2.12(a)** Program from Fig. 2.11 with object code.

```
Line            Source statement

   5   COPY      START    0            COPY FILE FROM INPUT TO OUTPUT
  10   FIRST     STL      RETADR       SAVE RETURN ADDRESS
  15   CLOOP     JSUB     RDREC        READ INPUT RECORD
  20             LDA      LENGTH       TEST FOR EOF (LENGTH = 0)
  25             COMP     #0
  30             JEQ      ENDFIL       EXIT IF EOF FOUND
  35             JSUB     WRREC        WRITE OUTPUT RECORD
  40             J        CLOOP        LOOP
  45   ENDFIL    LDA      =C'EOF'      INSERT END OF FILE MARKER
  50             STA      BUFFER
  55             LDA      #3           SET LENGTH = 3
  60             STA      LENGTH
  65             JSUB     WRREC        WRITE EOF
  70             J        @RETADR      RETURN TO CALLER
  92             USE      CDATA
  95   RETADR    RESW     1
 100   LENGTH    RESW     1            LENGTH OF RECORD
 103             USE      CBLKS
 105   BUFFER    RESB     4096         4096-BYTE BUFFER AREA
 106   BUFEND    EQU      *            FIRST LOCATION AFTER BUFFER
 107   MAXLEN    EQU      BUFEND-BUFFER  MAXIMUM RECORD LENGTH
 110   .
 115   .                  SUBROUTINE TO READ RECORD INTO BUFFER
 120   .
 123             USE
 125   RDREC     CLEAR    X            CLEAR LOOP COUNTER
 130             CLEAR    A            CLEAR A TO ZERO
 132             CLEAR    S            CLEAR S TO ZERO
 133            +LDT      #MAXLEN
 135   RLOOP     TD       INPUT        TEST INPUT DEVICE
 140             JEQ      RLOOP        LOOP UNTIL READY
 145             RD       INPUT        READ CHARACTER INTO REGISTER A
 150             COMPR    A,S          TEST FOR END OF RECORD (X'00')
 155             JEQ      EXIT         EXIT LOOP IF EOR
 160             STCH     BUFFER,X     STORE CHARACTER IN BUFFER
 165             TIXR     T            LOOP UNLESS MAX LENGTH
 170             JLT      RLOOP         HAS BEEN REACHED
 175   EXIT      STX      LENGTH       SAVE RECORD LENGTH
 180             RSUB                  RETURN TO CALLER
 183             USE      CDATA
 185   INPUT     BYTE     X'F1'        CODE FOR INPUT DEVICE
 195   .
 200   .                  SUBROUTINE TO WRITE RECORD FROM BUFFER
 205   .
 208             USE
 210   WRREC     CLEAR    X            CLEAR LOOP COUNTER
 212             LDT      LENGTH
 215   WLOOP     TD       =X'05'       TEST OUTPUT DEVICE
 220             JEQ      WLOOP        LOOP UNTIL READY
 225             LDCH     BUFFER,X     GET CHARACTER FROM BUFFER
 230             WD       =X'05'       WRITE CHARACTER
 235             TIXR     T            LOOP UNTIL ALL CHARACTERS
 240             JLT      WLOOP         HAVE BEEN WRITTEN
 245             RSUB                  RETURN TO CALLER
 252             USE      CDATA
 253             LTORG
 255             END      FIRST
```

**Figure 2.11**   Example of a program with multiple program blocks.

Disp : T0 - (PC)

  = 0069 - 0009

  = 0060

PC = 0000 + 0009 = 0009

( starting address of dy
    block ) + 0009 = 0009

$$\overset{n \quad i \quad x \quad b \quad p \quad e}{\boxed{0000 \ 0U} \ \boxed{1 \ 1 \ 0 \ 0 \ 1 \ 0} \ \boxed{060}} \longrightarrow 032060$$

SymTAB

| Label name | Block number | Address | Flag |
|---|---|---|---|
| Length | 1 | 0003 | |

note: line 107

1000    MAXLEN    EQU    BUFEND - BUFFER

→ shown without a block number indicates that
  MAXLEN is an absolute symbol, whose value is
  not relative to the start of any program block.

object Program :
  → It is not necessary to physically rearrange the
    generated code in the object program. The assembler
    just simply inserts the proper load address in
    each text record. The loader will load their
    codes into correct place

** 

  → Header record as before
  → Text records : the first 2 text records generated
    from lines through 70.

→ when <u>USE</u> statement on line 92 is encountered, the assembler writes the new Text record eventhough there is room (space) in the previous text record.

→ The process continues till the end of the program

H ∧ COPY ∧ 000000 ∧ 001071

T ∧ 000000 ∧ 1E ∧ 172063 ∧ HB 2021 ∧ · · · — · · · ∧ 010003

T ∧ 00001E ∧ 09 ∧ 0F2048 ∧ HB 2029 ∧ 3E203F

T ∧ 000027 ∧ 1D ∧ BH10 ∧ BH00 ∧ · · · · · · ∧ B850 ∧

T ∧ 0000 HLI ∧ 09 ∧ 3B9FEA ∧ 13201F ∧ HF0000

T ∧ 00006C ∧ 01 ∧ F1 · · · · ∧ HF0000

T ∧ 0000HD ∧ 19 ∧ BH10 ∧ 772017 ∧ · · ·

T ∧ 000060 ∧ 0H ∧ H5HFH6 ∧ 0-5

E ∧ 000000



Fig: Program blocks loaded in memory

```
begin
  block number = 0 LOCCTR[i] = 0 for all i
  read the first input line
  if OPCODE = 'START' then
  begin
    write line to intermediate file ;
    read next input line
  end {if START}
  while OPCODE ≠ 'END' do
  if OPCODE = 'USE'
  begin
    if there is no OPEREND name then
      set block name as default
    else block name as OPERAND name
    if there is no entry for block name then
      insert (block name, block number ++) in block table
    i = block number for block name
    if this is not a comment line then
      begin
      if there is a symbol in the LABEL field then
        begin
        search SYMTAB for LABEL
        if found then
          set error flag (duplicate symbol)
        else
          insert (LABEL, LOCCTR[i]) into SYMTAB
        end {if symbol}
      Search OPTAB for OPCODE
      if found then
        add 3 instruction length to LOCCTR[i]
      else if OPCODE = 'WORD' then
        add 3 to LOCCTR[i]
      else if OPCODE = 'RESW' then
        add 3 * #[OPERAND] to LOCCTR[i]
      else if OPCODE = 'RESB' then
        add #[OPERAND] to LOCCTR[i]
      else if OPCODE = 'BYTE' then
      begin
        find length of constant in bytes
        add length to LOCCTR[i]
      end {if byte}
  else
```

**Figure 2.12(b)**   Pass 1 of program blocks.

```
    Set error flag
    end {if not a comment}
  write line to intermediate file
  read Text input line
  end {while not END}
write last line to intermediate file
save Length[i] as LOCCTR[i] for all i
Address[o] = starting address
Address[i] = address(i - 1) + Length(i - 1)
              {for i = 1 to max(block number)}
insert(address[i], Length[i]) in block table for all i
end {Pass 1}
```

**Figure 2.12(b)** *(cont'd)*

```
If OPCODE = 'USE' then
  set block number for block name with OPERAND field
  search SYMTAB for OPERAND
  store symbol value + address [block number] as operand address
end {Pass 2}
```

**Figure 2.12(c)**  Pass 2 of program blocks.

## Loading

Assembler — object program → Loader

Loader → contacts os & os says start loading at particular address

Assembler generates the object program (Header Record, Text record, modify record and End record). Assembler interacts with loader through object program. loader contacts operating system to load at particular address. Then os checks if that particular space is free if so it starts loading. If not it will tell the loader to either wait or remove unnecessary space

∴ "loader loads the object program residing in hard disk to main memory and start executing".

→ when there is no enough space, somebody has to instruct the loader to change its address and try loading. It is done by assembler not by os.

∴ Assembler instructs the loader to change the address.

ie Line 15/35/65.          + JSUB RPRG C

line 15 →: address is HB10_1035_ at 0006.

      01036 starts from 0007 (middle). ie we
can access 1 byte but not 1 nibble.

   0006 - HB
   0007 → 10 → go here and modify the record. This
      is done by assembler ∴ we have modification
record m∧00000700.5

→ Loader should listen to both assembler and OS

→ Assembler says goto 0007 and modify but OS
   says it is loaded at 5000 ~~and~~ ∴ modify at 5007.

**→ + JSUB #4096 → doesnot need modification record
   ∴ it is immediate addressing. Irrespective of
   relocation it remains same.

→ Execution is part of microprocessor

Relocation
   → All instructions works except FH instructions.

E₁   + JSUB RPREL  HB101036

   [ h | 10110001 | 01036 ]

      TA = 01036 → loaded at 0000

** if it is loaded at 5000, it will not work properly.
loader stops functioning ∴ TA = 01036.
      + JSUB 06036    ∴ go for modification record
             X

→ start adding length as

• memory address - 3 bytes
• Register-to-register - 2 bytes
• Extended - 4 bytes

→ note: All literals should be placed where LTORG appears in the program if LTORG is not present all literals will be inserted at the end of the program. (line no 255)

→ At line to 95, block 1 (CDATA) starts ∴ It stores the LOCCTR value ie 0027 in LOCCTR-0. Then starts assigning 0000 to block 1.

→ At line 105, block 2 (CBLCK) starts ∴ it saves the LOCCTR value - 0006 in LOCCTR-1 column.

→ line 125, block 0 starts again. It restores the LOCCTR value 0027 and starts over till line no. 185 (LOCCTR =004D)

→ line 185, block-1 (CDATA) restarts by restoring the saved LOCCTR value & saves LOCCTR = 0007

→ line 210, restores 004D and starts over till line no 245 having LOCCTR = 0066 which is stored in LOCCTR-0 column

# THREE LOCATION COUNTERS

| LOCCTR-0 (default block) | LOCCTR-1 (CDATA Block) | LOCCTR-2 (CBLKS) |
|---|---|---|
| 0000 | 0000 | 0000 |
| 0027 | 0006 | 1000 |
| 004D | 0007 | |
| 0066 | 000B | |

# LITERAL TABLE

| Literal name | value of literal | length of literal | address of literal |
|---|---|---|---|
| = C'EOF' | 454F46 | 03 | 0007 |
| = X'05' | 05 | 01 | 000A |
| | | | 000B |

# BLOCK TABLE

| Block name | Block number | Address | Length |
|---|---|---|---|
| Default | 0 | 0000 | 0066 |
| CDATA | 1 | 0066 | 000B |
| CBLks | 2 | 0071 | 1000 |

Program Length
= 66 + 0B + 1000
= 1071

# SYMBOL TABLE (with block number)

| Symbol name | value | Block no | Symbol name | value | Block no |
|---|---|---|---|---|---|
| FIRST | 0000 | 0 | MAXLEN | 1000 | |
| CLOOP | 0003 | 0 | RDREC | 0027 | 0 |
| ENDFIL | 0015 | 0 | RLOOP | 0031 | 0 |
| RETADR | 0000 | 1 | EXIT | 0047 | 0 |
| LENGTH | 0003 | 1 | INPUT | 0006 | 1 |
| BUFFER | 0000 | 2 | CORREC | 004D | 0 |
| BUFEND | 1000 | 2 | WLOOP | 0052 | 0 |

→ line 253, we have USE CDATA (Block 1) and

253 → LTORG ∴ store LOCCTR = 0007 at line

253

ie    253    0007                LTORG

***↓***

all literals should be placed where LTORG

appears in the program. we have two literals

here  ie    = C 'EOF'   and   = X 'O5'   ⇒ 4 bytes

3 Byte                1 Byte

starts  at   0007 , 0008, 0009, 000A

(E)      (O)     (F)    (05)

so it stores   000B  in LOCCTR-1 column

→ line 105 , block 2 starts and it reserve 1000

bytes q memory ∴ It saves 1000 in LOCCTR-2

column

note→    BUFEND    EQU  100 ⇒ value q buffend is 100

BUFEND    EQU  *  ⇒ value will be current location

value   = 0000 + 1000 = 1000

Stores these value in literal table.

10    0000    0    ` STL    RETADR

⤷ present in block 1 : add

usr of block 0 (default block)

∴ displacement = sizeq (previous block) + TA - PC

= sizeq (B0) + RETADR - PC

= 0066 + 0000 - 0003

= 0063



n  i  x  b  p  e
1 | 0 | 1 | 1 | 1 | 0 0 | 1 | 0 | 063

1    7    2    063    ⟹ 172063

15    0003    0    JSUB    RDREC
                            ⤷ Block (0)

disp = sizeq previous block + TA - PC

= 0 + 0027 - 006 = 021



n  i  x  b  p  e
4 | 1 | 0 | 1 | 1 | 0 0 | 1 | 0 | 021

4    B    2 021    ⟹ 4B2021

20    0006    0         LDA       LENGTH
                                  ‿‿‿‿
                                  block 1

disp: size of previous block + TA - PC

= 0066 + 0003 - 0009 = 0060



0    3    2060    = 032060


35    000F            JSUB      CORREC
                      ‿‿        ‿‿‿‿‿ belongs to block 0
                      48

⟹ As before


** 45    0015            LDA       = C'EOF'
                                  ‿‿‿‿‿
                                    ↓
                        belongs to CDATA not to default block
                        since it is literal which is placed after
                        
                        LTORG.

disp = size of previous block + TA - PC

     = size of B0 + c'EOF' - PC

     = 0066 + 0007 - 0018 = 0055



0    3    2-055    ⟹    032055

.50        001B           STA   BUFFER

                    $\underset{01}{\underbrace{\phantom{xx}}}$   $\underbrace{\phantom{xxxxx}}$  belongs to block 2 (CBLKS)

$$disp = size \ of \ (B_0 + B_1) + BUFFER - PC$$

$$= size \ of \ (default \ block + CDATA) + BUFFER - PC$$

$$= 0066 + 000B + 0000 - 001B$$

$$= 0071 - 001B = \not{0}056$$



$$0 \quad F \qquad 2 \ 056 \implies 0F2056$$

objeet Program

HCOPY    ^ 000000 ^ 001071

1   T^ 000000 ^16^ 172063 ^ 4B2021^D32060^ 290000^332006 ^ 4B203B^ 3F2FEE

                                    ^ 032055^ 0F2056^ 010003

2   T^ 00001E^ 09 ^0F2049^ 4B2029^ 3E203F

3   T^ 000027 ^1D^ B410 ^ B410^ B400^ B440 ^ 75101000 ^ E32039^ 332FFA ^

                    ^ DB2032^ A004 ^332008 ^ 57A02F ^B850

4   T^ 000044^09^ 3B2FEA ^13201F^ 4F0000

5   T^ 00006C ^ 01^F1

6   T^ 000010 ^ 19^ B410^ 772017^ E3201B^ 332FFA^ 53A016^ DF2012^ B850

                                    ^ 3F2FEE^ 4F0000

7   T^ 00006D ^04^ 454F46^05

    E^ 000000

# Loading the object program into memory

| Address | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0000 | 17 | 20 | 63 | 4B | 20 | 21 | 03 | 20 | 60 | 29 | 00 | 0D | 33 | 20 | 06 | 4B | copy H(i) T1+T2 |
| 0010 | 20 | 3B | 3F | 2F | EE | 03 | 20 | 55 | 0F | 20 | 56 | 01 | 00 | 03 | 0F | 20 | |
| 0020 | 48 | 4B | 20 | 29 | 3E | 20 | 3F | B4 | 10 | B4 | 00 | B4 | 40 | 75 | 10 | 10 | copy (b) T2+Th |
| 0030 | 60 | 63 | 20 | 38 | 33 | 2F | FA | DB | 20 | 32 | A0 | 04 | 33 | 20 | 08 | 57 | |
| 0040 | A0 | 2F | B9 | 50 | 3B | 2F | EA | 13 | 20 | 1F | 4F | 00 | 00 | B4 | 10 | 17 | copy (c) T0 |
| 0050 | 20 | 17 | 63 | 20 | 1B | 33 | 2F | FA | 53 | A0 | 16 | 0F | 20 | 12 | B8 | 50 | |
| 0060 | 3B | 2F | EE | 4F | 00 | 00 | RETADR | | LENGTH | | | F1 | 45 | 4F | 46 | | |
| 0070 | 05 | | | | | | | | DATA(i) | | | copy by | | | | | |
| 0080 | COPY | | | | | | | | | | | | | | | | |
| 0090 | | | | | | | BUFFER | | | | | | | | | | |
| ⋮ | | | | | | | | | | | | | | | | | |
| 1050 | | | | | | | | | | | | | | | | | |
| 1060 | | | | | | | | | | | | | | | | | |
| 1070 | | | | | | | | | | | | | | | | | |

How does microprocessor execute an instruction

Ex:  $10^{0000}$  STL    RETADR    172063        L = 666600

$\Rightarrow$ store the contents of linkage register into RETADR location

$\Rightarrow$ opcode for STL = 14 (known by microprocessor)

| 1 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 06 |
|---|---|---|---|---|---|---|---|---|---|---|---|

n i x b p e

1  7  2063

$TA = PC + disp$
$= 0003 + 063 = 0066$
RETADR

ie    STL    0066 $\Rightarrow$ copy the contents of linkage register (666600) into address 0066

note :- If starting address (location) is changed from
0000 to 5000, it works as usual :- we don't
have format in address. ∴ no need of modification
reused as before. ⟶ advantage of program block.

2.3.5    Control Sections and Program linking

→ Control section is a part of the program that contains its identity after assembly.

→ Each control section can be loaded and relocated independently of the other

→ Control sections are usually used for subroutines or other logical subdivisions of a program.

→ The programmer can assemble, load and manipulate each of these control section **separately**

→ It uses a assembler directive : CSECT which indicates the beginning of the control section. where each control section starts its location too counter separately

→ when control sections form logically related parts of a program, it is necessary to provide some means for linking them together. This is because instructions in one control section may need to refer to instructions or data located in another control section. and assembler has no idea where exactly the control sections will be located at execution time.

→ such references between control sections are called external references.

→ The assembler generate information for each external reference that will allow the loader to perform the required linking.

→ There are two types of exter reference external symbols

(1) External Definition ( EXTDEF )
   • Symbols that are defined in one section and are used by other sections

   syntax: EXTDEF name [;name]

   Ex: EXTDET BUFFER, BUFEND

(2) External Reference (EXTREF)
   • Symbols that are used in this control sections but are defined in some other control sections.

   syntax: EXTREF name [,name]

   Ex: EXTREF RDREC, WRREC

note: To reference a external symbols, extended format instruction (Format 4) is needed

| Line | Loc | Source statement | | Object code | |
|------|-----|------|------|------|------|
| 5 | 0000 | COPY | START | 0 | |
| 6 | | | EXTDEF | BUFFER,BUFEND,LENGTH | |
| 7 | | | EXTREF | RDREC,WRREC | |
| 10 | 0000 | FIRST | STL | RETADR | 172027 |
| 15 | 0003 | CLOOP | +JSUB | RDREC | 4B100000 |
| 20 | 0007 | | LDA | LENGTH | 032023 |
| 25 | 000A | | COMP | #0 | 290000 |
| 30 | 000D | | JEQ | ENDFIL | 332007 | 𝑇ᵢ (1D) |
| 35 | 0010 | | +JSUB | WRREC | 4B100000 |
| 40 | 0014 | | J | CLOOP | 3F2FEC |
| 45 | 0017 | ENDFIL | LDA | =C'EOF' | 032016 |
| 50 | 001A | | STA | BUFFER | 0F2016 |
| 55 | 001D | | LDA | #3 | 010003 |
| 60 | 0020 | | STA | LENGTH | 0F200A | 𝑇₂ |
| 65 | 0023 | | +JSUB | WRREC | 4B100000 |
| 70 | 0027 | | J | @RETADR | 3E2000 |
| 95 | 002A | RETADR | RESW | 1 | |
| 100 | 002D | LENGTH | RESW | 1 | |
| 103 | | | LTORG | | |
| | 0030 | * | =C'EOF' | | 454F46 | 𝑇₈ (9) |
| 105 | 0033 | BUFFER | RESB | 4096 | |
| 106 | 1033 | BUFEND | EQU | * | |
| 107 | 1000 | MAXLEN | EQU | BUFEND-BUFFER | |
| 109 | 0000 | RDREC | CSECT | | |
| 110 | | . | | | |
| 115 | | . | SUBROUTINE TO READ RECORD INTO BUFFER | | |
| 120 | | . | | | |
| 122 | | | EXTREF | BUFFER,LENGTH,BUFEND | |
| 125 | 0000 | | CLEAR | X | B410 |
| 130 | 0002 | | CLEAR | A | B400 |
| 132 | 0004 | | CLEAR | S | B440 |
| 133 | 0006 | | LDT | MAXLEN | 77201F |
| 135 | 0009 | RLOOP | TD | INPUT | E3201B |
| 140 | 000C | | JEQ | RLOOP | 332FFA | 𝑇ᵢ (19) |
| 145 | 000F | | RD | INPUT | DB2015 |
| 150 | 0012 | | COMPR | A,S | A004 |
| 155 | 0014 | | JEQ | EXIT | 332009 |
| 160 | 0017 | | +STCH | BUFFER,X | 57900000 |
| 165 | 001B | | TIXR | T | B850 |
| 170 | 001D | | JLT | RLOOP | 3B2FE9 |
| 175 | 0020 | EXIT | +STX | LENGTH | 13100000 |
| 180 | 0024 | | RSUB | | 4F0000 | 𝑇₂ (06) |
| 185 | 0027 | INPUT | BYTE | X'F1' | F1 |
| 190 | 0028 | MAXLEN | WORD | BUFEND-BUFFER | 000000 |
| 193 | 0000 | WRREC | CSECT | | |
| 195 | | . | | | |
| 200 | | . | SUBROUTINE TO WRITE RECORD FROM BUFFER | | |
| 205 | | . | | | |
| 207 | | | EXTREF | LENGTH,BUFFER | |
| 210 | 0000 | | CLEAR | X | B410 |
| 212 | 0002 | | +LDT | LENGTH | 77100000 |
| 215 | 0006 | WLOOP | TD | =X'05' | E32012 |
| 220 | 0009 | | JEQ | WLOOP | 332FFA |
| 225 | 000C | | +LDCH | BUFFER,X | 53900000 | 𝑇ᵢ (1C) |
| 230 | 0010 | | WD | =X'05' | DF2008 |
| 235 | 0013 | | TIXR | T | B850 |
| 240 | 0015 | | JLT | WLOOP | 3B2FEE |
| 245 | 0018 | | RSUB | | 4F0000 |
| 255 | | | END | FIRST | |
| | 001B | * | =X'05' | | 05 |

**Figure 2.16** Program from Fig. 2.15 with object code.

| Line | | Source statement | | |
|------|-----|-----|-----|-----|
| 5 | COPY | START | 0 | COPY FILE FROM INPUT TO OUTPUT |
| 6 | | EXTDEF | BUFFER,BUFEND,LENGTH | |
| 7 | | EXTREF | RDREC,WRREC | |
| 10 | FIRST | STL | RETADR | SAVE RETURN ADDRESS |
| 15 | CLOOP | +JSUB | RDREC | READ INPUT RECORD |
| 20 | | LDA | LENGTH | TEST FOR EOF (LENGTH = 0) |
| 25 | | COMP | #0 | |
| 30 | | JEQ | ENDFIL | EXIT IF EOF FOUND |
| 35 | | +JSUB | WRREC | WRITE OUTPUT RECORD |
| 40 | | J | CLOOP | LOOP |
| 45 | ENDFIL | LDA | =C'EOF' | INSERT END OF FILE MARKER |
| 50 | | STA | BUFFER | |
| 55 | | LDA | #3 | SET LENGTH = 3 |
| 60 | | STA | LENGTH | |
| 65 | | +JSUB | WRREC | WRITE EOF |
| 70 | | J | @RETADR | RETURN TO CALLER |
| 95 | RETADR | RESW | 1 | |
| 100 | LENGTH | RESW | 1 | LENGTH OF RECORD |
| 103 | | LTORG | | |
| 105 | BUFFER | RESB | 4096 | 4096-BYTE BUFFER AREA |
| 106 | BUFEND | EQU | * | |
| 107 | MAXLEN | EQU | BUFEND-BUFFER | |
| | | | | |
| 109 | RDREC | CSECT | | |
| 110 | . | | | |
| 115 | . | | SUBROUTINE TO READ RECORD INTO BUFFER | |
| 120 | . | | | |
| 122 | | EXTREF | BUFFER,LENGTH,BUFEND | |
| 125 | | CLEAR | X | CLEAR LOOP COUNTER |
| 130 | | CLEAR | A | CLEAR A TO ZERO |
| 132 | | CLEAR | S | CLEAR S TO ZERO |
| 133 | | LDT | MAXLEN | |
| 135 | RLOOP | TD | INPUT | TEST INPUT DEVICE |
| 140 | | JEQ | RLOOP | LOOP UNTIL READY |
| 145 | | RD | INPUT | READ CHARACTER INTO REGISTER A |
| 150 | | COMPR | A,S | TEST FOR END OF RECORD (X'00') |
| 155 | | JEQ | EXIT | EXIT LOOP IF EOR |
| 160 | | +STCH | BUFFER,X | STORE CHARACTER IN BUFFER |
| 165 | | TIXR | T | LOOP UNLESS MAX LENGTH |
| 170 | | JLT | RLOOP | HAS BEEN REACHED |
| 175 | EXIT | +STX | LENGTH | SAVE RECORD LENGTH |
| 180 | | RSUB | | RETURN TO CALLER |
| 185 | INPUT | BYTE | X'F1' | CODE FOR INPUT DEVICE |
| 190 | MAXLEN | WORD | BUFEND-BUFFER | |
| | | | | |
| 193 | WRREC | CSECT | | |
| 195 | . | | | |
| 200 | . | | SUBROUTINE TO WRITE RECORD FROM BUFFER | |
| 205 | . | | | |
| 207 | | EXTREF | LENGTH,BUFFER | |
| 210 | | CLEAR | X | CLEAR LOOP COUNTER |
| 212 | | +LDT | LENGTH | |
| 215 | WLOOP | TD | =X'05' | TEST OUTPUT DEVICE |
| 220 | | JEQ | WLOOP | LOOP UNTIL READY |
| 225 | | +LDCH | BUFFER,X | GET CHARACTER FROM BUFFER |
| 230 | | WD | =X'05' | WRITE CHARACTER |
| 235 | | TIXR | T | LOOP UNTIL ALL CHARACTERS |
| 240 | | JLT | WLOOP | HAVE BEEN WRITTEN |
| 245 | | RSUB | | RETURN TO CALLER |
| 255 | | END | FIRST | |

**Figure 2.15** Illustration of control sections and program linking.

In Fig 2.16, there are three control sections.

1) man program → COPY from line 5 to line 107

2) read subroutine → RDREC from line 109 to line 170

3) write subroutine → WRREC from line 193 to 255.

→ Assembler establishes a separate location counter (beginning at 0) for each control section

→ Control section named COPY, RDREC, WRREC are not named in EXTDEF because they are automatically considered to be external symbols.

→ Assembler handles the external references as follows

a)   15    0003    CLOOP    $+\underbrace{JSUB}_{48}$    $\underbrace{RDREC}_{EXTREF}$

• Assembler is unaware of RDREC address, so it inserts an address of 3000 and pass this to loader, which is taken take care during loading.

• The address of RDREC will have no predictable relationship to anything in the control section by name COPY, therefore relative addressing is not possible. Thus an extended format instruction must be p used to provide room for the actual address to be inserted.

nibble table

h | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 00000

h    B    1    00000

b)  160   0017   +STCH   BUFFER, X      57900000

→ BUFFER is used in RDREC control section

but defined in <u>COPY</u> control section.

c)  190   0028   MAXLEN   WORD   BUFEND-BUFFER   000000

→ two external references in the expression
   BUFEND and BUFFER in WRREC section.

→ The assembler inserts an address of 300,
   it passes information to the loader to add
   to this data area the address of BUFEND
   and subtract from this data area the
   address of BUFFER, which results in the
   desired value.

d)  107   1000   MAXLEN   EQU   BUFEND-BUFFER

→ Both expressions looks same but the
        (107 &190)
   difference is here BUFEND and BUFFER
   are defined and used in the same control
   section. so value can be calculated
   immediately.

   MAXLEN   EQU   1033 - 0033 .
                = 1000

→ A reference to MAXLEN in the copy control section will use the definition on line 107, whereas a reference to MAXLEN in RDREC control section will use the definition on line 190.

→ object program

. Along with header record, text record, modification record, two more records are added

a) **Define** Record → information of symbols defined in this control section

| col. 1 | D |
|--------|---|
| col. 2-7 | Name of external symbol defined on this control section |
| col. 8-13 | Relative address within this control section (Hexadecimal) |
| col. 14-73 | Repeat information in col. 8-13 for other external symbol |

b) **Refer** Record → symbols that are used as external references in this control section

| col. 1 | R |
|--------|---|
| col 2-7 | Name of external symbol referred in this control section |
| col. 8-73 | Name of other external reference symbol |

c) **modification** Record

| col. 1 | m |
|--------|---|
| col 2-7 | starting address of the field modified, in half-bytes (hexadecimal) |
| col. 11-16 | External symbol whose value is to be added to, or subtracted from the indicated field |

COPY

H∧COPY ∧000000∧001033

D∧BUFFER∧000033∧BUFEND∧001033∧LENGTH∧000026

R∧RDREC ∧WRREC

T∧000000∧1D∧172027∧4B100000∧032023∧290000∧332007∧4B100000∧3F2FEC∧
032016∧0F2016

T∧00001D∧0D∧0100003∧0F200A∧4B100000∧3E2000

T∧000030∧03∧454F46

m∧000004∧05∧+RDREC

m∧000011∧05∧+WRREC

m∧000024∧05∧+WRREC

E∧000000

RDREC

H∧RDREC ∧000000∧00002B

R∧BUFFER∧LENGTH∧BUFEND

T∧000000∧1D∧B410∧B400∧B440∧77201F∧E3201B∧332FFA∧DB2015∧A004∧332009∧
579000000∧B850

T∧00001D∧0E∧3B2FE9∧131000000∧4F0000∧F1∧000000

m∧000018∧05∧+BUFFER

m∧000021∧05∧+LENGTH

m∧000028∧06∧+BUFEND

m∧000028∧06∧−BUFFER

E

WRREC

H∧WRREC ∧000000∧00001C

R∧LENGTH∧BUFFER

T∧000000∧1C∧B410∧77100000∧E32012∧332FFA∧53900000∧DF2008∧B850∧
∧3B2FEE∧4F000005

m∧000003∧05∧+LENGTH

m∧00000D∧05∧+BUFFER

E

3.4 Assembler Design Options

. we will learn two alternatives of the standard two-pass assembler.

a) One-pass Assemblers

b) Multi-pass Assemblers

a) One-pass Assemblers

→ As we know already, assembling the forward references is very difficult ∵ we don't know the address. This can be eliminated very easily for data items. That is data items are defined in the source program before they are referenced.

( Ei. A~~ssign~~ declaration of data variable in c).

→ It is not the same for labels on instructions.

Ei: If the program has a forward jump ie escaping from a loop after testing some condition → here we can't define before itself. Therefore the assembler has to provide the way to handle forward references.

→ Two types of one-pass assemblers.

(1) load-go assembler : Assembler produces object program code directly in memory for immediate execution.

→ Here object program is not written out and no loader is needed.

→ Application : program development and testing.

→ Ei :- A university computing system for student. Here a large fraction of the total workload consists of program translation. Because programs are re-assembled nearly every time they are run, efficiency of the assembly process is an important consideration.

→ load-and-go assembler avoids the overhead of writing the object program out (secondary storage) and reading it back in. ∴ forward references can be handled easily.

→ Avoids usage of forward references

→ Used on systems where external working-storing devices are either slow or not available.

| Line | Loc | Source statement | | | Object code |
|------|------|------|------|------|------|
| 2 | 1000 | COPY | START | 1000 | |
| 1 | 1000 | EOF | BYTE | C'EOF' | 454F46 |
| 2 | 1003 | THREE | WORD | 3 | 000003 |
| 3 | 1006 | ZERO | WORD | 0 | 000000 |
| 4 | 1009 | RETADR | RESW | 1 | |
| 5 | 100C | LENGTH | RESW | 1 | |
| 6 | 100F | BUFFER | RESB | 4096 | |
| 9 | | . | | | |
| 10 | 200F | FIRST | STL | RETADR | 141009 |
| 15 | 2012 | CLOOP | JSUB | RDREC | 48203D |
| 20 | 2015 | | LDA | LENGTH | 00100C |
| 25 | 2018 | | COMP | ZERO | 281006 |
| 30 | 201B | | JEQ | ENDFIL | 302024 |
| 35 | 201E | | JSUB | WRREC | 482062 |
| 40 | 2021 | | J | CLOOP | 302012 |
| 45 | 2024 | ENDFIL | LDA | EOF | 001000 |
| 50 | 2027 | | STA | BUFFER | 0C100F |
| 55 | 202A | | LDA | THREE | 001003 |
| 60 | 202D | | STA | LENGTH | 0C100C |
| 65 | 2030 | | JSUB | WRREC | 482062 |
| 70 | 2033 | | LDL | RETADR | 081009 |
| 75 | 2036 | | RSUB | | 4C0000 |
| 110 | | . | | | |
| 115 | | . | | SUBROUTINE TO READ RECORD INTO BUFFER | |
| 120 | | . | | | |
| 121 | 2039 | INPUT | BYTE | X'F1' | F1 |
| 122 | 203A | MAXLEN | WORD | 4096 | 001000 |
| 124 | | . | | | |
| 125 | 203D | RDREC | LDX | ZERO | 041006 |
| 130 | 2040 | | LDA | ZERO | 001006 |
| 135 | 2043 | RLOOP | TD | INPUT | E02039 |
| 140 | 2046 | | JEQ | RLOOP | 302043 |
| 145 | 2049 | | RD | INPUT | D82039 |
| 150 | 204C | | COMP | ZERO | 281006 |
| 155 | 204F | | JEQ | EXIT | 30205B |
| 160 | 2052 | | STCH | BUFFER,X | 54900F |
| 165 | 2055 | | TIX | MAXLEN | 2C203A |
| 170 | 2058 | | JLT | RLOOP | 382043 |
| 175 | 205B | EXIT | STX | LENGTH | 10100C |
| 180 | 205E | | RSUB | | 4C0000 |
| 195 | | . | | | |
| 200 | | . | | SUBROUTINE TO WRITE RECORD FROM BUFFER | |
| 205 | | . | | | |
| 206 | 2061 | OUTPUT | BYTE | X'05' | 05 |
| 207 | | . | | | |
| 210 | 2062 | WRREC | LDX | ZERO | 041006 |
| 215 | 2065 | WLOOP | TD | OUTPUT | E02061 |
| 220 | 2068 | | JEQ | WLOOP | 302065 |
| 225 | 206B | | LDCH | BUFFER,X | 50900F |
| 230 | 206E | | WD | OUTPUT | DC2061 |
| 235 | 2071 | | TIX | LENGTH | 2C100C |
| 240 | 2074 | | JLT | WLOOP | 382065 |
| 245 | 2077 | | RSUB | | 4C0000 |
| 255 | | | END | FIRST | |

**Figure 2.18** Sample program for a one-pass assembler.

iv) Assembler generating object code:

Working process of load go assembler

↳ fig 2.18 shows an example for one-pass assembler

↳ Here all data item definitions are placed ahead of the code that references them.

Ex:
```
01   1000   EOF     BYTE   c 'EOF'
02   1003   THREE   WORD   3
03   1006   ZERO    WORD   0
              .
06   100F   BUFFER  RESB   4096.
```

↳ The assembler generates object code as it scans the source program

↳ If an instruction operand is a symbol that has not yet been defined (forward reference), the operand address is ommitted during assembly.

↳ The symbol is entered into symbol table if not exists along with a flag indicating undefined.

↳ The address of the operand field of the instruction that refers to the undefined symbol is added to a list of forward references associated with the symbol table entry.

Ex:
```
15   2013   CLOOP   JSUB   RDREC   4800000
```
2013 → 48
2013 → ....

| RDREC | * | • |→ | 2013 | 0 |
|-------|---|---|---|------|---|
| LENGTH | 100C | | | | |

↳ address where it has to start later

↳ when the definition for a symbol is encountered, the forward reference list for that symbol is scanned (if exists) and the proper address is inserted into any instructions previously generated.

↳ Fig 2.19 shows the object code and symbol table entries as they were scanned. till line no. 40

↳ 15  2012  CLOOP  JSUB  RDREC
         ~~undefined~~
                              ↳ undefined

∴ symbol table entry is    | RDREC | * | → | 2013 | 0 |

2013 is the address location where it has to load once found further.

↳ Same for line no. 30, 35.

| memory Address | Contents | | | |
|---|---|---|---|---|
| 1000 | HSHFHC00 | 00030000 | C0XXXXXX | XXXXXXXX |
| 1010 | XXXXXXXX | XXXXXXX | XXXXXXXX | XXXXXXXY |
| ⋮ | | | | |
| 2000 | XXXXXXXN | XXXXXXN | XXXXXXXX | XXXXXYXh |
| 2010 | 1009H8-- | --00100C | 28100630 | ----H8-- |
| 2020 | --3C2012 | | | |
| ⋮ | | | | |

SYMBOL TABLE

| | | |
|---|---|---|
| LENGTH | 100C | |
| RDREC | * | → 2013 0 |
| THREE | 1003 | |
| ZERO | 1006 | |
| WRREC | * | → 201F 0 |
| EOF | 1000 | |
| ENDFIL | * | → 201C 0 |
| RETADR | 1009 | |
| BUFFER | 100F | |
| CLOOP | 2012 | |
| FIRST | 200F | |

Fig 2.19: Object code in memory and symbol table entries after scanning line no 40

Memory
Address                                    Contents                                              SYMBOL TABLE

| Memory Address | Contents | | | | Symbol Table | |
|---|---|---|---|---|---|---|
| 1000 | H5HFH600 | 0003000C | 00XXXXXX | XXXXXXXX | LENGTH | 100C |
| 1010 | XXXXXXXX | XXXXXXXX | XXXXXXXX | XXXXXXXX | RDREC | 203D |
| . | | | | | THREE | 1003 |
| . | | | | | ZERO | 1006 |
| 2000 | XXXXXXXX | XXXXXXX | XXXXXXXX | XXXXXXXH | WRREC | * • → 201F |
| 2010 | 1009H820 | 3D00100C | 28100630 | 2024H8-- | BOF | 1000 |
| 2020 | --3C2012 | 00100000 | 100F0010 | 0300100C | ENDFIL | 2024 |
| 2030 | H8----08 | 1009HC00 | 00F10010 | 0004100C | RETADR | 1009 |
| 2040 | 00100660 | 20893050 | H31852039 | 28100630 | BUFFER | 100F |
| 2050 | -----SH90 | 0F | | | CLOOP | 2012 |
| | | | | | FIRST | 200F |
| | | | | | MAXLEN | 203A |
| | | | | | INPUT | 2037 |
| | | | | | EXIT | * • → 2050 |
| | | | | | RLOOP | 2043 |

Fig : object code in memory and symbol table entries after scanning line 160 of fig 2.18

| Memory Address | Contents | | | | Symbol Table | |
|---|---|---|---|---|---|---|
| 1000 | H5HFH600 | 0003000D | 00XXXXXX | XXXXXXXX | LENGTH | 100C |
| 1010 | XXXXXXXX | XXXXXXXX | XXXXXXXX | XXXXXXXX | RDREC | 203D |
| . | | | | | THREE | 1003 |
| . | | | | | ZERO | 1006 |
| 2000 | XXXXXXXX | XXXXXXXX | XXXXXXXX | XXXXXXXH | WRREC | 2062 |
| 2010 | 1009H820 | 3D00100C | 28100630 | 2024H820 | EOF | 1000 |
| 2020 | 63C2012 | 0010000C | 100F0010 | 0300100C | ENDFIL | 2024 |
| 2030 | H8206208 | 1009HC00 | 00F10010 | 000H1006 | RETADR | 1009 |
| 2040 | 00100GE0 | 20393030 | H3D82039 | 28100630 | BUFFER | 100F |
| 2050 | 205B5H90 | 0F2C203A | 38504310 | 100CH100 | CLOOP | 2012 |
| 2060 | 00050H10 | 06E02061 | 30206550 | 9001DC30 | FIRST | 200F |
| 2070 | 612C100C | 38206SHC | 0000 | | MAXLEN | 203A |
| 2080 | | | | | INPUT | 2037 |
| | | | | | EXIT | 205B |
| | | | | | RLOOP | 2043 |
| | | | | | OUTPUT | 2061 |
| | | | | | WLOOP | 2065 |

Fig : complete object program in memory and symbol table entries

note: If any symbols in the SYMTAB are still marked with * should be flagged by the assembler as errors. (undefined symbol error in c long after compiling completely)

↳ The assembler searches SYMTAB for the value of the symbol named END statement and jumps to this location to begin execution of the assembled program.

↳ In load-and-go assembler, the actual address must be known at assembly time.

```
begin
  read first input line
  if OPCODE = 'START' then
    begin
      save #[OPERAND] as starting address
      initialize LOCCTR as starting address
      read next input line
    end (if START)
  else
    initialize LOCCTR to 0
while OPCODE ≠ 'END' do
    begin
      if there is not a comment line then
        begin
          if there is a symbol in the LABEL field then
            begin
              search SYMTAB for LABEL
                if found then
              begin
                  if symbol value as null
                  set symbol value as LOCCTR and search
                      the linked list with the corresponding
                      operand
                  PTR addresses and generate operand
                      addresses as corresponding symbol
                      values
                  set symbol value as LOCCTR in symbol
                      table and delete the linked list
              end
              else
                  insert (LABEL, LOCCTR) into SYMTAB
            end
              search OPTAB for OPCODE
                if found then
                  begin
                    search SYMTAB for OPERAND address
                  if found then
                    if symbol value not equal to null then
                      store symbol value as OPERAND address
                    else
                      insert at the end of the linked list
                        with a node with address as LOCCTR
                  else
                      insert (symbol name, null)
```

Figure 2.19(c)   Algorithm for One pass assembler.

```
                        add 3 to LOCCTR
                    end
                else if OPCODE = 'WORD' then
                    add 3 to LOCCTR & convert comment to
                        object code
                else if OPCODE = 'RESW' then
                    add 3 #[OPERAND] to LOCCTR
                else if OPCODE = 'RESB' then
                    add #[OPERAND] to LOCCTR
                else if OPCODE = 'BYTE' then
                    begin
                        find length of constant in bytes
                        add length to LOCCTR
                        convert constant to object code
                    end
                if object code will not fit into current
                    text record then
                    begin
                        write text record to object program
                        initialize new text record
                    end
                add object code to Text record
            end
        write listing line
        read next input line
    end
    write last Text record to object program
    write End record to object program
    write last listing line
end (Pass 1)
```

**Figure 2.19(c)**  *(cont'd)*

references that could not be handled by the assembler. Of course, the objec

iii) Assembler generating object code.

→ This type of one-pass assembler also works in the same manner as before (load-and-go) except where the definition of symbol is encountered

→ when a symbol definition is encountered, instructions that made forward references to that symbol may no longer be available in memory for modification. ⟹ means they have already been written out as part of a Text Record on the object program.

↳ In such situations, assembler generates another Text record with the correct operand address.

↳ when the program is loaded, this address will be inserted into the instructions by the loader.

↳ The object program for fig 2.18 is shown below during one-pass-one.

H∧COPY ∧001000∧001078

1. T∧001000∧09∧4SLFUB∧000003∧000000

T∧001000∧09∧4SLFUB∧000003∧000000∧2B1006∧300000∧H80000∧3C2012

2. T∧00200F∧15∧H1009∧H80000∧001000∧2B1006∧300000∧H80000∧3C2012

3. T∧00201C∧02∧2024

4. T∧002024∧19∧001000∧0C100F∧001003∧0C100C∧H80000∧081009∧HC0000∧

E∧001000

5. T∧002013∧02∧203D

6   T^002030^1E^0H1006^001006^6E02039^3020G3^DF2039^2S1006^
    300005^5H900F^2C205A^382063

7   T^002050^62^2056B

8   T^00205B^07^101006^H50000^0C

9   T^00201F^02^2062

10   T^002031^02^2062

11   T^002062^1E^0H1006^E02061^302065^50900F^DC2061^2C100C^
    382065^HC0000

E^00200E

→ The second text record contains the object code generated from lines 10 through 40 in fig 2.18.

↳ The operand addresses for instructions on line 15,30 and 45 has been generated as C0C0.

↳ when definition of ENDFIL on line 45 is encountered, the assembler generates the third Text Record. It indicates that the value 2039h (address of ENDFIL) has to be loaded at location 201C.

↳ This continues for all the forward references encountered.

## b) multi-pass Assemblers

↳ We know that whenever we use EQU assembly directive, any symbol used on the right-hand side should be defined previously in the source program. This is not true always.

↳ eg: ALPHA   EQU   BETA

        BETA   EQU   DELTA

        DELTA   RESW   1

↳ As we see above, we have multiple forward references ie Alpha depends on value of Beta, Beta depends on value of delta.

↳ Any assembler that makes only two sequential passes over the source program cannot resolve such a sequential sequence of definitions.

↳ To overcome this we go for multi-pass assembler which makes as many passes as needed to process the definitions of symbols.

↳ It is not necessary for multi-pass assembler to make more than two passes over the entire program.

** ↳ Instead, the portions of the program that involve forward references in symbol definitions are saved

during pass-1. Additional passes through their stored definitions are made as the assembly program. This process is followed by a normal Pass-2

↳ SYMTAB stores the symbol definition, symbols which are dependent on this, op of symbols dependent on this symbol

Ex:-

| | LOC | | | |
|---|---|---|---|---|
| 1 | | HALFSZ | EQU | MAXLEN/2 |
| 2 | | MAXLEN | EQU | BUFEND - BUFFER |
| 3 | | PREVBT | EQU | BUFFER -1 |
| 4 | 1,2h | BUFFER | RESB | 4096 ⇒ (1000)₁₆ |
| 5 | 203h | BUFEND | EQU | * — value the current location counter value |

→ below fig shows the symbol table entry when it reads line no. 1 indicating that Halfsz depends on maxlen value

| BUFEND | * | | | MAXLEN | φ |
|---|---|---|---|---|---|

| HALFS2 | 41 | MAXLEN|2 | | φ |
|---|---|---|---|---|

| PREVBT | 1033 | | φ |
|---|---|---|---|

| MAXLEN | 41 | BUFEND - BUFFER | | HALFS2 | φ |
|---|---|---|---|---|---|

| BUFFER | 1034 | | φ |
|---|---|---|---|

(d)

⟶ Fig (d) shows the symbol table entry after scanning line no. ⑪ whose location counter value is 1034. MAXLEN dependency value 42 is reduced to 41.

| BUFEND | 2034 | | φ |
|---|---|---|---|

| HALFS2 | 800 | | φ |
|---|---|---|---|

| PREVBT | 1033 | | φ |
|---|---|---|---|

| MAXLEN | 1000 | | φ |
|---|---|---|---|

| BUFFER | 1034 | | φ |
|---|---|---|---|

(c)

⟶ Fig (c) indicates the complete symbol table entry proces

| BUFEND | * |  |  |  | → MAXLEN | $\phi$ |
| HALFS2 | $\leq 1$ | MAXLEN/2 |  | $\phi$ |
| MAXLEN | $\leq 2$ | BUFEND — BUFFER |  |  | → HALFS2 | $\phi$ |
| BUFFER | * |  |  |  | → MAXLEN | $\phi$ |

(b)

⌊→ fig (b) shows symbol table entry after reading line ro 2
As we see, MAXLEN depends on 2 symbols BUFEND &
BUFFER. ∴ MAXLEN ≤ 2.



| BUFEND | * |  |  |  | → MAXLEN | $\phi$ |
| HALFS2 | $\leq 1$ | MAXLEN/2 |  | $\phi$ |
| MAXLEN | $\leq 2$ | BUFEND — BUFFER |  |  | → HALFS2 | $\phi$ |
| PREVBT | $\leq 1$ | BUFFER-1 |  | $\psi$ |
| BUFFER | * |  |  |  | → MAXLEN | → PREVBT | $\phi$ |

(c)

Thus two are dependent
on value of Buffer.

# CHAPTER 4

# Macro Processors

Chapter : Microprocessor

→ We are going to study definition of macro

→ what is the need for macro

→ Data structures used in macro invocation and expansion

Macro : It is a single instruction that expands automatically into a set of instructions to perform a particular task. Thus macro instructions allow the programmer to write a shorthand version of a program, and leave the mechanical details to be handled by the macroprocessor.

Ex: In sic/xE, 7 instructions (STA, STB, etc) is required to save the contents of all registers before calling a subprogram, but by using a macro instruction, the programmer can write a single short instruction like SAVEREGS. The SAVEREGS macro instruction would be expanded into seven instruction required to save the contents of all registers

↳ LOADREGS macro instruction would be used to reload the register contents after returning from the subprogram

* → Macro processor performs no analysis of the text it handles and is not concerned about the meaning of the involved statements during macro expansion. ∴ The design of a macro processor is machine independent.

### 4.1 Basic macro processor functions.

→ Macros refers to a set of statements which will replace every invocation to it. The two concepts associated with macros are:

(i) Macro Definition
(ii) Macro Expansion

(i) Macro Definition:

→ Consists of macro prototype, one or more module and macro preprocessor.

→ macro definition is a set of statements present inbetween a macro header statement (MACRO) and a macro end statement (MEND). MACRO and MEND are two assembly directives used in macro definition.

* **Syntax** of macro prototype:

&lt;macro name&gt; [&lt; formal parameter specification&gt;[,...]]

where
    &lt;macroname&gt; : The mnemonic field of a statement
    &lt;formal parameter&gt;: &lt;parameter name&gt;[&lt;parameter kind&gt;]
    ... each parameter begins with '&'

Syntax: Macro call

        &lt;macro name&gt; [ &lt;actual parameter specification&gt; [,...] ]

&#10230; In general, macro definition is given as

~~NAME   MACRO PARAMETERS~~

NAME   MACRO   PARAMETERS

      .
      .
      .

body; the statements which are generated as the expansion
of the macros

     .
     .
     .

MEND

// macro invocation is as

      .
      .
      .

NAME   PARAMETERS

     .
     .

&#10230; As in fig 4.1, macro definition is at line 10

10   RDBUFF   MACRO   &INDEV, &BUFADR, &RECLTH

     .
     .

95   MEND

&#10230; macro invocation in fig 4.1

190   RDBUFF   F1, BUFFER, LENGTH

1.8) Macro Expansion

→ Main invocation statements are the statements of the macro body that are expanded each time the macro is invoked.

→ The program in fig 4.1 is supplied as input to a macro processor.

→ Fig 4.2 shows the output that would be generated

→ In expanding the macro invocation on line 190, argument F1 is substituted for the parameter &INDEV, BUFFER is substituted for &BUFADR, &LENGTH is substituted for &RECLTH.

→ Line 190a through 190m show the complete expansion of the macro invocation on line 190.

→ Same in line 210a through 210h for WRBUF MACRO

→ As we see the macro body does have any labels ∵ for example, line 140 → JEQ *-3 and line 155 JLT *+11. If we put label, it would be generated twice on line 210d and 220d resulting in an error (duplicate label definition) when the program is assembled.

→ *-3, *+9, ... indicate pc-relative addressing

178

**Figure 4.1** — Use of macros in a SIC/XE program.

| Line | LABEL | OPCODE | OPERAND | |
|---|---|---|---|---|
| 5 | COPY | START | 0 | COPY FILE FROM INPUT TO OUTPUT |
| 10 | RDBUFF | MACRO | &INDEV,&BUFADR,&RECLTH | |
| 15 | . | | | MACRO TO READ RECORD INTO BUFFER |
| 20 | . | | | |
| 25 | | CLEAR | X | CLEAR LOOP COUNTER |
| 30 | | CLEAR | A | |
| 35 | | CLEAR | S | |
| 40 | | +LDT | #4096 | SET MAXIMUM RECORD LENGTH |
| 45 | | TD | =X'&INDEV' | TEST INPUT DEVICE |
| 50 | | JEQ | *-3 | LOOP UNTIL READY |
| 55 | | RD | =X'&INDEV' | READ CHARACTER INTO REG A |
| 60 | | COMPR | A,S | TEST FOR END OF RECORD |
| 65 | | JEQ | *+11 | EXIT LOOP IF EOR |
| 70 | | STCH | &BUFADR,X | STORE CHARACTER IN BUFFER |
| 75 | | TIXR | T | LOOP UNLESS MAXIMUM LENGTH |
| 80 | | JLT | *-19 | HAS BEEN REACHED |
| 85 | | STX | &RECLTH | SAVE RECORD LENGTH |
| 90 | | MEND | | |
| 95 | WRBUFF | MACRO | &OUTDEV,&BUFADR,&RECLTH | |
| 100 | . | | | MACRO TO WRITE RECORD FROM BUFFER |
| 105 | . | | | |
| 110 | | CLEAR | X | CLEAR LOOP COUNTER |
| 115 | | LDT | &RECLTH | |
| 120 | | LDCH | &BUFADR,X | GET CHARACTER FROM BUFFER |
| 125 | | TD | =X'&OUTDEV' | TEST OUTPUT DEVICE |
| 130 | | JEQ | *-3 | LOOP UNTIL READY |
| 135 | | WD | =X'&OUTDEV' | WRITE CHARACTER |
| 140 | | TIXR | T | LOOP UNTIL ALL CHARACTERS |
| 145 | | JLT | *-14 | HAVE BEEN WRITTEN |
| 150 | | MEND | | |
| 155 | . | | | MAIN PROGRAM |
| 160 | . | | | |
| 165 | FIRST | STL | RETADR | SAVE RETURN ADDRESS |
| 170 | CLOOP | RDBUFF | F1,BUFFER,LENGTH | READ RECORD INTO BUFFER |
| 175 | | LDA | LENGTH | TEST FOR END OF FILE |
| 180 | | COMP | #0 | |
| 185 | | JEQ | ENDFIL | EXIT IF EOF FOUND |
| 190 | | WRBUFF | 05,BUFFER,LENGTH | WRITE OUTPUT RECORD |
| 195 | | J | CLOOP | LOOP |
| 200 | ENDFIL | WRBUFF | 05,EOF,THREE | INSERT EOF MARKER |
| 205 | | J | @RETADR | |
| 210 | EOF | BYTE | C'EOF' | |
| 215 | THREE | WORD | 3 | |
| 220 | RETADR | RESW | 1 | |
| 225 | LENGTH | RESW | 1 | LENGTH OF RECORD |
| 230 | BUFFER | RESB | 4096 | 4096-BYTE BUFFER AREA |
| 235 | | END | FIRST | |

---

**Figure 4.2** — Program from Fig. 4.1 with macros expanded.

| Line | | Source statement | | |
|---|---|---|---|---|
| 5 | COPY | START | 0 | COPY FILE FROM INPUT TO OUTPUT |
| 180 | FIRST | STL | RETADR | SAVE RETURN ADDRESS |
| 190 | CLOOP | CLEAR | X | CLEAR LOOP COUNTER |
| 190a | | CLEAR | A | |
| 190b | | CLEAR | S | |
| 190c | | +LDT | #4096 | SET MAXIMUM RECORD LENGTH |
| 190d | | TD | =X'F1' | TEST INPUT DEVICE |
| 190e | | JEQ | *-3 | LOOP UNTIL READY |
| 190f | | RD | =X'F1' | READ CHARACTER INTO REG A |
| 190g | | COMPR | A,S | TEST FOR END OF RECORD |
| 190h | | JEQ | *+11 | EXIT LOOP IF EOR |
| 190i | | STCH | BUFFER,X | STORE CHARACTER IN BUFFER |
| 190j | | TIXR | T | LOOP UNLESS MAXIMUM LENGTH |
| 190k | | JLT | *-19 | HAS BEEN REACHED |
| 190l | | STX | LENGTH | SAVE RECORD LENGTH |
| 195 | | LDA | LENGTH | TEST FOR END OF FILE |
| 200 | | COMP | #0 | |
| 205 | | JEQ | ENDFIL | EXIT IF EOF FOUND |
| 210 | | CLEAR | X | CLEAR LOOP COUNTER |
| 210a | | LDT | LENGTH | |
| 210b | | LDCH | BUFFER,X | GET CHARACTER FROM BUFFER |
| 210c | | TD | =X'05' | TEST OUTPUT DEVICE |
| 210d | | JEQ | *-3 | LOOP UNTIL READY |
| 210e | | WD | =X'05' | WRITE CHARACTER |
| 210f | | TIXR | T | LOOP UNTIL ALL CHARACTERS |
| 210g | | JLT | *-14 | HAVE BEEN WRITTEN |
| 215 | | J | CLOOP | LOOP |
| 220 | ENDFIL | CLEAR | X | CLEAR LOOP COUNTER |
| 220a | | LDT | THREE | |
| 220b | | LDCH | EOF,X | GET CHARACTER FROM BUFFER |
| 220c | | TD | =X'05' | TEST OUTPUT DEVICE |
| 220d | | JEQ | *-3 | LOOP UNTIL READY |
| 220e | | WD | =X'05' | WRITE CHARACTER |
| 220f | | TIXR | T | LOOP UNTIL ALL CHARACTERS |
| 220g | | JLT | *-14 | HAVE BEEN WRITTEN |
| 225 | | J | @RETADR | |
| 250 | EOF | BYTE | C'EOF' | |
| 255 | THREE | WORD | 3 | |
| 260 | RETADR | RESW | 1 | |
| 265 | LENGTH | RESW | 1 | LENGTH OF RECORD |
| 270 | BUFFER | RESB | 4096 | 4096-BYTE BUFFER AREA |
| 275 | | END | FIRST | |

179

```
1   MACROS  MACRO                     {Defines SIC standard version macros}
2   RDBUFF  MACRO   &INDEV,&BUFADR,&RECLTH
3           .                         {SIC standard version}
4           MEND                      {End of RDBUFF}
    WRBUFF  MACRO   &OUTDEV,&BUFADR,&RECLTH
            .                         {SIC standard version}
5           MEND                      {End of WRBUFF}

6           MEND                      {End of MACROS}
```

(a)

```
1   MACROX  MACRO                     {Defines SIC/XE macros}
2   RDBUFF  MACRO   &INDEV,&BUFADR,&RECLTH
3           .                         {SIC/XE version}
4           MEND                      {End of RDBUFF}
    WRBUFF  MACRO   &OUTDEV,&BUFADR,&RECLTH
            .                         {SIC/XE version}
5           MEND                      {End of WRBUFF}

6           MEND                      {End of MACROX}
```

(b)

Figure 4.3   Example of the definition of macros within a macro body.

```
DEFTAB
RDBUFF   &INDEV,&BUFADR,&RECLTH
CLEAR    X
CLEAR    A
CLEAR    S
+LDT     #4096
TD       =X'?'
JEQ      *-3
RD       =X'?'
COMPR    A,S
JEQ      *+11
STCH     &BUFADR,X
TIXR     T
JLT      *-19
STX      #3
MEND
```

(a)

```
ARGTAB
1   F1
2   BUFFER
3   LENGTH
```

(b)

Figure 4.4   Contents of macro processor tables for the program in Fig. 4.1: (a) entries in NAMTAB and DEFTAB defining macro RDBUFF, (b) entries in ARGTAB for invocation of RDBUFF on line 190.

→ For designing two pass macro processor, all macro definitions are processed during pass-1 and all macro invocation statements are expanded during pass-2.

→ The two pass macro processor would not allow the body of one macro instruction to contain definition of other macro, ∵ all macros defined during the pass before any macro invocation were expanded.

→ Example of recursive macro definition is shown in Fig 4.3 (a) for sic machine and Fig 4.3 (b) for sic/xe machine

→ The same program can be run on either a sic machine or sic/xe machine. Invocation of MACROS or MACROX is only changed for use.

→ A one-pass macro processor that alternate between macro definition and macro expansion in a recursive way is able to handle recursive macro definition provided that a macro definition of a macro should appear before the invocation.

...substitution for one pass macro processor. (by hand)

→ There are three main data structures involved

(i) **Definition table (DEFTAB)**

   ↳ It stores the macro definitions which contains the macro prototype and the statements that make up the macro body

   ↳ Comment lines are omitted ∵ they are not part of the macro expansion

   ↳ References to the macro instruction parameters are converted to a positional notation for efficiency in substituting arguments.

(ii) **Name Table (NAMTAB)**

   ↳ It stores the macro names, which serves as index to DEFTAB.

   ↳ For each macro instruction defined, NAMTAB contains pointers to the beginning and end of the definition in ~~DEFTAB~~ DEFTAB.

(iii) **Argument Table (ARGTAB)**

   ↳ Used during the expansion of macro invocations.

   ↳ When a macro invocation stmt is recognized, the arguments are stored in ARGTAB according to their position in the argument list.

   ↳ When it is expanded, arguments from ARGTAB are substituted for the corresponding parameters in the macro body.

*Algorithm for a one-pass macro processor*



**ONE PASS PROGRAM**

*Procedure of GETLINE*

*Procedure:*

**GETLINE:**
If EXPANDING then
get the next line to be processed from DEFTAB
else
read next line from input-file

**EXPAND**
setup the original into with args
figure a macro invocation idol
iterators of GETLINE

**PROCESSLINE**

**PROCESSLINE:**
DEFINE
EXPAND
output source line

**DEFINE:**
make appropriate entries in DEFTAB and NAMTAB

```
begin { macro processor }
  EXPANDING := FALSE
  while OPCODE ≠ 'END' do
    begin
      GETLINE
      PROCESSLINE
    end {while}
end {macro processor}

procedure PROCESSLINE
  begin
    search NAMTAB for OPCODE
    if found then
      EXPAND
    else if OPCODE = 'MACRO' then
      DEFINE
    else write source line to expanded file
  end {PROCESSLINE}
```

Figure 4.5 Algorithm for a one pass macro processor.

```
procedure DEFINE
  begin
    enter macro name into NAMTAB
    enter macro prototype into DEFTAB
    LEVEL := 1
    while LEVEL > 0 do
      begin
        GETLINE
        if this is not a comment line then
          begin
            substitute positional notation for parameters
            enter line into DEFTAB
            if OPCODE = 'MACRO' then
              LEVEL := LEVEL + 1
            else if OPCODE = 'MEND' then
              LEVEL := LEVEL - 1
          end {if not comment}
      end {while}
    store in NAMTAB pointers to beginning and end of definition
  end {DEFINE}

procedure EXPAND
  begin
    EXPANDING := TRUE
    get first line of macro definition (prototype) from DEFTAB
    set up arguments from macro invocation in ARGTAB
    write macro invocation to expanded file as a comment
    while not end of macro definition do
      begin
        GETLINE
        PROCESSLINE
      end {while}
    EXPANDING := FALSE
  end {EXPAND}

procedure GETLINE
  begin
    if EXPANDING then
      begin
        get next line of macro definition from DEFTAB
        substitute arguments from ARGTAB for positional notation
      end {if}
    else
      read next line from input file
  end {GETLINE}
```

Figure 4.5 (cont'd)

Handling nested macro definition within macros

→ In DEFTAB (define procedure), when a macro definition is being entered into DEFTAB, the normal approach is to continue until an MEND directive is reached. This will not work for nested macro defn : the first MEND encountered in the inner macro will terminate the whole macro definition process

→ To solve this problem, a DEFINE procedure is used which maintains a counter named LEVEL. The LEVEL value is incremented by 1 when MACRO directive is read. The value is determined by 1 when MEND directive is read. When LEVEL value becomes 0, the MEND that corresponds to the original MACRO directive has been found. This process is very much like matching left and right parenthesis when scanning an arithmetic expression.

# LOADERS AND LINKERS.

. 3 processes a system program performs →

1. Loading — bringing the object program into memory f execution

2. Relocation —
modify the object program so that it can be located at a different location from the original one

3. Linking —
combing 2 or more seperate object and programs and supply information needed to allow reference b/w them.

Loaders — system program that perform loading function.
— Can also support linking & relocating.

Linker — seperate system program f linking operation.

## LOADERS 3.1

→ Basic Loader Function. or fundamental
  ∘ The most basic /loader function is →
      bringing object program into memory and starting execution.

  ∘ Absolute Loader
     It is the most basic loader that just performs loading function.
     It performs all its functions in a single pass.
      - It checks Header record to verify that correct program is being loaded and that it is will fit in the memory space available
      - It reads each Text record. and the object code is moved to its corresponding memory location
      - The End record poo indicates end of object code and gives address of location from where execution starts.

ALGORITHM :: Absolute loader

```
begin
        read Header record
        verify program name and length
        read first Text record

        while record type ≠ `E´ do
              begin
                  {if object code is in character form, convert in
                   internal representation}
                  move object code to specified location in memory
                  read next object program record.
              end
        jump to address specified in End record
end.
```

Note →.

In the object program, each byte of assembled code is given using hexadecimal representation in character form.

eg - OP code for STL → 14

It is represented using pair of character '1' & '4'

So, when loader reads this, they occupy 2 bytes of memory. But, In the instruction loaded for execution that it is to be stored as 1 byte represented by hexadecimal 14.

⟹ Each pair of byte from object prog record must be packed together into 1 byte during loading.

* This method of representation is insufficient
* So, object program can be stored in binary form → each byte of object code stored in 1 byte of memory but they aren't easy to read for humans.

- Simple Bootstrap Loader.
  - Special absolute loader that is ~~toot~~ executed when computer is first started or restarted.
  - It loads the 1st program to be run on the computer, ie OS.

```
Line
0   BOOT      START         0         BOOTSTRAP LOADER FOR SIC/XE
1      .
2   . THIS BOOTSTRAP READS OBJECT CODE FROM DEVICE F1 AND
3   . ENTERS IT INTO MEMORY LOCATION STARTING FROM
4   . ADDRESS 80h. AFTER LOADING IS COMPLETE CONTROL JUMP
5   . TO 80h IS EXECUTED TO BEGIN EXECUTION OF PROGRAM.
6   . REGISTER X CONTAINS NEXT ADDRESS TO BE LOADED
7              CLEAR     A          CLEAR REG A TO 0
8              LDX       #128       INITIALIZE REG X TO 80h
9   LOOP       JSUB      GETC       READ HEX DIGIT FROM PROG
10             RMO       A,S        SAVE IN REG S.
11             SHIFTL    S,4        MOVE TO HIGH-ORDER 4 BITS
12             JSUB      GETC       GET NEXT HEX DIGIT
13             ADDR      S,A        COMBINE DIGITS TO 1 BYTE
14             STCH      0,X        STORE AT ADDR. IN X.
15             TIXR      X,X        ADD 1 TO MEMORY ADDRESS
16             J         LOOP       LOOP TILL EOF REACHED.
17      .
18   . SUBROUTINE TO READ FROM DEVICE AND CONVERT IT
19   . FROM ASCII TO HEXA DIGIT VALUE AND RETURN IT
20   . TO REG A. IF EOF ENCOUNTERED, CONTROL TRANSFERRED
21   . TO 80h
22   .
23  GETC       TD        INPUT      TEST INPUT DEVICE
24             JEQ       GETC       LOOP UNTIL READY
25             RD        INPUT      READ CHARACTER
26             COMP      #4         IF CHAR IS 04h (EOF)
27             JEQ       80         JUMP TO START OF PROG LOADED
28             COMP      #48        COMP TO 30h ('0')
29             JLT       GETC       SKIP CHAR < '0'
30             SUB       #48        SUBTRACT 30h FROM ASCII
31             COMP      #10        FOR 'A' TO 'F', RESULT <10 THEN
               JLT       RETURN.    CONVERSION COMPLETE, ELSE.
32             SUB       #7         SUBTRACT 7 MORE.
33  RETURN     RSUB                 RETURN TO CALLER
34  INPUT      BYTE      X'F1'      INPUT DEVICE.
35             END       LOOP
```

- The bootstrap begins at address 0. [Line 0]
- It loads the OS starting at address 80h by initializing register x (the pointer) to 80h [LINE x]

- As this is the 1st prog to be loaded, its loading is simple.

The object program from device F1 is
- represented as 2 hexadecimal digit for 1byte
- has no Header or End record or any other control information.

Here, the object code is loaded into consecutive bytes of memory starting at 80h.

- Subroutine GETC →
  It reads 1 char from device F1 and converts it from ASCII to the hexadecimal digit it represented.
  When it encounters EOF, the control moves to 80h (i.e start of loaded program).

So in this program,
the main loop keeps track of the next memory location for loading and reads the 2 characters & stores it as 1byte.
The subroutine reads the character and converts it from ASCII to represented hex value.

**3.2**

→Disadvantage of absolute loader.
  o The program needs absolute memory location for loading to be specified by the programmer.
  But in large & advanced machine, multiple independent program run together & shay memory. Here predefining memory for loading is impossible.

  o The subroutines of libraries aren't used efficiently.
  For efficient use only required subroutines should be loaded but this isn't possible with absolute addresses

MACHINE DEPENDENT LOADER FEATURES

→ In most modern computer, the loaders also perform the relocation and linking function, in addition to the basic loading function.

○ Relocation
- Loaders that allow relocation are called relocating loader or relative loader.
- Methods for specifying relocation as part of object program

(i). Modification record is used to describe each part of object code that must be changed when program relocates
And the instructions whose value is affected by relocation are ones that use extended format.

The modification record specify the start address & length of fields to be altered. It then describes the modification to be performed

But, this method isn't suited for all machines
eg- In a SIC machine, there is no relative addressing & so, almost all instructions need to be modified during relocation. This leads to a lot of Modification record that dramatically increases object code size.

(ii) There is a relocation bit associated with ~~Text Record~~ each word of object code in ~~Text Record~~
It is used in machines that primarily use direct addressing & fixed instruction format
In SIC machine, each instruction occupies 1 word, i.e. one relocation bit per instruction
The relocation bits gathered together to a bit mask which is present in the Text Record following the ~~record~~ length indicator.
eg- T$_\wedge$001057, 0A$_\wedge$800$_\wedge$1000 36$_\wedge$4C0000$_\wedge$F1$_\wedge$001000

Scanned by CamScanner

If relocation bit correspondy to a word is

 — 1 → modification required
    prog's staty addr is to be
    added to this word dury
    relocation

 — 0 → no modification required.

If Text record has fewa than 12 words, then for unused words the Corresponding words relocation bit = 0.

eg - FFC ( IIII IIII IIOO)
  First 10 words need to be modified.

~~iii) Some~~

(ii) Some compute have hardware relocation capability that eliminates need of loader to relocate program.

The SIC/XE machine usually use the Modification record scheme for relocation.

ALGORITHM: SIC/XE relocation loader
begin.
 get PROGADDR from operaty system

 while not end of input do
  begin
   read next record
   while record type ≠ 'E' do
    begin
     read next input record
     while record type = 'T' then
     begin.
      more object code from a
      record to location ADDR
      + specified address
     end
     while record type = 'M'
     add PROGADDR at location
     PROGADDR+ specified address
    end.
  end
end.

The SIC machines usually use the modification bit scheme.

ALGORITHM: SIC reloaction loader algorithm.

```
begin
    get PROGADDR from operating system
    while not end of input do
        begin
            read next record
            while end ≠ record type ≠ 'E' do
            while record type = 'T'
                begin
                    get length = second data.
                    mask bits(M) as third data.

                    for (i=0, i<length, i++)
                        if Mᵢ = 1 then
                            add PROGADDR at the
                            location PROGADDR + specified
                            address.
                        else
                            move object code from record
                            to location PROGADDR +
                            specified address
                    read next record.
                end
        end
end.
```

• Program Linking

- in programs made up of multiple control sections
  can be assembled in 2 ways
    ○ all control sections together
            ie in same invocation of assembly
    ○ each independently.
  In both case they will appear as separate segments
  of object code after assembly
  Assembler sees code only as control sections
  that are to be loaded, relocated & linked. It
  doesn't need to know which control sections

were assembled at same time

Consider 3 programs each containing single control section:

- LIST A, LISTB, LISTC ──→ list of items of each prog.
  END A, ENDB, ENDC ──→ marks end of lists
  Reference to external symbol in
  REF1 to REF3 ──→ as instruction operands
  REF4 to REF8 ──→ values of data word.

### REF1
In PROGA, REF1 is a reference to label within the program so, no modification for relocation or linking needed

In PROGB & PROGC ──→ REF1 is reference to an external symbol so, assembler uses extended format instruction with address-field ──→ 0000 & Modification record required to tell loader that add value of LISTA is to be added after linking

### REF2
Similiar to REF1 but here PROGB has local reference & PROGA & PROG C have external symbol.

### REF3
It is an immediate operand whose value is ENA−LISTA In PROGA it can be directly computed but in the other 2 prog, the value is unknown

The expression is assembled as external reference & final result is an absolute value independant of location of where program is loaded.

General approach ──→
   assembler evaluate as much of the expression as it can & remaiy is passed on to loader via Modification record.

### eg- REF4
In PROGA assembler can evaluate all expression except for LISTC
The result is an initial value of 000014h and
                                    └→ (0054-0040)
  1 Modification record

In PROGB no terms can be evaluated by the assembler The result is an initial value of 000000h & 3 Modification record-

In PROGC assembler can supply value of LISTC but rest is unknown.

Initial value is relative address of LISTC and & Modification record telly to add value of ENDA & subtract value of LISTA.

— Consider the 3 progs have been loaded into memory with PROGA start at address 4000, with PROGB & PROGC immediately following

*REF4 to REF8 will end up with same value in each of the 3 program after relocation and linking.

eg – value of reference REF4 in PROGA.
Located at 4054 ( 4000 + relative address of REF4(0054))

Initial value of REF4 → 000014h (from the Text record)

To this we add address assigned to LISTC (4112) [beginning of PROGC + 30]

⇒ value in memory 4054
→ 000014 + 004112
= 004126.



Object Program

PROGA | HPROGA ... o
:
T 000054 OF [000014] ...   REF4
:
M 000054, + LISTC

PROGC | HPROGC ...
:
DLISTC [000030] ,
:

Memory

0000
:
0045 4054 ... [004126] ...   (REF4)

(+) → 4112
Actual address of LISTC

Load Address
PROGA   004000
PROGB   004063
(PROGC)   0040E2

In PROG B fr REF4

ini
located at relative address 70
so, memory location $(4063 + 70 \rightarrow 40D3)$
initial value $\rightarrow$ 000000
$$+ ENDA \rightarrow 4054 \quad (4000 + 54)$$
$$+ LISTC \rightarrow 4112 \quad (40E2 + 30)$$
$$- LISTA \rightarrow 4040 \quad (4000 + 40)$$
$$= 004126$$

$\rightarrow$ same as in PROGA

Similarly fr PROGC , REF4 also results in 004126

* REF 1 - REF3 $\rightarrow$ which are reference that are instruction operand, calculated values after loading arent always equal as additional address calculation step involved in case of base or PC relative instruction

eg - REF1 $\rightarrow$

Fr PROGA $\rightarrow$ target address 4040.
            displacement 01D + PC (4023)

Fr PROGA $\rightarrow$ REF1 is extended format instruction with direct address which is 4040
            (LISTA locatn $\rightarrow$ 4000 + 40 - 4040)

$\rightarrow$ Algorithm & Data Structure fr Linking Loader.

- Algorithm for linking & relocatg loader that uses Modification record for relocation so that linking & relocation function are performed using same mechanism.

- i/p to loader is set of object programs that are to be linked together

    Programs may contain external reference to symbol whose definition come later & so linky operation can't be performed till the external symbol is assigned ant address.

∴ Linking loader makes 2 passes over its i/p.
    Pass 1 $\rightarrow$ assigns address to all external symbols
    Pass 2 $\rightarrow$ performs actual loading, relocation &
            linking

Data structure needed,

- ESTAB → external symbol table
  - → it stores name & address of each external symbol in the set of control sections (programs) that are loaded
  - → Hashed organization is used for this table.

**PROGRAM**

Important variable needed,

- PROGADDR → program load address
- CSADDR → control section address

PROGADDR is the beginning address where linking program is to be loaded. It's value is supplied by the OS

CSADDR — starting address of control section currently being scanned by loader

## Pass1

→ loader only concerned with Header & Define record types.

→ Value for PROGADDR is obtained from OS, which is the CSADDR for the 1st control section.

→ Control section name is obtained from Header record and is entered in ESTAB with its corresponding value given by CSADDR.

→ All external symbols that appear in Define record also entered in ESTAB. Their address is relative address + CSADDR.

→ When end record reached, control section length (SDTH) added to CSADDR → this is CSADDR for next section.

**Pass 1:**

```
begin
    get PROGADDR from operating system
    set CSADDR to PROGADDR (for first control section)
    while not end of input do
        begin
            read next input record (Header record for control section)
            set CSLTH to control section length
            search ESTAB for control section name
            if found then
                set error flag (duplicate external symbol)
            else
                enter control section name into ESTAB with value CSADDR
            while record type ≠ 'E' do
                begin
                    read next input record
                    if record type = 'D' then
                        for each symbol in the record do
                            begin
                                search ESTAB for symbol name
                                if found then
                                    set error flag (duplicate external symbol)
                                else
                                    enter symbol into ESTAB with value
                                        (CSADDR + indicated address)
                            end (for)
                end (while ≠ 'E')
            add CSLTH to CSADDR (starting address for next control section)
        end (while not EOF)
end (Pass 1)
```

**Pass 2.**

- Here actual loading, relocation & linking is done
- As Each Text Record is read, object code is moved to its specified address which is,
  relative address + CSADDR.

- When Modification Record is encountered, symbol required for modification is looked up in ESTAB & its value is added or subtracted from intended location.

**Pass 2:**

```
begin
    set CSADDR to PROGADDR
    set EXECADDR to PROGADDR
    while not end of input do
        begin
            read next input record (Header record)
            set CSLTH to control section length
            while record type ≠ 'E' do
                begin
                    read next input record
                    if record type = 'T' then
                        begin
                            (if object code is in character form, convert
                                into internal representation)
                            move object code from record to location
                                (CSADDR + specified address)
                        end (if 'T')
                    else if record type = 'M' then
                        begin
                            search ESTAB for modifying symbol name
                            if found then
                                add or subtract symbol value at location
                                    (CSADDR + specified address)
                            else
                                set error flag (undefined external symbol)
                        end (if 'M')
                end (while ≠ 'E')
            if an address is specified (in End record) then
                set EXECADDR to (CSADDR + specified address)
            add CSLTH to CSADDR
        end (while not EOF)
    jump to location given by EXECADDR (to start execution of loaded program)
end (Pass 2)
```

Last step performed by loader,
transfer of control to loaded program
to Begin execution. The End record for each
control section may contain address of 1st instruction
in that control section to be executed.
If more than 1 control section specifies transfer
address, loader uses the last one encountered
If no control section specifies transfer address, loader
uses beginng of linked program (i.e PROG ADDR).

→ Algorithm can be made more efficient if we
use refrence na for external symbol in Modification
record, instead of the symbol name
Then we will need to add a Refer record that
specifies the symbol & its reference no
eg. $R_\wedge$ 02 LISTB 03 ENDB 04 LISTC 05 ENDC.
→ Refrence record in PROGA.
So the modification record will be of the form,
$M_\wedge$ 000024 $_\wedge$ 05 $_\wedge$ +02

Advantage of this method →
o avoids multiple searches of ESTAB for
same symbol while loady of control section
Now, only 1 lookup in ESTAB required for each
external refrenc symbl.

MACHINE INDEPENDENT LOADER FUNCTION

o Automatic Library Search →

    — Many linking loaders can automatically incorporate
subroutines from program libraries into the program
bely loaded.
    — Some std. libraries are used in such a way, other
libraries may be specified by control statemt o
by parameters to loader.

- Subroutines called by program being loaded are automatically fetched from the library and linked to the program while loading. This is known as
→ Automatic library call (or) library search.

## How is it done?

The linking loader that supports this must be able to keep track of external symbols used that aren't part of the input.

To do this the loader enters all external symbol it encounters into the ESTAB, ~~When it~~ if the symbol isn't already not present. When it encounters the external symbol's definition it complete its entry (if present) by filling in its address.

If at the end of Pass 1, some unresolved symbols present in ESTAB then loader searches for them in the libraries.

It is possible that subroutines featched from libraries may also contain external symbols so library search needs to be repeated till all external references have been resolved.

This process allows programmers to override the standard library's subroutines by providing our own subroutine as if input to loader so when loader goes to search library for unresolved symbol reference, the overrided subroutine reference is already defined & resolved.

## How libraries are searched?

The libraries themselves have assembled or compiled version of subroutines. It is possible to search them using their Define records, but it is inefficient.

Special structure called directory used to search libraries. It contains name of each routine & a pointer to its address within the file.

If a subroutine referred to by multiple names; there is an entry for each name and all point to same location.

- This same technique applies to resolution of external reference to data items.

○ <u>Loader Options</u>

- Loaders allow options that modify standard processing of the loader. Many loaders have a special command language that is used to specify options. Sometimes there is a separate i/p file that contains such control statement, Sometimes the statements are embedded in the primary input stream b/w object programs a can be included in the same program.

- On some systems options are specified as part of Job control language that is processed by the OS. Here, OS incorporates the options specified into a control block that is made available to loader when its invoked.

- Some options -
   ○ to select alternative sources of i/p
   eg- INCLUDE program-name (lib-name)
   This directs loader to read object program from a library & treat it as primary loader i/p's part.

- to allow users to delete external symbols or entire sections

  RDELELE csect-name

deletes control section(s) from set of progs being loaded

- to change external reference within prog being loaded & linked

  CHANGE name 1, name 2

name 1 is changed to name2 wherever it appears in the object prog.

eg –
Consider a main program say COPY that has 2 subprograms – RDREC : to read records
      WRREC : to write records

Each has its own control section

Suppose utility subroutine available such that it contains subroutines – READ & WRITE and it is more favourable for COPY to use them

As a temp measure, first we use some loader commands to make these changes without recently reassembling the program, to test the new routine

| | | |
|---|---|---|
| INCLUDE | READ (UTLIB) | } tells loader to include control section READ & WRITE from UTLIB library |
| INCLUDE | WRITE (UTLIB) | |
| DELETE | RDREC, WRREC | → tells not to load RDREC & WRREC |
| CHANGE | RDREC, READ | } → changes all external reference to RDREC to refer to READ & reference to WRREC to refer to WRITE. |
| CHANGE | WRREC, WRITE | |

- LIBRARY MYLIB
   → it automatically includes library routines to satisfy external reference

- NOCALL SYMBOLS
   → tells loader not that these external references are to remain unresolved.

- option to specify that no external reference is to be resolved

    Usefull when programs are to be linked but not immediately executed

- option to specify where execution should begin

- Option to control whether or not loader should execute program if error is detected during load

# LOADER DESIGN OPTIONS

Organisation of loading function

    — linking & relocation take place at load time
      (used by linking loader)

    — linkage editors — linking is performed prior to load time

    — dynamic linking — linking is performed at execution time

→ Linkage Editors



a) Linkage editor

In linkage editor, the source program is first assembled or compiled

. Linkage Editor vs Linking Loader.

→ Linking Loader performs all linking & relocation function
   & loads linked program directly into memory for
   execution

- Linkage editor produces a linked version of program
  called load module or executable image, which is written
  into a file or library for later execution.

→ Linkage too editor is useful for programs that
  need to be executed multiple times without reassembling
  everytime.
  For execution, relocation loads loads program into
  memory. Only the object code modification required is
  getting the actual address for loading, rest is done
  during linking. So, now loading can be done in 1 Pass.

→ Linking Loader is better when program needs to
  be reassembled for nearly every execution.

• The linked program produced by linkage editor is in
  a form that is suitable for processing by relocating loader.
  - All external references are resolved
  - relocation is indicated by some mechanism like Modification record
    & bit mask.
  Information about external references are often retained in the
  linked program as it allows subsequent relinking of program to
  replace control sections, modify external references, etc.

• If actual address for loading is known, then linkage editor can
  perform the relocation, i.e result is linked program that is
  exact image of way program will appear in memory.
  But,
      flexibility of loading program at any location is preferred
  over the reduction of overhead for performing relocation at run time
• Other useful functions ⟶
      → modification of a linked program without having to
        process the entire program.
        eg- Consider a program PLANNER that has multiple

subroutines. One of its subroine PROJECT had to be changed due to error a to improve efficieny. After new vesion of PROJECT is assembled or compiled, linkage editor can replace this subroutine in the linked version of PLANNER. usy some linkage editor commands.

```
INCLUDE   PLANNER (PROGLIB)
DELETE    PROJECT
INCLUDE   PROJECT (NEWLIB)
REPLACE   PLANNER (PROGLIB)
```

→ linkage editor can be used to build packages of subroutines or other control sections that are generally used together.

This is useful while dealy with subroutin libraries that support high level programmy lang.

eg- In a typical implementation of FORTAN, there are large number of subroutines that are used to handle formatted input & output. There are large no. of cross-reference b/w these subprograms because they are closely related

But, it is desirable to keep them as seperate modules for program modularity & maintability.

But, same set of cross-reference will be processed for almost every FORTAN program linked. This represents a substantial overhead.

We can use the linkage editor to combine the subroutines into a package using commands like,

```
INCLUDE   READR (FTNLIB)
INCLUDE   WRITER (FTNLIB)
INCLUDE   ENCODE (FTNLIB)
.
.
SAVE      FTNIO (SUBLIB)
```

The linked module FTNIO can be ndexed in directory of SUBLIB under same name as original subroutines. Thus, search of SUBLIB before FTNLIB would retreive FTNIO instead of seperate routines.

And as FTNIO would already have all cross-reference b/w subroutines resolved, these linkage wouldn't need to be reprocessed when user's program is linked.

→ linkage editor allow user to specify that external references are not to be resolved by automatic library search,

in FTNIO

eg - If 100 FORTAN program using I/o routines are to be stored in a library, the library will store 100copies of FTNIO if all external reference were resolved
This wastes a lot of library space

We can use commands to specify that no library search is to be performed during linkage editing and so they can only be resolved during execution.
This will require slightly more overhead due to 2 linkage operation but it results in large saving of library space.

a Linkage editors are in general more flexible than linking loader & also offer more control.
But they also are more complex and have greater overhead.

→ Dynamic Linking. (or dynamic loading or load on call)
   ∘ Here the linking is performed during execution time.
   i.e a subroutine is loaded & linked to rest of program when it is first called.
   ∘ It is used to allow several executing programs to share 1 copy of a subroutine or library.
   eg - run-time support routines for a high level lang like C could be stored in a dynamic link library.
      A single copy of the routines could be loaded into memory & all executing C programs could be linked to this copy instead of having separate copy for each.

- In object-oriented system, dynamic linking is used for references to software objects.

This allows implementation of object & its methods to be determined during run-time. The implementation can be changed anytime without affecting program that uses the object.

- Advantage of dynamic linking —
  * it provides ability to load routine only when they are needed.
  This results in saving of time & memory space.

  eg — Consider program contains subroutines that correct or diagnose error in i/p data during execution. If no error occurs (which can be common) then these subroutines will not be used and so will not be loaded & linked.

  - If program has many subroutines but uses only a few depending on its input, then only the subroutines required can be loaded & linked during execution.

- How to accomplish loading & linking of called subroutine?
  - The routine that must be dynamically loaded must be called via OS service request, ie the request is to the part of the loader that is kept in memory.
  So a JSUB instruction ~~that refers to an external~~ symbol
  - So instead of executing a JSUB instruction that refers to an ~~ereit~~ external symbol, program makes a load & call request to OS with symbolic name of routine as the parameter.



| Load-and-call ERRHANDL | Dynamic loader (part of OS) |
| | User program |

Here, the user program sends a load-and-call request for ERRHANDL subroutine.

- The OS examines internal table to determine whether or not routine is already loaded

If not, routine is loaded from specified user of or system library
[Load]

and then control is passed to the routine being called,
[Call]

When subroutine completes its processing, it returns to its caller (i.e OS routine that handles load-and-call request).

The OS then returns the control to the user programs

After subroutine is completed, the memory that was allocated for loading may be released & used for other purpose. But, this isn't done immediately as if a 2nd call to it occurs, another load operation won't be required. So, it is desirable to keep the subroutine till memory isn't required by an.

If subroutine called is still in memory, control is directly passed to it from the dynamic loader.

- In dynamic loading, binding of symbolic name to actual address is delayed from load time until execution time which results in greater flexibility

• But, this also requires more overhead as OS intervenes in the calling process

→ **Bootstrap Loaders**

In a idle computer with no program in memory, how do thing start?

— When machine is empty and idle there is no need for relocation, only absolute address for program being $1^{st}$ loaded is needed. (this program is usually the OS). For this we need an absolute loader loaded.

— Early computers required operatn to enter into memory the object code of absolute loader using switches on computer console. But, this is too inconvenient & error-prone.

— In some computer, absolute loader program is permanantly present in a ROM. When some hardware signal occurs indicating start up of the system, the machine begins executing this ROM program.
   In some computer, program is executed in the ROM on others, program is copied to to main memory & executed. But, it is inconvenient to change the ROM program if modification necessary.

— Intermediate solution,
      have a built-in hardware function (or small ROM program) that reads fixed length records from some device into memory at fixed location

   After reading operation is complete, control is transferred to address in memory where records is stored. These records contain address in machine instructions that absolute loader loads the absolute program that follows.

   If the instructions can't be fit in 1 record, then record causes reading of other records & they in turn cause reading of more records.
                           → hence the term bootstrap.

   $1^{st}$ record(s) ⟶ bootstrap loader.
   This loader added to begining of all object programs that

are to be loaded into empty & idle system.

## IMPLEMENTATION EXAMPLE ⟶

→ MS-DOS Linker for Pentium & other x86 system.

• Most MS-DOS compiler & assembler produce object modules, not executable machine language programs.

○ These object modules have extension .OBJ and they contain binary image of translated instructions & data of program. It also describes structure of program.

• MS-DOS LINK — linkage editor that combines one or more object modules to produce a complete executable program.

The executable program have extension .EXE. LINK can also combine the translated program with other module from object code libraries.

• A typical MS-DOS object module,

| Record Type | Description |
|---|---|
| THEADR | Translator Header |
| TYPDEF<br>PUBDEF<br>EXTDEF | External symbol & references |
| LNAMES<br>SEGDEF<br>GRPDEF | Segment definition and grouping |
| LEDATA<br>LIDATA | Translated instructions & data. |
| FIXUPP | Relocation & linking information |
| MODEND | End of object module. |

similar to Header & End record of SIC/XE {

• THEADR record — specifies name of object module

MODEND record — marks end of module & contains reference to entry point of program

* PUBDEF record — contains list of external symbol called public names that are defined in the object module.

* EXTDEF record — contains list of external symbols that are reffered to in the object module.

Similar to Define & Refer record of SIC/XE
Both PUBDEF & EXTDEF contain info abt data type designated by an external name.

* TYPEDEF record — defines the types

* SEGDEF record → describes segment in object module includy their name, length & alignment

GRPDEF record — specify how these segments are combined into groups

LNAMES record — contains list of all segment & class names used in program.

SEGDEF & GRPDEF refer to segment by givy the position of its name in the LNAMES records.

~~* LE DATA~~
* LEDATA record — contains translated informations & date from source program
It is similar to Text record of SIC/XE

LIDATA record — specify translated instructions & date that occur in repeating pattern.

* FIXUPP record — used to resolve external references & carry out address modifications that are associated with relocation & groupy of segment within the program.
It's similar to Modification record of SIC/XE
But FIXUPP records are more complicated.
A FIXUPP record must immediately follow the LEDATA a LIDATA record to which it applies.

• LINK performs its functions in two Passes.

Pass1 — computes starting address for each segment in the program

It constructs a symbol table that associates an address with each segment (us'y LNAMES, SEGMEN SEGDEF & GRPDEF records) and each external symbol (us'y EXTDEF & PUBDEF records).

If unresolved external symbols remain after all object modules are processed, LINK searches the specified libraries.

Pass 2 — LINK extracts translated instruction & data from object modules & build an image of executable program in memory.

This is because,
executable program is organised by segment & not by order of object modules.
Building a memory image, most efficient way to handle rearrangements caused by combining & concatenating segment.

If enough memory isn't available, LINK uses temp disk file in addition.

Here LINK process each LEDATA & LIDATA record along with corresponding FIXUPP records & places binary data from LEDATA & LIDATA record into memory image at locations reflecting segment address computed duiy Pass1.

Relocation & resolving of external reference is done here. A table of segment fixups is maintained that is used to perform relocation that reflects actual segment address when program is executed.

Once memory image is complete LINK writes it to .EXE file, which also contains a header that contains table of segment fixups & information about memory requirement & entry points & also initial contents of CS & SP registers

→ Sunos Linkers for SPARC system.

• Sunos provides 2 different linkers
       — run-time linker
       — link-editer.

• Link-editer is most commonly used in process of compiling a program.
 It takes 1 or more object module produced by assemblers & compilers & combines them to produce a single o/p module.

   ○ Types of output module →
     1. Relocatable object module.
       It is suitable for further link-edity
     2. Static executable
       It has all symbolic references bound & ready to run
     3. Dynamic executable
       It has some symbolic references that are to be bound at run-time
     4. Shared Object
       It provides services that can be bound at run time to 1 or more dynamic executables.

• Object module contain multiple sections which represent instructions & data areas from source program.
 These sections have a set of attributes such as "executable", "writable".
 Object modules also include list of relocation & linky operations that need to be performed & a symbol table that describes the symbols used.

• Sun-OS link-editer reads the object modules that are given to it to process. Sections that have same attributes are concatenated to form new section in o/p file.

- Symbol table from i/p files are processed to match symbol definations & references, and relocation & linking operations are performed within o/p file.

  Linking generates new symbol table & new set of relocation instruction in output file. They represent symbols that need to be bound at run-time & relocations that need to be performed during loading.

- Relocation & linking operation are specified using set of processor-specific code

  The codes reflect instruction format & addressing modes that are found in the machine as they describe the size of the field to be modified & calculations that need to be performed.

- Symbolic references from i/p file that aren't resolved are processed by referring to <u>archives</u> & shared objects

  $\downarrow$

  collection of relocatable object modules.

  Directory within archives associate symbol name with object module that contains its definition. & selected module from archive is included to resolve the references

  ~~Share~~

  Shared object is an indivisible unit that was generated by link-edit operation.
  some previous

  If reference symbol is defined in a shared object, entire content of shared object becomes logical part of o/p file.

  Link-editor records dependency to shared object, actual inclusion of the shared object happens at run time

- SunOS run-time linker,
  used to build dynamic executable & shared objects at execution time
  It determines what shared objects are required by dynamic executable & ensures that they are included.
  It also resolves any additional dependencies on other shared objects.

After locating & including necessary objects, linker performs relocation & linking to prepare program for execution.

They bind symbol to actual memory address to which segment is loaded &. Then control is passed to executable program after binding data reference.

Binding of procedure call is done during execution. Doing link-edit, calls to globally defined procedure is converted to reference to a procedure linkage table. When procedure is called for the 1st time, control is passed to own-time linker via the table. The linker looks up the actual address of the procedure & includes it to linkage table.

So, subsequent call will directly go to called

     procedure ⟶ lazy binding.

○ Run-time linker provides ≈ flexibility.

     During execution, prog can dynamically bind to new shared objects, by this allows prog to choose b/w no. of shared objects.

If a shared object isn't needed it isn't binded.

→ **Cray MPP Linker**      for Cray T3E system.

○ T3E system contain large no. of processing elements (PEs).

Each PE has its own local memory & can access memory of all other PEs.

○ An application program on a T3E system is allocated a partition that consist of several PEs. (to take advantage of ||el architecture of machine)

- Work to be done is divided b/w the PEs

  eg - partition contain consists of 16 PEs, 2 elements of a 1D array is distributed

PE0                    PE1                          PE15

```
┌──────┐         ┌──────┐                    ┌──────┐
│ A[1] │         │ A[17]│                    │ A[241]│
├──────┤         ├──────┤                    ├──────┤
│ A[2] │         │ A[18]│                    │ A[242]│
├──────┤         ├──────┤      . . . .       ├──────┤
│  ⋮   │         │  ⋮   │                    │      │
│  ⋮   │         │  ⋮   │                    │      │
├──────┤         ├──────┤                    ├──────┤
│ A[16]│         │ A[32]│                    │ A[256]│
└──────┘         └──────┘                    └──────┘
```

If prog contains loop that process all 256 elements, PE0 can execute loop fr A[1] to A[16] PE1 can execute loop fr A[17] to A[32] & so on.

- Shared data → data that is divided among no. of PEs.

  Private data → data that isn't shared by dividing it, each PE contains a copy of the data.

  Or PE has private data that exists only in its own local memory

- When program is loaded,

  each PE gets a copy of executable code, its private data & its portion of shared data.

- MPP linker organizes blocks of code or data from object program into lists.

  The blocks on a given list all share some same property.

  The blocks on each list is collected, address is assigned to each block & relocation and linking operations are performed.

  The linker then writes a executable file that contains relocated & linked blocks. It also specifies no. of PEs required & other control information.

- Distribution of shared data depends on no. of PEs.

  If no. of PEs is specified at compile time, it can't be overridden later.

  If no, either

  — linker can create executable file that targets for a fixed no. of PEs

  or

  partition size can be chosen at run time. This is called plastic executable

Plastic executable is often larger than one targeted for fixed no. of PEs as,

  it must contain copy of all relocatable object module & all linker directives that are needed to produce final executable.

# Compiler Design - 10CS63

## UNIT 1 : Introduction

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

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

Translator

Source prog → | Compiler | → Target prog

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

Source prog → | Interpretter | → output
Input →

Language Pre-processing system :

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

o A Hybrid compiler :

Src prog

↓

[ Translator ]

↓ (bytecodes)

Intermediate prog ──→ [ Virtual Machine ] ──→ output
input ──→

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

Structure of a compiler :

↓ character stream

[ Lexical Analyzer ]

↓ token stream

[ Syntax Analyzer ]

↓ Syntax tree

[ Semantic Analyzer ]

↓ Syntax tree

[ Intermediate code Generator ]

↓ Intermediate representation

[ Machine Independent code optimizer ]

↓ intermediate representation

[ Code Generator ]

↓ target machine code

[ Machine Dependent code optimizer ]

↓ target machine code

Front End Pass

Back End Pass

[ Symbol table ]

[ Error Table ]

→ 2 main parts:

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

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

→ 7 phases:

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

$$< token\text{-}name, attribute\ value >$$

abstract symbol
used in parser ←         └→ points to entry in the
                              symbol table for this token

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

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

```
      +
     / \
    2   *
       / \
      3   5
```

interior node: operators
exterior node: arguments

(3) Semantic Analysis

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

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

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

- coercions - Lang specification may permit some type conversion

(4) Intermediate Code Generation (ICG)

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

- properties :
  → Should be easy to produce
  → Should be easy to translate into target M/c.

(5) Code Optimization (M/c independent)

- to improve intermediate code to get better target code
→ Better in terms of : faster, shorter, less power consuming code

- Instead of using int to float operation,
  replace integer by its floating-point value directly

(6) Code Generation

- Input from intermediate representation maps to target lang.

- If target lang is M/c code - the instructions are translated into sequences of M/c instruction to perform same task.

- Judicious assignment of registers to hold variables is done

→ Compiler construction Tools :

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

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

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

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

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

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


Application of compiler Technology :

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

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

  → Memory Hierarchy
  Building very Large or Fast storage, but not both

③ Design New computer Architectures
→ RISC - reduces complex memory addressing, support data structure Access, procedure invocation ...

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

④ Program Translations
(i) Binary Translation - Increases S/w availability
(ii) Hardware Synthesis - Verilog, VHDL - reduces time & effort
(iii) Database Query Interpretter - SQL queries effective retrieval
(iv) compiled simulation - model run, to validate design.
(v) Reduce redundancy in code

⑤ Software Productivity Tools
(i) Type checking - to catch program inconsistency
(ii) Bounds checking - Lang. provides range checking like for the buffer overflow, security, optimize range check, sophisticated analyses, error detection tools.
(iii) Memory Management Tools - (Garbage collection)
   • Automatic memory management tracks all memory related errors - leaks...

1. Write the difference between Compiler and Interpreter

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

10.

Source program → [Compiler] → Target program

Input → [Target program] → output

10.

Source pgm
Input → [Interpreter] → output

→ Examples showing detail phases of compiler :

① position = initial + rate * 60

↓

Lexical Analysis

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

↓

Syntax Analysis

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

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

Syntax
Tree

↓

Semantic Analysis

```
        =
      /   \
  id,1     +
          / \
      id,2   *
            / \
        id,3   int to float
                |
                60
```

↓

Intermediate code generation

$t1 = \text{int to float} (60)$
$t2 = id_3 * t_1$
$t3 = id_2 + t_2$
$id_1 = t_3$

↓

M/c independent code optimization

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

⟶ Code generation

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

② a[index] = 4 + 2 + index

↓

Lexical Analysis

<id,1> <[> <id,2> <]> <=> <4> <+> <2> <+> <id,2>

↓

Syntax Analysis

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

```
          =
        /   \
      [ ]      +
     /  \    /   \
    a  index 4    +
                 / \
                2   index
```

↓

Semantic Analysis

```
          =
        /   \
      [ ]      +
     /  \    /   \
    a  index 4    +
                 / \
                2   index
```

↓

Intermediate code Generation

$t_1 = 4 + 2$

a[index] = $t_1$ + index

↓

M/C Independent code optimization

a[index] = 6 + index

↓

Code Generation

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

→ Environments and States:

environment                    state

names            locations              values
                 (variables)

○ Environment is mapping from names to locations in the store
● State is mapping from locations in store to their values

Dynamic Mapping Exceptions:

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

Ex:     int i;                      // global i          (global)
        void fun (...) {
            int i;                  // local i
        }                                          Data
                                                   Segment

(ii) Static Binding of Locations to values - declared constants

Ex:    # define ARRAYSIZE 1000        // static bind

Static scope and Block Structure:

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

→

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

② main ()
{
　int w,x,y,z;
　int i=4 ; int j=5;
```
{
    int j=7; i=6;
    w = i+j;
    printf (w);
}
```
　　　　　　　　　　　　　　　　　　→ 6+7 = 13

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

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

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

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

o/p
3
2

(4)
```
int  w, x, y, z ;
int  i = 3 ;
int  j = 4 ;
{
    int  i = 5 ;
    w = i + j ;              5 + 4 = 9
}

x = i + j ;                  3 + 4 = 7
{
    int  j = 6 ;
    i = 7 ;
    y = i + j ;              7 + 6 = 13
}

z = i + j ;                  7 + 4 = 11
```

10. what is printed by the following C cod.

a) #define a (x+1)

```
int x=2;
void b() { x=0; printf("%d\n", 1); }  → 3
void c() { int a=1; printf("%d\n", 0); }  → 2
void main() { b(); c(); }
```

b) #define a (x+1)

```
int x=2;
void b() { x=0; printf("%d\n", x); }  → 3
void c() { printf("%d\n", a); }  → 4      ∵ redesignment for this
void main() { b(); c(); }                    some variable x
                                                  x=3, a=3+1=4
```

c) #define a (x+1)

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

d) #define a (x+1)

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

→ Parameter Passing Mechanisms :

(1) Actual parameters - parameters used in call of procedure
(2) Formal parameters - parameters used in procedure definition

① Call by value

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

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

② Call by reference

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

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

• changes to formal parameter ⇒ Appear as changes in actual param

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

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

③ Call by name

• used in early prog - Algol 60.

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

→ Examples :

① call by value

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

② call by reference

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

③ call by Name - Aliasing

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

• Aliasing :

→ Interesting consequence of call by reference parameter passing where references to objects are passed by value.

• It is possible that two formal parameters refer to the same location — such variables are ALIAS to one another.

• Though they may be _distinct_ formal parameters, they may be Alias of one other.

Ex: Let a be array in procedure p

```
p
{
    || q(x,y)  call
    q(a,a) ;
}
```

array names are references to location ⟹ Alias

$x[l] = y[l]$

Questions                 Chapter -1 Introduction

1. Define Compilers?
2. Differentiate b/w compilers & Interpreter?
3. Explain The long processor System?
4. Describe the analysis-Synthesis model of the compiler or Explain in detail the Various phases of Compiler with an example?
5. Explain in detail the Various phases of Compilation for The i/p string
   a. $P = 1 + n * 60$          c. $a = (b+c) * (b+c) * 2$
   b. $x = a * b + a * b$       d. $a[index] = 4 + 2 + index$
6. why is it necessary to group phases of Compiler
7. what Is the purpose of Compiler Const$^n$ tool. Describe The different Compiler Construction tool we used?
8. Analyse the s/w productivity toal and explain
9. Explain The different parameter passing technique with an example?

# Chapter-3 - Lexical Analysis

→ Lexical Analysis :

Interaction between Lexer and Parser :



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

→ Lexer - cascade of 2 processes :
1. Scanning consists of simple processes that do not require tokenization of input, such as deletion of comments & compaction of consecutive whitespace characters into one.

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

Lexer versus Parser : Seperate phases because :
1. simplicity of design - important consideration
   (compiler)
2. compiler efficiency improved - use specialized technique for lexical Analysis (Input Buffering)
3. compiler portability is enhanced.

→ Tokens, Patterns, Lexemes :

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

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

The token names are the input symbols that parser processes.

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

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

③ Lexeme : Sequence of characters in source program that Matches the pattern of a token and is identified by the lexer as an instance of that token.
Ex: token name : Keyword
pattern : $[i][f]$
lexeme : if

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

→ Lexical Errors : Recovery options

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

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

Examples :

fi ( a < b)  ⟹  if (a < b)
int a, ;  ⟹  int a;   or  int a, b;

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

① Single buffer / 1-Buffer Technique

We use only one single buffer to store processed character from large no. of characters from source prog.

Main overhead is that if,

| lexeme size > Buffer size |

we **loose** the lexeme

· It reloads data, removes old data.

World

5 Bytes

$$\boxed{W \mid o \mid r \mid l}$$  d

4 Bytes

② 2-Buffer Technique ⟨ without sentinel / with sentinel

We use two buffers that are alternately reloaded,
Each buffer of same size N, N = Size of a <u>disk block</u>. (4096 Byte)

· Using read system call, N characters are read.

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

→ special char <u>eof</u> marks end of src file and this char is different from any other char of src prog.

→ Two pointers maintained :

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

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

- We need 2 checks in 2 Buffer without sentinels :

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

2) Before returning token check whether valid or not.

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

Buffer 1                          Buffer 2

| | | |E| |=| |M|*|eof|c|*|*|2|eof| | | |eof|

←── 4096 Bytes ──→    ↑         ↑
              lexemebegin    forward

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

→ <u>Look Ahead code with sentinel :</u>

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

              else if ( forward is at end of Buffer 2) {
                    reload Buffer 1 ;
                    forward = Beginning of Buffer 1 ; }
```

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

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

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

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

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

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

⑥ Proper Suffix - suffix which is not $\epsilon$ or equal to string itself
   Ex: anana, na

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

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

$\rightarrow$ Operations on Languages:

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

$\rightarrow$ Regular Definition:

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

$d_1 \rightarrow r_1$

$d_2 \rightarrow r_2$

$\vdots$

$d_n \rightarrow r_n$     where

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

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

Ex ① C identifiers:

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

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

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

Ex ② unsigned numbers:

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

digits $\rightarrow$ digit digit*

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

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

number $\rightarrow$ digits optional fraction optional exponent

# Algebraic Laws for Regular Expressions

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

→ Recognition of Tokens:

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

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

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

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

```
Code:

    state = 0 ;

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

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

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

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

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

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

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

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

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

② for identifier

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



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

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

    case 2 :  retract ();
              InstallID ();
              return ( retToken );
              break;
       }
   }
```

③ unsigned number

digit → [0-9]



④ Keywords

Ex:     Keyword → IF | THEN | ELSE



⑤ Delimiter / whitespace

delim → space | tab | newline

chapter-2    Lexical Analysis

Questions

1. Explain Lexical Analysis in detail with block diagram

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

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

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

5. Why 2-buffer technique is used in LA? write the algorithm for lookahead code with sentinel.

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

7. list the algebraic laws for Regular Expression.

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

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

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

Topics
1) Introductions
2) Context-free Grammars
3) Writing a Grammar
4) Top down parsing
5) Bottom up parsing

1) Introduction :

i) The Role of the parser/Block diagram for Syntax Analysis.



Block diagram:
Source program → Lexical Analysis → (token) → parser → (parse true) → Semantic Analysis → Front End
parser → (get next token) → Lexical Analysis
Lexical Analysis ↔ Symbol table ↔ parser

## The General types of parses for grammars



## Syntax - Error Handling

## Common programming Errors

i) Lexical Errors:

There Include Misspellings of Identifiers, keywords or operators

eg: Use of an Identifier elipseSize Instead of ellipseSize

ii) Syntactic Errors:

These errors Include misplaced Semicolon/Extora or missing braces;

iii) Semantic Errors

These Include type mismatches b/n operators and operands.
An example: return statement in a Java method with
result type Void

iv) Logical Errors:
Can be Anything from Incorrect reasoning on the part of
the programmer
eg: Using '=' instead of '==' in C programming.

## Error Recovery Techniques:

i) panic Mode Recovery
ii) phrase level Recovery
iii) Error productions
iv) Global Corrections

i) panic Mode Recovery:
→ In the panic Mode Recovery, keep deleting one
character at a time untill we find Synchronization
tokens ((;) and (}))

* Synchronization tokens → Semicolon (;)
→ Epilog (})

eg: int a,; // Error

ii) phrase level Recovery:
→ It Include Insert, delete, update
→ On discovering an error, a parser may perform local
Correction on the remaining i/p : that is, It may
replace a prefix of the remaining i/p by some
string that allows the parser to Continue

→ It Includes → Replacing a Comma, by a Semicolon
  ↘ delete an extra Semicolon
  → Inserting a missing Semicolon

eg: int a, ;
    Replace , by ; and delete extra ;

iii) Error productions:
 → By anticipating Common errors that might be encountered,
   we can augment the grammar for the language at
   hand with productions that generate Incorrect Constructs
 → A parser Constructed from a grammar augmented by
   these error productions detects the anticipated errors
   when an error productions is used during parsing

iv) Global Corrections:
 → Ideally, we would like a Compiler to make as few
   changes as possible in processing an incorrect input
   string. There are algorithm for choosing a minimal
   Sequence of changes to obtain a globally least Cost
   Correction.
 → Given an Incorrect input String x and Grammar G,
   these algorithms will find a parse tree for a related
   string y, such that the number of Insertions, deletions,
   and changes of tokens required to transform x into
   y is as small as possible.

drawback of Global Corrections:
 → These generally too Costly to Implement in terms of
   time and space, So these are currently only a theoretical
   Interest.

# CONTEXT FREE GRAMMARS

defn: Content free grammar is a 4-tuple defined as
(V, T, P, S), where

V: Set of Variable

T: set of Terminals

P: Set of production

S is the start symbols

Differentiate b/n CFG and RE

| CFG | RE |
|---|---|
| 1. It is the part of the Syntax Analysis | 1. It is the part of the lexical Analysis |
| 2. Usefull for describing nested gramatical structure Such as balanced paranthesis and so, on. | 2. Usefull for describing the structure of construct / lexical construct Such as Identifiers, keywords etc |
| 3. CFG's are Combined using pushdown automata | 3. Regular Expressions are Combined using Finite Automata |
| 4. CFG Can keep track of no. of Symbols Seen so far | 4. RE Cannot keep track of no. of symbols Seen so far |
| 5. Every CFG need not be RE | 5. Every RE Is a CFG |
| 6. CFG are more powerfull | 6. RE are less powerfull as Compound to CFG |
| 7. Eg: letter → [A-Z a-z] digit → [0-9] id → letter (letter/digit) | 7. Eg: [a-z A-Z 0-9][0-9]* |

Q) For the following CFG
a. Give the LMD for the string
b. Give the RMD for the string
c. Give the parse tree for the string
d. Is the grammar ambigous / Unambigous? Justify

1. $S \rightarrow SS+ | SS* | a \Rightarrow aa+a*$

2. $S \rightarrow 0S1 | 01 \Rightarrow 000111$

3. $S \rightarrow +SS | *SS | a \Rightarrow +*aaa$

4. $S \rightarrow S(S)S | \epsilon \Rightarrow (()())$

5. $S \rightarrow S+S | SS | (S) | S* | a \Rightarrow (a+a)*a$

6. $S \rightarrow (L) | a \qquad L \rightarrow L,S | S \Rightarrow (a,a)$

7. $S \rightarrow aSbS | bSaS | \epsilon \Rightarrow aabbab$

8. $bexpr \rightarrow bexpr \text{ or } bterm | bterm$
   $bterm \rightarrow bterm \text{ and } bfactor | bfactor$
   $bfactor \rightarrow \text{not } bfactor | (bexpr) | true | false$
   $\Rightarrow \text{not (true | false)}$

9. $E \rightarrow E+E | E*E | -E | (E) | id \Rightarrow id+id+id$

10. $S \rightarrow iEtS | iEtSeS | a \Rightarrow If E_1 \text{ then If } E_2 \text{ then } S_1 \text{ else } S_2$

11. $R \rightarrow R'|'R | RR | R* | (R) | a | b | c \Rightarrow a|b*c$

1) $S \rightarrow SS+ \mid SS* \mid a \Rightarrow aa+a*$

$S \xrightarrow{lm} SS*$          $S \xrightarrow{rm} SS*$

$\Rightarrow SS+S*$          $\Rightarrow Sa*$

$\Rightarrow aS+S*$          $\Rightarrow SS+a*$

$\Rightarrow aa+S*$          $\Rightarrow Sa+a*$

$\Rightarrow aa+a*$          $\Rightarrow aa+a*$

parse tree:



The Grammar Is Unambigous
because It has only one LMD and one RMD

2) $S \rightarrow 0S1 \mid 01 \Rightarrow 000111$

$S \xrightarrow{lm} 0S1$          $S \xrightarrow{rm} 0S1$

$\Rightarrow 00S11$          $\Rightarrow 00S11$

$\Rightarrow 00S111$          $\Rightarrow 000111$

parse tree:



The Grammar Is Unambigous
because It has only one LMD and only one RMD

3)

3) $S \rightarrow +SS | *SS | a \longrightarrow +*aaa$

$S \xrightarrow{lm} +SS$

$\xrightarrow{lm} +*SSS$

$\Rightarrow +*aSS$

$\Rightarrow +*aaS$

$\Rightarrow +*aaa$

$S \xrightarrow{rm} +SS$

$\xrightarrow{rm} +*SSS$

$\Rightarrow +*SSa$

$\Rightarrow +*Saa$

$\Rightarrow +*aaa$



The Grammar Is Unambigous
because It has only one LMD and one RMD

4) $S \rightarrow S(S)S | \varepsilon \Rightarrow (()())$

$S \xrightarrow{lm} S(S)S$

$\Rightarrow (S)S$

$\Rightarrow (S(S)S)S$

$\Rightarrow (()S(S)S)S$

$\Rightarrow (()S(S)S)S$

$\Rightarrow (()()S)S$

$\Rightarrow (()())S$

$\Rightarrow (()())$

(i)

$S \xrightarrow{rm} S(S)S$

$\xrightarrow{lm} (S)S$

$\xrightarrow{lm} (S(S)S)S$

$\xrightarrow{lm} (S(S)S(S)S)S$

$\Rightarrow ((S)S(S)S)S$

$\Rightarrow (()S(S)S)S$

$\Rightarrow (()(S)S)S$

$\Rightarrow (()()S)S$

$\Rightarrow (()())S$

$\Rightarrow (()())$  (ii)

RMD:

$S \xrightarrow{rm} S(S)S$

$\implies S(S)$

$\implies S(S(S)S)$

$\implies S(S(S)S(S))$

$\implies S(S(S)S())$

$\implies S(S(S)())$

$\implies S(S()())$

$\implies S(()())$

$\implies (()())$

(i)

$S \xrightarrow{rm} S(S)\underline{S}$

$\implies S(S)$

$\implies S(S(S)S)$

$\implies S(S(S))$

$\implies S(\underline{S}())$

$\implies S(S(S)S())$

$\implies S(S(S)())$

$\implies S(S()())$

$\implies S(()())$

$\implies (()())$

(ii)

parx true for LMD (i) and (ii)



The Grammar Is ambigous
Since It has 2 LMD and 2 RMD

5) $S \rightarrow S+S | SS | (S) | S*|a \implies (a+a)*a$

$S \xrightarrow{lm} SS$　　　　　　　　$S \xrightarrow{rm} SS$

$\implies S*S$　　　　　　　　　$\implies S*S$

$\implies (S)*S$　　　　　　　　$\implies S*a$

$\implies (S+S)*S$　　　　　　$\implies (S)*a$

$\implies (a+S)*S$　　　　　　$\implies (S+S)*a$

$\implies (a+a)*S$　　　　　　$\implies (S+a)*a$

$\implies (a+a)*a$　　　　　　$\implies (a+a)*a$



The Grammer Is Unambigous
because It has only one RMD and LMD

6) $S \rightarrow (L) | a$

$L \rightarrow L, S | S \implies (a.a)$

$S \xrightarrow{lm} (L)$　　　　　　　$S \xrightarrow{rm} (L)$

$\implies (L,S)$　　　　　　　　$\implies (L,S)$

$\implies (S,S)$　　　　　　　　$\implies (L,a)$

$\implies (a.S)$　　　　　　　　$\implies (S,a)$

$\implies (a.a)$　　　　　　　　$\implies (a.a)$



The Grammar Is Unambigous
because It has only one LMD and RMD

5.) $S \to S+S \mid SS \mid (S) \mid S*\mid a \quad \Rightarrow (a+a)*a$

$S \xrightarrow{lm} SS$          $S \xrightarrow{rm} SS$

$\Rightarrow S*S$               $\Rightarrow S*S$

$\Rightarrow (S)*S$             $\Rightarrow S*a$

$\Rightarrow (S+S)*S$           $\Rightarrow (S)*a$

$\Rightarrow (a+S)*S$           $\Rightarrow (S+S)*a$

$\Rightarrow (a+a)*S$           $\Rightarrow (S+a)*a$

$\Rightarrow (a+a)*a$           $\Rightarrow (a+a)*a$



The Grammar Is Unambigous
because It has only one RMD and LMD

6.) $S \to (L) \mid a$

$L \to L, S \mid S \quad \Rightarrow (a, a)$

$S \xrightarrow{lm} (L)$          $S \xrightarrow{rm} (L)$

$\Rightarrow (L, S)$             $\Rightarrow (L, S)$

$\Rightarrow (S, S)$             $\Rightarrow (L, a)$

$\Rightarrow (a, S)$             $\Rightarrow (S, a)$

$\Rightarrow (a, a)$             $\Rightarrow (a, a)$



The Grammar Is Unambigous
because It has only one LMD and RMD

7) $S \rightarrow asbs \mid bsas \mid \varepsilon \quad \Rightarrow aabbab$

LMD:

$S \xrightarrow{lm} asbs$

$\Rightarrow aasbsbs$

$\Rightarrow aabsbs$

$\Rightarrow aabbsasbs$

$\Rightarrow aabbasbs$

$\Rightarrow aabbabs$

$\Rightarrow aabbab$

$S \xrightarrow{lm} asbs$

$\Rightarrow aasbsbs$

$\Rightarrow aabsbs$

$\Rightarrow aabsbs$

$\Rightarrow aabbasbs$

$\Rightarrow aabbabs$

$\Rightarrow aabbab$

RMD:

$S \xrightarrow{rm} asbs$

$S \Rightarrow asb$

$\Rightarrow aasbsb$

$\Rightarrow aasbbsasb$

$\Rightarrow aasbbsasb$

$\Rightarrow aasbbsab$

$\Rightarrow aasbbab$

$\Rightarrow aabbab$

$S \xrightarrow{rm} asbs$

$\xrightarrow{rm} asbasbs$

$\xrightarrow{rm} asbasb$

$\Rightarrow asbab$

$\Rightarrow aasbsbab$

$\Rightarrow aasbbab$

$\Rightarrow aabbab$

parse tree:



The grammar is ambigous
since It has 2 LMD and 2RMD

8) bexpr → bexpr or bterm | bterm

bterm → bterm and bfactor | bfactor

bfactor → not bfactor | (bexpr) | true | false

not (true or false)

LMD :

bexpr $\xrightarrow{lm}$ bterm

$\Rightarrow$ bfactor

$\Rightarrow$ not bfactor

$\Rightarrow$ not (bexpr)

$\Rightarrow$ not (bexpr or bterm)

$\Rightarrow$ not (bterm or bterm)

$\Rightarrow$ not (bfactor or bterm)

$\Rightarrow$ not (true or bterm)

$\Rightarrow$ not (true or bfactor)

$\Rightarrow$ not (true or false)

RMD :

bexpr $\xrightarrow{rm}$ bterm

$\Rightarrow$ bfactor

$\Rightarrow$ not bfactor

$\Rightarrow$ not (bexpr)

$\Rightarrow$ not (bexpr or bterm)

$\Rightarrow$ not (bexpr or bfactor)

$\Rightarrow$ not (bexpr or false)

$\Rightarrow$ not (bterm or false)

$\Rightarrow$ not (bfactor or false)

$\Rightarrow$ not (true or false)

9) $E \rightarrow E+E \mid E*E \mid -E \mid (E) \mid id$

$$\Rightarrow id + id * id$$

LMD

$E \xrightarrow{lm} E+E \; (E \rightarrow E+E)$

$\Rightarrow E+E*E$

$\Rightarrow id+E*E$

$\Rightarrow id+id*E$

$\Rightarrow id+id*id$

$E \xrightarrow{lm} E*E$

$\Rightarrow E+E*E$

$\Rightarrow id+E*E$

$\Rightarrow id+id*E$

$\Rightarrow id+id*id$

RMD

$E \xrightarrow{rm} E+E$

$\Rightarrow E+E*E$

$\Rightarrow E+E*id$

$\Rightarrow E+id*id$

$\Rightarrow id+id*id$

$E \xrightarrow{rm} E*E$

$\Rightarrow E*id$

$\Rightarrow E+E*id$

$\Rightarrow E+id*id$

$\Rightarrow id+id*id$



The Grammar Is Ambigous

Since we got moore than 1 LMD and moore than 1 RMD for This Grammar

10. $S \rightarrow iEtS / iEtSeS / a \qquad E \rightarrow b$

$\Rightarrow$ If $E_1$ then If $E_2$ then $S_1$ Else $S_2$

1)



2)



The given grammar is ambiguous ∴ it has two different parse trees

11) $R \rightarrow \dot{R} | R | RR | R* | (R) | a | b | c \Rightarrow a | b * c$

Lmp 1

$R \rightarrow R | R$
$\overset{lm}{\Rightarrow} a | R$
$\overset{lm}{\Rightarrow} a | RR$
$\overset{lm}{\Rightarrow} a | R * R$
$\overset{lm}{\Rightarrow} a | b * R$
$\Rightarrow a | b * c$

Lmp 2

$R \rightarrow RR$
$\overset{lm}{\Rightarrow} R | RR$
$\Rightarrow a | RR$
$\Rightarrow a | R * R$
$\Rightarrow a | b * R$
$\Rightarrow a | b * c$

The given grammar is ambiguous ∴ it has LMDs.

# Eliminating ambiguity

Can be done in 2 methods

i) Dis-ambiguity rule    ii) Using precedence &
                              associativity of operators

## 1) Dis-ambigouty Rule

→ Some grammars corresponding the statements are ambigous, This is due to dangling-else.

The dangling else problem can be eliminated and thus ambiguity of the grammar can also be eliminated

Q) what is dangling else problem??

→ Consider the following grammar

$$S → iCtS \mid iCtSeS \mid a \qquad i/p \ string : ibtibtaea$$

$$C → b$$

where

    i → Stands for Keyword 'if'

    C → Stands for 'Condition' to be satisfied,
          and C is nonterminal

    t → Stands for keyword 'then'

    s → Stands for 'Statement' for non terminal

    e → Stands for keyword 'else'

    a → Stands for other statement

    b → Stands for other statement

Since the above grammar is ambigous, we get two different parse tree for the string ibtibtaea



(i)                                              (ii)

Since there are 2 parse true for the same string ibtibtaea the given grammar is ambigous.
obsove the fallowing points
→The first parse tree associates else with 2nd statement
→The second parse tree associates else with first If stmt

The ambiguity wheather to associate else with first If statement / second If-statement is called dangling-else problem.

Eg> Q> Eliminate ambiguity from the fallowing ambigous grammar:
$$S \rightarrow iCts | iCtSeS | a$$
$$C \rightarrow b$$

→ In all programming languages when If-statements are nested, the first parse tree is prefered. So, the general rule is "Match each else with closest unmatched then". This rule can be diructly incorporated into grammar and ambiguity can be eliminated as shown below:

step1) The matched stmt M is an If-else statement where the statement S before else and after else keyword is matched. This can be expressed as:
$$M \rightarrow iCtMeM$$

step 2): An Unmatched statement U is the one consisting of:
a) Simple If-statement where the statement S is matched statement / Unmatched statement. ∴ The equivalent production is → $U \rightarrow iCtS$

b) If-else statement where the statement before else is matched and statement after else is Unmatched. ∴ The equivalent production is: $U \rightarrow iCtMeU$

Step 3: The matched statement M and Un-matched statement U Can obtained using the statement S as shown below:

$$S \rightarrow M|U$$

So, the final grammar which is un-ambigous is shown below:

$$S \rightarrow M|U$$
$$M \rightarrow iCtMeM/a$$
$$U \rightarrow iCtS$$
$$U \rightarrow iCtMeU \quad c \rightarrow b$$

observe that the above grammar associates else with closest then and eliminates ambiguity from the grammar

## Eliminating ambiguity using precedence and Associativity

This method is explained using the following example:

eg:Q) Convert the ambigous grammar into Unambigous grammar:

$$E \rightarrow E*E|E-E$$
$$E \rightarrow E^{\wedge}E|E/E$$
$$E \rightarrow E+E$$
$$E \rightarrow (E)|id$$

The grammar Can be Converted Into unambigous grammar using the precedence of operators as shown well as associativity operators as shown below:

step1: Arrange the operators in increasing order of the precedence along with the associativity as shown below:

| operators | Associativity | non-terminal used |
|-----------|---------------|-------------------|
| +, - | LEFT | E |
| *, / | LEFT | T |
| ^ | RIGHT | P |

Since there are three levels of precedence, we associate three non-terminals: E, T and P. Also an extra non-terminal F, generating basic units in an arithmatic expression

step 2: The basic units in expression are id (identifier) and paranthesized expressions. the production corresponding to this can be written as:

$$F \rightarrow (E) \mid id$$

step3: The next highest priority operator is ^ and it is right associative. So, the production must start from the non-terminal P and it should have right recursion as shown below:

$$P \rightarrow F \wedge P \mid F$$

step4: The next highest priority operators are * and / and they are left associative. So, the production must start from the non-terminal T and it should have left recursion as shown below:

$$T \rightarrow T * P \mid T/P \mid P$$

step5: The next highest priority operators are + and − and they are left associativity. So, the production must start from the non-terminal E and it should have left recursion as shown below:

$$E \rightarrow E + T \mid E - T \mid T$$

step6: The final grammar which is unambigous grammar can be written as shown below:

$$E \rightarrow E + T \mid E - T \mid T$$
$$T \rightarrow T * P \mid T/P \mid P$$
$$P \rightarrow F \wedge P \mid F$$
$$F \rightarrow (E) \mid id$$

8) Convert the following ~~Ambigous~~ grammar into Unambigous grammar by considering * and - operators lowest priority and they are left associative, / and + operators have the highest priority and are right associatituve and ∧ operator have the highest priority and are right associativity and ∧ operator has precedence in between and It is left associativity.

$$E \rightarrow E + E$$
$$E \rightarrow E - E$$
$$E \rightarrow E \wedge E$$
$$E \rightarrow E * E$$
$$E \rightarrow E / E$$
$$E \rightarrow (E) | id$$

→ The grammar can be converted into unambigous grammar using the precednce of operators as well as associativity operators as shown below:

step1: Arrange the operators in Increasing order of the precedence along with associativity as shown below:

| precednce | operators | Associativity | non-terminal used |
|---|---|---|---|
| (lowest) | *, - | LEFT | E |
| | ∧ | LEFT | P |
| (highest) | /, + | RIGHT | T |

Since there are three levels of precednce we associate three non-terminals: E, P and T. Also use an extra non-terminal F generating basic units in an arith--matic expression

step2: The basic units in expression are id (identifier) and parenthesized expressions. The production corresponding to this can be written as:

$$F \rightarrow (E) | id$$

step 3: The next highest priority operators are + and /
They are right associative, So, The production
must start from the nonterminal T and it should
be right recursive in RHS q the production as
shown below:

$T \rightarrow F + T \mid F / T \mid F$

step 4: The next highest priority operator is $\wedge$ and It
Is left associative. So, The production must start
from the non-terminal P and It should be left
recursive in RHS of the production as shown
below:

$P \rightarrow P \wedge T \mid T$

step 5: The next highest priority operators are $*$
and $-$ and They are left associative. So, the
production must start from the non-terminal
E and it should be left recursive in RHS of
the production as shown below;

~~$P \rightarrow P \wedge T \mid T$~~     $E \rightarrow E + P \mid E - P \mid P$

~~step 6: The next highest priority operators are $*$~~
~~and $-$ and they~~

step 6: The final grammar which Is unambigous can
be written as shown below:

$E \rightarrow E + P \mid E - P \mid P$
$P \rightarrow P \wedge T \mid T$
$T \rightarrow F + T \mid F / T \mid F$
$F \rightarrow (E) / id$

Q> Define Ambiguity ? Show that the following grammar
    is ambigous.

    R → R'|'R |RR | R* | (R) | a|b|c for input storing
    a|b*c

Give an unambigous grammar for the above grammar
such that precedence order from lowest to highest are
Concatenation, *, |, (), identifier and all are left to
right associativity

Ans: The grammar is said to be ambigous if it has
    more than one LMD / more than one RMP.

                    on
    If there are two different parse trees for the input
    string by applying LMD / by applying RMD

    → i|p storing: a|b*c                          parse tree:

    LMD 1          ҍLMD 2

    R → R'|'R       R → RR
    ⟹ a'|'R        ⟹ R'|'RR
    ⟹ a'|'RR       ⟹ a'|'RR
    ⟹ a'|'R*R      ⟹ a'|'R*R
    ⟹ a'|'b*R      ⟹ a|b*R
    ⟹ a'|'b*c      ⟹ a|b*c

    It has two LMD's ∴ the
    given grammar is Unambigous

→ Unambigous grammar

1. Arrange the operators in the ascending order with the precednce and associativity

| operators | Associativity | non-terminal used |
|---|---|---|
| . | LEFT | R |
| * | LEFT | S |
| \| | LEFT | T |

2. The basic units in expression are (R) and a,b,c we use additional non-terminal U for generating those    U → (R)|a|b|c

3. The nxt highest priority operators | and It ls left associative. So the production must start from the non-terminal T and It must be left recursive is RHS of the production

$$T → T'|'U|U$$

4. The next highest priority operators is * and ls left associative. So the production must start from non-terminal S and it is a unary operators

$$S → T*|T$$

5. The next highest priority operator ls Concatenation. and ls left associative. So the production must start from the non-terminal R and It should have left recursion as

$$R → RS|S$$

Step 6: The final grammar which is unambigous can be written as

$$R \rightarrow RS|s$$
$$S \rightarrow T*|T$$
$$T \rightarrow T'|'U|U$$
$$U \rightarrow (R)|a|b|c$$

The parse tree for the i/p a|b*c is

R
├── R
│   └── s
│       └── T
│           ├── T
│           │   └── U
│           │       └── a
│           ├── '
│           └── U
│               └── b
│   *
└── S
    └── T
        └── U
            └── C

## Left Recursion :

<u>defn:</u>

If Non-terminal Symbol and the 1st symbol of the production are Same, then It is left recursion.

General form: $A \longrightarrow A\alpha_1 \mid A\alpha_2 \mid A\alpha_3 \mid \beta$

$$\Downarrow$$

$$A \longrightarrow \beta A'$$
$$A' \longrightarrow \alpha_1 A' \mid \alpha_2 A' \mid \varepsilon$$

## Algorithm for Left Recursion Elimination

Algorithm Left-Recursion

i/p : Grammar $G$ with no cycles or $\varepsilon$-production

o/p : An equivalent grammar with no left-recursion

Method : Apply the algorithm to $G$. Note that the resulting non-left-recursive grammar may have $\varepsilon$-productions.

1. Arrange the non-terminals in Some order $A_1, A_2 \dots A_n$
2. For (each i from 1 to n) {
3.     for (each j from 1 to i-1) {
4.         replace each production of the form $A_i \rightarrow A_j \gamma$ by the productions $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$ where $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$ are all current $A_j$-productions
5.     }
6.     Eliminate the Immediate left recursion among the $A_i$-production
7. }

Example on Removing/Eliminating left Recursion

1. $E \longrightarrow \underset{A}{E} + \underset{\alpha}{T} \mid \underset{\beta}{T}$

   $\Downarrow$

   $E \longrightarrow TE'$  $\longrightarrow$ Epsilon
   $E' \longrightarrow +TE' \mid \varepsilon$

   [ here $E \longrightarrow E+T \mid T$
   (both are Same)
   $\therefore$ The grammar contains
   left recursion ]

2. $T \longrightarrow \underset{A}{T} * \underset{\alpha}{F} \mid \underset{\beta}{F}$  $\Longrightarrow$  $\boxed{\begin{array}{l} A \longrightarrow \beta A' \\ A' \longrightarrow \alpha A' \mid \varepsilon \end{array}}$

   $\Downarrow$

   $T \longrightarrow FT'$
   $T' \longrightarrow *FT' \mid \varepsilon$

3. $S \longrightarrow \underset{A}{S} \underset{\alpha}{(S)S} \mid \underset{\beta}{\varepsilon}$

   $\Downarrow$

   $S \longrightarrow S'$
   $S' \longrightarrow (S)SS' \mid \varepsilon$

4. $S \longrightarrow SS+ \mid SS* \mid a$

   $\Downarrow$

   $S \longrightarrow aS'$
   $S' \longrightarrow S+S' \mid S*S' \mid \varepsilon$

5. $E \longrightarrow E+T \mid T$
   $T \longrightarrow T*F \mid F$
   $F \longrightarrow (E) \mid id$

   $\Downarrow$

   $E \longrightarrow TE'$
   $E' \longrightarrow +TE' \mid \varepsilon$
   $T \longrightarrow FT'$
   $T' \longrightarrow *FT' \mid \varepsilon$
   $F \longrightarrow (E) \mid id$

## Left factoring

General form: $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \alpha\beta_3 \cdots \mid \alpha\beta_n \mid \gamma$

$$\Downarrow$$

$$A \rightarrow \alpha A' \mid \gamma$$
$$A' \rightarrow \beta_1 \mid \beta_2 \cdots \mid \beta_n$$

## Algorithm for left factoring

Algorithm left_factoring

i/p: Grammar $G$

o/p: An equivalent left-factored grammar

method: For each nonterminal $A$, find the longest prefix $\alpha$ common to two or more of its alternatives.

If $\alpha \neq \varepsilon$ – i,e there is a nontrivial common prefix – replace all of the A-productions

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \cdots \mid \alpha\beta_n \mid \gamma,$$

where $\gamma$ represents all alternatives that do not begin with $\alpha$, by

$$A \rightarrow \alpha A' \mid \gamma$$
$$A' \rightarrow \beta_1 \mid \beta_2 \mid \cdots \mid \beta_n$$

Here $A'$ is new non-terminal. Repeatedly apply this transformation until no two alternatives for a non-terminal have a common prefix

Examples on left factoring

1. $S \rightarrow \underset{\alpha}{SS} \underset{\beta_1}{+} \mid \underset{\alpha}{SS} \underset{\beta_2}{*} \mid \underset{\gamma}{a}$

   $S \rightarrow SSS' \mid a$

   $S' \rightarrow + \mid *$

   [The common terminals we have to take as $\alpha$ and the remaining term we have to take as $\beta_1, \beta_2$ -- So on in each production

2. $S \rightarrow 0S' \mid S$   $S \rightarrow \underset{\alpha}{0S} \underset{\beta_1}{1} \mid \underset{\alpha}{0} \underset{\beta_2}{1}$
   $S'$

   $\Downarrow$

   $S \rightarrow 0S' \mid \epsilon$

   $S' \rightarrow S1 \mid 1$

3. $S \rightarrow iEts \mid iEtses \mid a$
   $E \rightarrow b$
   $\Downarrow$

   $S \rightarrow \underset{\alpha}{iEts} \underset{\beta}{} \mid \underset{\alpha}{iEts} \underset{\beta_2}{es} \mid \underset{\gamma}{a}$

   $\Downarrow$

   $S \rightarrow iEtss' \mid a$    [Since $\beta_1$ is empty we'll take $\beta_1$ as $\epsilon$
   $S' \rightarrow \epsilon \mid es$
   $E \rightarrow b$

**Top down parsers :**

→ dfn: Is a parser of an i/p string of token by tracing out the steps in a left most derivation, it derives the string from the start symbol

→ It is termed as topdown because the parse tree is traversed in a preorder way that is from the root node to the leaf node

→ It has various types.

```
            Top-down parser
           /              \
   Backtracking      Nonback-tracking
                      /          \
                 Recursive      table
                 descent        driven
```

**i) Backtracking**

→ Backtracking tries different possibilities for parsing an i/p string by backing up on arbitrary amount in the i/p if any possibilities fails

→ These are more powerfull but very much slower, as they require exponential time to parse, hence they're not available suitable for practical compilers

eg:  S → CAd
     A → abla
     i/p string → cabd
     i/p string → cad

ii) Non-backtracking parsers:

i) Recursive descent          ii) Table driven

i) Recursive descent parser:

→ Recursive descent parsers are more versatile &
  Suitable for handwritten parsers

→ It helps to study the method for parsing and
  Serves as basis for topdown parses

$$S \rightarrow CAd$$
$$A \rightarrow ab|a$$

here, the grammar rule of a non-terminal A is
give as a defn of procedure call which will recognize A

a) The right hand side of a grammar rule, specifies the
   structure of the code for the procedure.

b) The sequence of terminals on the right hand side
   Corresponds to the i/p matches while the sequence
   of non-terminals are calls with the corresponding
   procedure

$$NT = \{S, A\}$$
$$T = \{a, b, c, d\}$$

```
procedure S()
  If (input == 'c')
  {
          Advance();
          A();
          If (input == 'd')
          {
```

```
                    Advance();
                    return (true);
            }
            else
            {
                return(false); }
        }
        else {
            return(false); }
    }

    procedure A()
    {
        isave = in-ptr;
        If (input == 'a')
        {
            Advance();
            If (input == 'b')
            {
                Advance();
                return (true);
            }
            else {
                in-ptr = isave
                If (input == 'a')
                {
                    Advance();
                    return (true);
                }
            }
            return (false);
        }
    }
```

isave:

Saves the ilp pointer position before each alternate production to facilitate backtracking whenever a terminal is encountered the ilp pointer

Advances the next position If alternate phase the in pointer retoraces to the previous position to trace the next alternate

Advance():

advance Is a procedure that Is written to advance the ilp pointer to the next position on a successfull completion of the parsing action the parser returns a true value

drawbacks of Recursive descent parsings

1. left recursion → It has the procedure production of the form A → Aα
   the parsers goes into Infinite loop
   eg: A → Abla
       ilp → abb
   The device storing abb There Is an ambiguity as to how many times the nonterminals has to be expanded

2. Backtracking: It occurs where there Is more than one alternate In the production to be tried while parsing the ilp string

$S → CAd$
$A → abla$    ilp: Cabd


no backt

i/p string : cad



⇒ Backtrack ∵ i/p symbol is 'd' but ptr is pointing to b

a.3. It is Very difficult to Identify the posn of the errors

Example : Recursive descent

Q> Write a recursive descent parser for the following grammar

$E \rightarrow TE'$

$E' \rightarrow +TE' | \epsilon$

$T \rightarrow FT'$

$T' \rightarrow *FT' | \epsilon$

$F \rightarrow (E) | id$

i/p : $id + id * id$

→ 
```
→ procedure E()
   {  If (input == 'T')
      T();
   }  Eprime();
   procedure T()
   {  F();
   }  Tprime();
   procedure F().
   {
      If (input == "C")
      {  Advance();
         E();
      If (input == ")")
         return(True);
      else
         return(False);
   }
```

```
elseif(input == "id")
{
    Advance();
    return(True);
}
else
    return(False);
}
procedure Eprime()
{
    If (input == "+")
    {
        Advance();
        T();
        Eprime();
        return(True);
    }
    return(False);
}
procedure Tprime()
{
    If (input == "*")
    {
        Advance();
        F();
        Tprime();
        return(True);
    }
    else
        return(False);
}
```

ii) Table driven :

→ Table driven is also called as predictive parsing

→ predictive parser is a recursive decent parser, which production has the capability to predict which production is to be used to replace the i/p storing

→ The predictive parser does not suffer from backtracking

Predictive Parsers / LL(1) parsers / Table driven parsers.

Q) Why do we need a FIRST and FOLLOW set

Consider the below given, top down approach for the example:



here, T' has two options
i.e T' → *FT' and T' → ε we are
in ambiguity situation to replace
with which RHS of the production.
If we have n alternate productions
for the same Non-terminal then
we have to backtrack (n-1) times
so, It is very tedious process.
So, we find FIRST & FOLLOW set.

FIRST AND FOLLOW SETS:

i) FIRST:

    i) If $x$ is a terminal, then $FIRST(x) = \{x\}$

    ii) If $X$ is a nonterminal, and $X \to Y_1 Y_2 \cdots Y_k$ is a production for som $k \geq 1$, Then place $a$ in $FIRST(X)$ If for some $i$, $a$ is in $FIRST(Y_i)$ and $\varepsilon$ is in all of $FIRST(Y_1) \cdots FIRST(Y_{i-1})$ That is $Y_1$

    $Y_{i-1} \xRightarrow{x} \varepsilon$, If $\varepsilon$ is in $FIRST(y_j)$ for all $j = 1 \cdots k$ then add $\varepsilon$ to $FIRST(X)$, For example, every-thing in $FIRST(y_i)$ is Surely in $FIRST(X)$ If $y_i$ doesn't drive $\varepsilon$, Then we add nothing more to $FIRST(X)$, but If $y_r \xRightarrow{*} \varepsilon$, then we add

        $FIRST(Y_2)$ and So on

    iii) If $X \to \varepsilon$ is a production, then add $\varepsilon$ to $FIRST(x)$.

Some of the general forms on how to find the FIRST(x) are given below:

1) $X \to aB$
    $FIRST(X) = \{a\}$
    $X \to \varepsilon$
    $FIRST(X) = \varepsilon$

2) $X \to ABC$
    $A \to a$
    $B \to b$
    $C \to c$
    $FIRST(X) = FIRST(A)$
             $= \{a\}$

3) $X \to ABC$
    $A \to a | \varepsilon$
    $B \to b$
    $C \to c$
    $FIRST(X) = FIRST(A)$
             $= \{a, \varepsilon\}$
                 $\downarrow$
                  $FIRST(B)$
                $= \{a, b\}$

4) $X \to ABC$
   $A \to a | \varepsilon$
   $B \to b | \varepsilon$
   $C \to c$

$FIRST(X) = FIRST(A)$
$\qquad = \{a, \varepsilon\}$
$\qquad \downarrow$
$\qquad FIRST(B)$
$\qquad\quad \{b, \varepsilon\}$
$\qquad = \{a, b, c\} \longrightarrow FIRST(C)$
$\qquad\qquad\qquad\qquad \{c\}$

5) $X \to ABC$
   $A \to a | \varepsilon$
   $B \to b | \varepsilon$
   $C \to c | \varepsilon$

$FIRST(X) = FIRST(A)$
$\qquad = \{a, \varepsilon\}$
$\qquad \downarrow$
$\qquad FIRST(B)$
$\qquad\quad \{b, \varepsilon\}$
$\qquad\qquad \longrightarrow FIRST(C)$
$\therefore = \{a, b, c, \varepsilon\} \qquad \longrightarrow \{c, \varepsilon\}$

ii) FOLLOW :

i) place $ in FOLLOW(S), where S is the start symbol, and $ is the input right endmarker

ii) If there is a production $A \to \alpha B \beta$, then everything in FIRST($\beta$) except $\varepsilon$ is in FOLLOW(B)

iii) If there is a production $A \to \alpha B$, on a production $A \to \alpha B \beta$, where FIRST($\beta$) contains $\varepsilon$, then everything in FOLLOW(A) is in FOLLOW(B).

How to find FOLLOW :

1) $A \to \underset{\alpha}{X} \underset{B}{B} \underset{\beta}{CD}$
   $C \to x$
   fallow(B) = {FIRST(CD)}
   $\qquad = \{c\}$

2) $A \to \underset{\alpha}{X} \underset{B}{BCD}{\beta}$
   Since $\beta = CD$
   $\therefore$ FOLLOW(B) = {FIRST($\beta$)}
   $\qquad = FIRST(CD)$
   $\qquad = \{c, d\}$

3> $A \longrightarrow \overset{\alpha}{\widetilde{X}} \overset{B}{\widetilde{B}} \overset{\beta}{\widetilde{CD}}$

$C \longrightarrow c | \varepsilon$

$D \longrightarrow d | \varepsilon$

FOLLOW $(B) =$ first $(\beta)$

$\qquad =$ first $(CD)$

$\qquad = \{c, d\} +$ FOLLOW $(A)$

4> $A \longrightarrow \overset{\alpha}{\widetilde{X}} \overset{B}{\underset{B}{\widetilde{B}}} \overset{\beta}{\widetilde{\phantom{x}}}$  (since $\beta = \varepsilon$, here)

fallow $(B) =$ FALLOW (left most nonterminal)

$\qquad =$ FALLOW $(A)$

5> $A \longrightarrow \overset{\alpha}{\widetilde{X}} \overset{B}{\widetilde{B}} \overset{\beta}{\overbrace{CD Be}}$ ——①

$C \longrightarrow c | \varepsilon$  $\overset{\alpha \quad B \quad \beta}{}$ ——②

$D \longrightarrow d$

FALLOW $(B) =$ FIRST $(CDBe)$

$\qquad = (C, D)$ ————①

$+$

FALLOW $(B) =$ FIRST $(e)$

$\qquad = \{e\}$

$\therefore$ FAL $(B) =$ ① $+$ ②

FAL $(B) = \{c, d, e\}$

Find the FIRST and FOLLOW set for the following grammars

1. $E \rightarrow TE'$

   $E' \rightarrow +TE' | \varepsilon$

   $T \rightarrow FT'$

   $T' \rightarrow *FT' | \varepsilon$

   $F \rightarrow (E) | id$

2. $S \rightarrow iEts | iEtses | a$

   $E \rightarrow b$

3) S →, G₁H;
   G₁ → aF
   F → bF | ε
   H → KL
   K → m | ε
   L → n | ε

5) S → aBDh
   B → ec
   C → bc | ε
   D → EF
   E → g | ε
   F → f | ε

4) S → aB | ac | sd | se
   B → bBc | f
   C → g

D → EF
E → g | ε → epsilon
F → f | ε

6. S → (L) | a
   L → L, s | s

7) S → L = R | R
   L → *R | id
   R → L

8. S → AaAb | BbBa
   A → ε
   B → ε

9) S → aABb
   A → c | ε
   B → d | ε

10. stmt_sequence → stmt stmt_sequence'
    stmt_sequence' → ; stmt_sequence | ε
    stmt → s

11) S → asbs | bsas | ε

12) S → a | ↑ | (T)
    T → T, S | s

13) S → As | b
    A → SA | a

Answer:

1) $E \rightarrow TE'$
   $E' \rightarrow +TE' | \varepsilon$
   $T \rightarrow FT'$
   $T' \rightarrow *FT' | \varepsilon$
   $F \rightarrow (E) | id$

| | E | E' | T | T' | F |
|---|---|---|---|---|---|
| FIRST | ( | + | ( | * | ( |
| | id | $\varepsilon$ | id | $\varepsilon$ | id |
| FOL | $ | $ | + | + | * |
| | ) | ) | $ | $ | |
| | | | ) | ) | |

here $FOLLOW(E) = F \rightarrow \underset{\alpha}{(}\underset{B}{E}\underset{\beta}{)}$

$$FOLLOW(E) = FIRST())$$
$$= \{ ) \}$$

$FOLLOW(E') = E \rightarrow \underset{\alpha}{T}\underset{B}{E'}\underset{\beta}{}$   and   $E' \rightarrow \underset{\alpha}{+T}\underset{B}{E'}\underset{\beta}{}$

$FOL(E') = FIRST(\beta) \downarrow \varepsilon$      $FOL(E') = FIRST(\beta) \downarrow \varepsilon$

$\therefore FOL(E') = FOL \cdot (E)$      $\therefore FOL(E') = FOL(E$

So, on   $FOLLOW(T) = E' \rightarrow \underset{\alpha}{+T}\underset{B}{E'}\underset{\beta}{}$      $E \rightarrow \underset{\alpha}{T}\underset{B}{E'}\underset{\beta}{}$

$$= FIRST(\beta)$$      $$= FIRST(\beta)$$
$$= FIRST(E')$$      $$= FIRST(E')$$
$$= \{ + \}$$      $$= \{ + \}$$

<u>Note</u>: we should not work $\underline{\varepsilon}$ in the <u>FOLLOW</u> set

2) $S \rightarrow iEk \mid iEtses \mid a$
$E \rightarrow b$

| | S | E |
|---|---|---|
| FIRST | i<br>a | b |
| FOLLOW | $ <br>e | t |

FOLLOW(S) = $S \rightarrow iEk$
　　　　　　　　　　$B \bar{P}$

$S \rightarrow iEtses$
　　　　$\alpha B \beta$

FOLLOW(S) = FIRST($\beta$)
　　　　　　　= FOLLOW(S) $\overset{\leftarrow}{e}$

FOLLOW(S) = FIRST($\beta$)
　　　　　　　= FIRST(es)
　　　　　　　= $\{e\}$

3) $S \rightarrow , G_1 H ;$
$G_1 \rightarrow aF$
$F \rightarrow bF \mid \epsilon$ — Epsilon
$H \rightarrow kL$
$K \rightarrow m \mid \epsilon$
$L \rightarrow n \mid \epsilon$

| | S | $G_1$ | F | H | K | L |
|---|---|---|---|---|---|---|
| FIRST | , | a | b<br><br>$\epsilon$ | m<br><br>n<br>$\epsilon$ | m<br><br>$\epsilon$ | n<br><br>$\epsilon$ |
| FOLLOW | $ | m<br>n<br>; | m<br>n<br>; | ; | n<br><br>;<br>; | ; |

4) S → aB | ac | Sd | Se
   B → bBC | f
   G → g

| | S | B | C |
|---|---|---|---|
| FIRST | a | b, f | g |
| FOLLOW | \$, d, e | \$, d, e, c | \$, d, e |

5) S → aBDh
   B → ec
   C → bC | ε
   D → EF
   E → g | ε
   F → f | ε

| | S | B | C | D | E | F |
|---|---|---|---|---|---|---|
| FIRST | a | e | b, ε | g, f, ε | g, ε | f, ε |
| FOLLOW | \$ | g, f, h | g, f, h | h | h, f | h |

6) S → (L) | a
   L → L, S | S

| | S | L |
|---|---|---|
| FIRST | (, a | (, a |
| FOLLOW | \$, ), , | ), , |

8) $S \rightarrow AaAb \mid BbBa$
$A \rightarrow \varepsilon$
$B \rightarrow \varepsilon$

|  | S | A | B |
|---|---|---|---|
| FIRST | a<br>b | $\varepsilon$ | $\varepsilon$ |
| FOLLOW | $ | a<br>b | b<br>a |

9) $S \rightarrow aABb$
$A \rightarrow c \mid \varepsilon$
$B \rightarrow d \mid \varepsilon$

|  | S | A | B |
|---|---|---|---|
| FIRST | a | c | d |
|  |  | $\varepsilon$ | $\varepsilon$ |
| FOLLOW | $ | d<br>b | b |

10)
12) $S \rightarrow a \mid \uparrow \mid (T)$
$T \rightarrow T, S \mid S$

|  | S | T |
|---|---|---|
| FIRST | a<br>$\uparrow$<br>( | a<br>$\uparrow$<br>( |
| FOLLOW | $<br>)<br>, | )<br>, |

**8)** S → As|b
A → SA|a

|        | S   | A   |
|--------|-----|-----|
| FIRST  | b<br>a | a |
| FOLLOW | \$<br>a<br>b | a<br>b |

**9)** S → L =R|R
L → *R|id
R → L

|        | S   | L   | R   |
|--------|-----|-----|-----|
| FIRST  | *<br>id | *<br>id | *<br>id |
| FOLLOW | \$ | =<br>\$ | \$<br>= |

**10)** stmt_sequence → stmt stmt_sequence'
stmt_sequence' → ; stmt_sequence'|ε
stmt → s

|        | stmt_sequence | stmt_sequence' | stmt |
|--------|---------------|----------------|------|
| FIRST  | s             | ;<br>ε         | s    |
| FOLLOW | \$            | \$             | ;<br>\$ |

18) $S \rightarrow asbs \mid bsas \mid \varepsilon$

|  | S |
|---|---|
| FIRST | a |
|  | b |
|  | $\varepsilon$ |
| Follow | \$ |
|  | a |
|  | b |

$FOLLOW(S) \Rightarrow S \rightarrow \underset{\alpha \; B \; \beta}{a \underline{s} b s}$  $\qquad S \rightarrow \underset{\alpha \; B \; \beta}{b \underline{s} a s}$  $\qquad S \xrightarrow{\quad} \varepsilon$

$\qquad\qquad\qquad = FIRST(\beta)$  $\qquad\qquad = FIRST(\beta)$

$\qquad\qquad\qquad = FIRST(b\$)$  $\qquad\qquad = FIRST(aS)$

$\qquad FOLLOW(S) = \{b\}$  $\qquad\qquad\qquad = \{a\}$

$FOL(S) \rightarrow \underset{\alpha \; B \; \beta}{a s \underline{b} s}$  $\qquad Fa(S) \rightarrow \underset{\alpha \; B \; \beta}{b s \underline{a} s}$

$\qquad\qquad\qquad = FIRST(\beta)$  $\qquad\qquad = FIRST(\beta)$
$\qquad\qquad\qquad\qquad\searrow \varepsilon$  $\qquad\qquad\qquad\searrow \varepsilon$

$\qquad FOL(S) = FOL(S)$  $\qquad FOL(S) = FOL(S)$

Top-down Parses

predictive parsing table/LL(1) grammar/Table Driven predictive parses

→ lookahead symbol

LL(1) grammar

↓ └→ Left most derivation

scan the i/p from left to right

## steps:

1) Eliminate left recursion from the grammar
2) perform left factoring
3) find the FIRST and FOLLOW set
4) Construct the predictive parsing table
5) check wheather the given i/p string is accepted/not

## Algorithm for Constructing predictive parsing table

INPUT : Grammar $G$

OUTPUT : parsing table M

METHOD: For each production $A \to \alpha$ of the grammar, do the following

1. For each terminal $\alpha$ in FIRST(A), add $A \to \alpha$ to M[A, $\alpha$]

2. If $\varepsilon$ is in FIRST($\alpha$), Then for each terminal b in FOLLOW(A), add $A \to \alpha$ to M[A, b], If $\varepsilon$ is in FIRST($\alpha$) & $ is in FOLLOW(A), add $A \to \alpha$ to M[A, $] as well

## Predictive parsing Algorithm

INPUT: A string w and a parsing table m for a grammar $G$.

OUTPUT: If w is in L(G) and LMD of w; otherwise an error condition

Input: 

| | | | a | + | b | $ |
|---|---|---|---|---|---|---|

METHOD: Initially, the parser is in a configuration with w$ in the i/p buffer and the start symbol S on top of the stack, above $. The pgm in fig. uses the predictive parsing table M to procedure a predictive parse for the i/p

set 'ip' to point to the frist symbol of w;
set 'x' to the top stack symbol;
while (x ≠ $) { /* stack is not empty */)
If (x is a) pop the stack & a advance ip:
else if (x is a terminal) error();
else if (M [x, a] is an error entry) error();
else if (M [x, a] = x → y, , y₀ --- yₖ) {
output the production X → y₁, y₂ -- yₖ;
pop the stack
push yₖ, yₖ₋₁ ... y₁ onto the stack, with y₁ on-top
}
Set x to the top stack symbol;



fig: model of a table driven predictive parser

Checking wheather the given grammar is LL(1) or not without using parsing table

A grammar is LL(1) iff whenever, $A \rightarrow \alpha | \beta$ are two distinct productions of $G$, the following conditions hold

i) For no terminal 'a' do with $\alpha$ and $\beta$ derive strings beginning with $a \Rightarrow FIRST(\alpha)$ and $FIRST(\beta)$ are disjoint.

ii) Atmost one of $\alpha$ and $\beta$ can derive the empty string $\Rightarrow$ either $FIRST(\alpha) \rightarrow \varepsilon$ or $FIRST(\beta) \rightarrow \varepsilon$ but not both.

iii) If $\beta \overset{*}{\Rightarrow} \varepsilon$, then $\alpha$ does not derive any string beginning with a terminal in FOLLOW (A). Likewise, if $\alpha \overset{*}{\Rightarrow} \varepsilon$, then $\beta$ does not derive any string beginning with a terminal in FOLLOW (B)

$\Rightarrow FIRST(\alpha)$ and FOLLOW (A) are disjoint or FIRST($\beta$) and FOLLOW(A) are disjoint

General forms:

① $A \rightarrow \overset{\alpha}{a B} | \overset{\beta}{a c}$

$FIRST(\alpha) = \{a\}$

$FIRST(\beta) = \{a\}$

are not disjoint, a/c to the algorithm in any production

eg: $A \rightarrow aB | bC$

$FIRST(\alpha) = \{a\}$

$FIRST(\beta) = \{b\}$ are disjoints

2) $A \rightarrow Bc \mid CD$      either $FIRST(\beta) \Rightarrow \varepsilon$ or
   $B \rightarrow b \mid \varepsilon$           $FIRST(\alpha) \Rightarrow \varepsilon$
   $C \rightarrow c \mid \varepsilon$          but not both

3) i) $A \rightarrow a \mid B$
    $B \rightarrow cAa \mid \varepsilon$
    $FIRST(\alpha) = \{a\}$
    $FOLLOW(A) = \{\$, a\}$ are not disjoint

   ii) $A \rightarrow B \mid a$
     $B \rightarrow cAa \mid \varepsilon$
     $FIRST(\beta) = \{a\}$ and $FOLLOW(A) = \{\$, a\}$
        are not disjoint

Examples:

1. $S \rightarrow iEtss' \mid a$
   $s' \rightarrow es \mid \varepsilon$
   $E \rightarrow b$

|  | S | $s'$ | E |
|---|---|---|---|
| FIRST | i | e | b |
|  | a | $\varepsilon$ |  |
| FOLLOW | \$ | \$ | t |
|  | e | e |  |

$$S \rightarrow \underset{\alpha}{iEtss'} \mid \overset{\beta}{a}$$

a. $FIRST(\alpha) \cap FIRST(\beta) = \phi$
    $\{i\} \cap \{a\} = \phi$

b. neither of $\alpha$ or $\beta$ are $\varepsilon$

1. $S \rightarrow \underset{\alpha}{i\underline{E tss'}} \mid \overset{\beta}{a}$

a. $FIRST(\alpha) \cap FIRST(\beta) = \phi$
    $\{i\} \cap \{a\} = \phi$

b. neither of $\alpha$ or $\beta$ are $\Rightarrow \varepsilon$

c. $\beta \Rightarrow \varepsilon$, then $FIRST(\alpha) \cap FOLLOW(A) = \phi$
                 $FIRST(es) \cap FOLLOW(s') = \phi$
                     $\{e\} \cap \{\$, e\} \neq 0$

∴ The given grammar is not   Condition fails
LL(1).

$S^1 \rightarrow e \overset{\alpha}{S} | \overset{\beta}{\varepsilon}$

a) $FIRST(\alpha) \cap FIRST(\beta) = \phi$

$\{e\} \cap \{\varepsilon\} = \phi$

b) The given grammar

$\beta \Rightarrow \varepsilon$ but $\alpha \not\Rightarrow \varepsilon$

c) $\beta \Rightarrow \varepsilon$, then $FIRST(\alpha) \cap FOL(A) = \phi$

$FIRST(eS) \cap FOL(S^1) = \phi$

$\{e\} \cap \{\$ e\} \neq \phi$

Condition fails

∴ The given grammar is not LL(1)

2) $S \rightarrow S(S)S | \varepsilon \implies S \rightarrow S^1$

$\quad\quad\quad\quad\quad\quad\quad\quad S^1 \rightarrow (S) SS^1 | \varepsilon$

| | S | S¹ |
|---|---|---|
| FIRST | $\varepsilon$ $\varepsilon$ | $\varepsilon$ $\varepsilon$ |
| FOLLOW | $\$$ ) ( | $\$$ ) ( |

1. $S \rightarrow S^1$
   not required because we
   don't have $\beta$ production

2. $S^1 \rightarrow \underset{\alpha}{(S) SS^1} | \underset{\beta}{\varepsilon}$

   a. $FIRST(\alpha) \cap FIRST(\beta) \neq \phi$
      $\{\varepsilon\} \cap \{\varepsilon\} = \phi$

   b. only $\beta \Rightarrow \varepsilon$ and $\alpha \not\Rightarrow \varepsilon$

   c. $\beta \Rightarrow \varepsilon$, then
      $FIRST(\alpha) \cap FOLLOW(A) = \phi$
      $\{(\} \cap \{\$ ( )\} \neq \phi$

∴ The given grammar is not LL(1)

3. $S \rightarrow SS+ |SS*|a \Rightarrow S \rightarrow aS'$
$$S' \rightarrow S+S | S*S' | \varepsilon$$

$\Downarrow$ left factoring

$$\left.\begin{array}{l} \text{final} \\ \text{production} \end{array}\right\{ \begin{array}{l} S \rightarrow aS' \\ S' \rightarrow SS'' | \varepsilon \\ S'' \rightarrow +S' | *S' \end{array}$$

|        | S  | S' | S'' |
|--------|----|----|-----|
| FIRST  | a  | a  | +   |
|        |    | ε  | *   |
| FOLLOW | \$ | \$ | \$  |
|        | +  | +  | +   |
|        | *  | *  | *   |

1. $S \rightarrow aS'$

   not required

2. $S'' \rightarrow +S' | *S'$

   a. $FIRST(+S') \cap FIRST(*S') = \phi$

   $\{+\} \cap \{*\} = \phi$

   b. neither $\alpha$ or $\beta \Rightarrow \varepsilon$

   all condition are satisfied

3. $S' \rightarrow SS'' | \varepsilon$

   a. $FIRST(SS'') \cap FIRST(\varepsilon)$

   $\{a\} \cap \{\varepsilon\} = \phi$

   b. $\beta \Rightarrow \varepsilon$ but not $\alpha$

   c. $\beta \Rightarrow \varepsilon$ then

   $FIRST(SS'') \cap FOLLOW(S') = \phi$

   $\{a\} \cap \{\$+*\} = \phi$

   all conditions are satisfied

   $\therefore$ The grammar is LL(1)

4) $E \to E+T \mid T$          $E \to TE'$
   $T \to T*F \mid F$   $\Rightarrow$   $E' \to +TE' \mid \varepsilon$
   $F \to (E) \mid id$          $T \to FT'$
                        $T' \to *FT' \mid \varepsilon$
                        $F \to (E) \mid id$

| | E | E' | T | T' | F |
|---|---|---|---|---|---|
| FIRST | ( | + | ( | * | ( |
| | id | $\varepsilon$ | id | $\varepsilon$ | id |
| FOLLOW | \$ | \$ | + | + | + |
| | ) | ) | \$ | \$ | * |
| | | | ) | ) | \$ |
| | | | | | ) |

i) $E \to TE'$

ii) $E' \to +E' \mid \varepsilon$

   a. FIRST$(+TE) \cap$ FIRST$(\varepsilon) = \phi$

      $\{+\} \cap \{\varepsilon\} = \phi$

   b. $\beta \Rightarrow \varepsilon$ but $\alpha \not\Rightarrow \varepsilon$

   c. $\beta \Rightarrow \varepsilon$ then FIRST$(+TE) \cap$ FOL$(E') = \phi$

      $\{+\} \cap \{\$\} = \phi$

iii) $T \to FT'$

iv) $T' \to *FH \varepsilon$

   a. FIRST$(*FT') \cap$ FIRST$(\varepsilon) = \phi$

      $\{*\} \cap \{\varepsilon\} = \phi$

   b. $\beta \Rightarrow \varepsilon$ but $\alpha \not\Rightarrow \varepsilon$

   c. $\beta \Rightarrow \varepsilon$ then FIRST$(*FT') \cap$ FOL$(T') = \phi$

      $\{*\} \cap \{+\$\} = \phi$

v) $F \to (E) \mid id$

   a. FIRST$((E)) \cap$ FIRST$(id) = \phi$

      $\{(\} \cap \{id\} = \phi$

   b. neither of them are not $\Rightarrow \varepsilon$

checking wheather the given grammar is LL(1) or not
with Constructing the predictive parsing table

1. $E \rightarrow E+T|T$
   $T \rightarrow T*F|F$    i/p: id+id*id
   $F \rightarrow (E)|id$

i) Remove left Recursion
   $E \rightarrow TE'$
   $E' \rightarrow +TE'|\varepsilon$
   $T \rightarrow FT'$
   $T' \rightarrow *FT'|\varepsilon$
   $F \rightarrow (E)|id$

ii) Remove left factoring
   $\rightarrow$ here, not required

iii) find FIRST and FOLLOW Set

|  | E | E' | T | T' | F |
|---|---|---|---|---|---|
| FIRST | ( id | + $\varepsilon$ | ( id | * $\varepsilon$ | ( id |
| FOLLOW | $ ) | $ ) | + $ ) | + $ ) | * + $ ) |

iv) Construct the parsing table

|  | ( | id | ) | + | * | $ |
|---|---|---|---|---|---|---|
| E | $E \rightarrow TE'$ | $E \rightarrow TE'$ |  |  |  |  |
| E' |  |  | $E' \rightarrow \varepsilon$ | $E' \rightarrow +TE'$ |  | $E' \rightarrow \varepsilon$ |
| T | $T \rightarrow FT'$ | $T \rightarrow FT'$ |  |  |  |  |
| T' |  |  | $T' \rightarrow \varepsilon$ | $T' \rightarrow \varepsilon$ | $T' \rightarrow *FT'$ | $T' \rightarrow \varepsilon$ |
| F | $F \rightarrow (E)$ | $F \rightarrow id$ |  |  |  |  |

v)

| Stack | Input | Action |
|---|---|---|
| E$ | id+ id*id $ | |
| TE'$ | id+ id*id$ | push E→TE' |
| FT'E'$ | id+id*id $ | push T→FT' |
| idT'E'$ | id+id*id $ | push F→id |
| T'E'$ | +id*id $ | matched "id" |
| E'$ | +id*id$ | T'→ε |
| +TE'$ | +id*id$ | push E'→+TE' |
| TE'$ | id*id$ | matched '+' |
| FT'E'$ | id*id$ | push T→FT' |
| idT'E'$ | id*id$ | push F→id |
| T'E'$ | *id $ | matched 'id' |
| *FT'E'$ | *id$ | push T'→*FT' |
| FT'E'$ | id$ | matched '*' |
| idT'E'$ | id$ | push F→id |
| T'E'$ | $ | matched 'id' |
| E'$ | $ | T'→ε |
| $ | $ | E'→ε |

∴ The grammar is accepted the i/p storing
i.e, The i/p is passed Successfully

The grammar is LL(1) ∵ no multiple entries

2) $S \rightarrow iEts \mid iEtses \mid a$
$E \rightarrow b$

i/p: IF $E_1$ then if $E_2$ then $S_1$ else $S_2$
    If $b$ then if $b$ then $a$ else $a$

i) No left Recursion

ii) Remove left factoring
$S \rightarrow iEtSS' \mid a$
$S' \rightarrow esS \mid \varepsilon$
$E \rightarrow b$

iii) FIRST and FOLLOW Set

|  | S | S' | E |
|---|---|---|---|
| FIRST | i | e | b |
|  | a | $\varepsilon$ |  |
| FOLLOW | $ | $ | t |
|  | e | e |  |

iv) parsing table

|  | $ | i | e | b | a | t | |
|---|---|---|---|---|---|---|---|
| S |  | S→iEtSS' |  |  |  |  | The grammar is not LL(1) because It has multiple entries for the same terminal in a table |
| S' | S'→$\varepsilon$ |  | S'→es<br>S'→$\varepsilon$ |  |  |  |  |
| E |  |  |  | E→b |  |  |  |

v)

| Stack | Input | Action |
|---|---|---|
| S$ | iEtEtsres₂<br>ibtibtaea$ |  |
| iEtSS'$ | ibtibtaea$ | S→iEtSS' |
| EtSS'$ | btibtaea$ | matched i |
| btSS'$ | btibtaea$ | E→b |

| Stack | input | Action |
|---|---|---|
| ts s' $ | tibtaea $ | matched 'b' |
| ss' $ | ibtaea $ | matched 't' |
| iEtss's' $ | ibtaea $ | s→iEtss' |
| Etss's' $ | btaea $ | matched 'i' |
| btss's' $ | btaea $ | E→b |
| tss's' $ | taea $ | matched 'b' |
| ss's' $ | aea $ | matched t |
| as's' $ | aea $ | s→a |
| s's' $ | ea $ | matched 'a' |

⌞→ ambiguity wheather to push s'→es on s→ε

⟹ i/p: If E Then S else S
           ibtaea

| stack | input | Action |
|---|---|---|
| s$ | ibtaea $ | s→iEtss' |
| iEtss'$ | ibtaea$ | matched 'i' |
| Etss'$ | btaea$ | push E→b |
| btss'$ | btaea $ | match 'b' |
| tss'$ | taea$ | match t |
| ss'$ | aea$ | s→a |
| as'$ | aea$ | match a |
| s'$ | ea$ | |

⌞→ ambiguity, whether to push s'→es on s→ε

3) $S \rightarrow SS+ | SS* | a$    i/p: $aa+a*$

**i) Remove left recursion**

$S \rightarrow aS'$

$S' \rightarrow S+S | S*S' | \varepsilon$

**ii) Remove left factoring**

$S \rightarrow aS'$

$S' \rightarrow SS'' | \varepsilon$

$S'' \rightarrow +S' | *S$

**iii) find FIRST and FOLLOW set**

|        | $S$ | $S'$ | $S''$ |
|--------|-----|------|-------|
| FIRST  | $a$ | $a$  | $+$   |
|        |     | $\varepsilon$ | $*$ |
| FOLLOW | $\$$ | $\$$ | $\$$ |
|        | $+$ | $+$  | $+$   |
|        | $*$ | $*$  | $*$   |

**iv) find the predictive parsing table**

|       | $\$$ | $a$ | $+$ | $*$ |
|-------|------|-----|-----|-----|
| $S$   |      | $S \rightarrow aS'$ |     |     |
| $S'$  | $S' \rightarrow \varepsilon$ | $S' \rightarrow SS''$ | $S' \rightarrow \varepsilon$ | $S' \rightarrow \varepsilon$ |
| $S''$ |      |     | $S'' \rightarrow +S'$ | $S'' \rightarrow *S'$ |

The given grammar is LL(1) because. There are no multiple production

**v) parse the i/p string**

| stack | input | action |
|-------|-------|--------|
| $S\$$ | $aa+a*\$$ |  |
| $aS'\$$ | $aa+a*\$$ | $S \rightarrow aS'$ |
| $S'\$$ | $a+a*\$$ | match 'a' |
| $SS''\$$ | $a+a*\$$ | $S' \rightarrow aS'$ |
| $aS'S''\$$ | $a+a*\$$ | $S \rightarrow aS'$ |
| $S'S''\$$ | $+a*\$$ | match 'a' |
| $S''\$$ | $+a*\$$ | $S' \rightarrow \varepsilon$ |

| Stack | input | Stack |
|---|---|---|
| +s'$ | +a*$ | s"→+s' |
| s'$ | a*$ | match '+' |
| Ss"$ | a*$ | push s'→ss" |
| as's"$ | a*$ | push s→as |
| s's"$ | *$ | match 'a' |
| s"$ | *$ | push s'→ε |
| *s'$ | *$ | push s"→*s' |
| s'$ | $ | match * |
| $ | $ | push s'→ε |

The I/p string Is Successfully parsed

4)  $S→OS1/O1$  i/p string: 000111

i) Remove left recursion
→ not needed

ii) Remove left factoring
$S→OS'$
$S'→S1/1$

iii) FIRST and FOLLOW Set

|  | S | $ |
|---|---|---|
| FIRST | O | O |
| FOLLOW | $ <br> 1 | $ <br> 1 |

iv) Construct the predictive parsing table

|  | $ | O | 1 |
|---|---|---|---|
| S |  | S→OS' |  |
| S' |  | S'→S1 | S'→1 |

The grammar is LL(1), since it does not have any multiple production

v) parse the input string

| Stack | input | action |
|-------|-------|--------|
| S$ | 000111$ | |
| 0S'$ | 000111$ | push s→0s' |
| s'$ | 00111$ | match 'o' |
| S1$ | 00111$ | push s→s1 |
| 0s'1$ | 00111$ | s→0s' |
| s'1$ | 0111$ | match 'o' |
| S11$ | 0111$ | s→s1 |
| 0s'11$ | 0111$ | push s→0s' |
| s'11$ | 111$ | match 0 |
| 111$ | 111$ | s'→1 |
| 11$ | 11$ | match 1 |
| 1$ | 1$ | match 1 |
| $ | $ | match 1 |

The input string is successfully parsed

5) $S \rightarrow +SS \mid *SS \mid a$   i/p: $+*aaa$

i) Remove left recursion $\Rightarrow$ not required

ii) Remove left factoring $\Rightarrow$ Not required

iii) Constonet FIRST and FOLLOW set-

| | S |
|---|---|
| FIRST | $+$ |
| | $*$ |
| | $a$ |
| FOLLOW | $\$$ |
| | $+$ |
| | $*$ |
| | $a$ |

iv) Construct the predictive parsing table

| | $\$$ | $+$ | $*$ | $a$ |
|---|---|---|---|---|
| S | | $+SS$ | $*SS$ | $a$ |

The grammar is LL(1). Since It Contains no more than 1 production

v) i/p string : $+*aaa$

| Steck | input | Action |
|---|---|---|
| $S\$$ | $+*aaa\$$ | |
| $+SS\$$ | $+*aaa\$$ | $S \rightarrow +SS$ |
| $SS\$$ | $*aaa\$$ | match '+' |
| $*SSS\$$ | $*aaa\$$ | $S \rightarrow *SS$ |
| $SSS\$$ | $aaa\$$ | match $*$ |
| $aSS\$$ | $aaa\$$ | push $S \rightarrow a$  i/p Is Successfully |
| $SS\$$ | $aa\$$ | match a     parsed |
| $aS\$$ | $aa\$$ | push $S \rightarrow a$ |
| $S\$$ | $a\$$ | match a |
| $a\$$ | $a\$$ | push $S \rightarrow a$ |
| $\$$ | $\$$ | match a |

6) $S \rightarrow S(S)S \mid \varepsilon$    i/p: $((X))$

i) Remove left Recursion

$S \rightarrow \varepsilon S'$

$S' \rightarrow (S)SS' \mid \varepsilon$

ii) Remove left factoring

→ not needed

7) $S \rightarrow S+S \mid SS \mid (S) \mid S* \mid a$

   i/p: $(a+a)*a$

i) Remove left Recursion

$S \rightarrow aS' \mid (S) S'$

$S' \rightarrow +SS' \mid SS' \mid *S' \mid \varepsilon$

iii) Find FIRST and FOLLOW

|  | S | S' |
|---|---|---|
| FIRST | a<br>( | +, ε<br>a<br>(<br>* |
| FOLLOW | \$   +<br>)   a<br>( <br>* | \$<br>)<br>+<br>a, C, * |

iv) Construct the predictive parsing table

|  | \$ | ( | a | ) | + | * |
|---|---|---|---|---|---|---|
| S |  | $S \rightarrow (S)S'$ | $S \rightarrow aS'$ |  |  |  |
| S' | $S' \rightarrow \varepsilon$ | $S' \rightarrow SS'$<br>$S' \rightarrow \varepsilon$ | $S \rightarrow SS'$<br>$S' \rightarrow \varepsilon$ | $S' \rightarrow \varepsilon$ | $S' \rightarrow +SS'$<br>$S' \rightarrow \varepsilon$ | $S' \rightarrow *S'$<br>$S' \rightarrow \varepsilon$ |

The grammar is not LL(1)

v) input : $(a+a)*a$

| stack | input | Action |
|---|---|---|
| S\$ | $(a+a)*a$ |  |
| (S)S'\$ | $(a+a)*a$ | $S \rightarrow (S)S'$ |
| S)S'\$ | $a+a)*a$ | match ( |
| aS')S'\$ | $a+a)*a$ | push $S \rightarrow aS'$ |
| S')S'\$ | $+a)*a$ | match a |

└→ Ambigous, whether to parse $S' \rightarrow +SS'$ or

$S' \rightarrow \varepsilon$

6) $S \rightarrow S(S)S \mid \varepsilon$   i/p: $(()())$

  i) Remove left Recursion
  @ $S \rightarrow \varepsilon S'$
     $S' \rightarrow (S)SS' \mid \varepsilon$

  ii) Remove left factoring
     $\rightarrow$ not needed

  iii) find FIRST and FOLLOW set

|        | S | S' |
|--------|---|----|
| FIRST  | ( <br> $\varepsilon$ | ( <br> $\varepsilon$ |
| FOLLOW | \$ <br> ) <br> ( | \$ <br> ) <br> ( |

  iv) parsing table

|    | \$ | ( | ) |
|----|----|---|---|
| S  | $S \rightarrow \varepsilon$ | $S \rightarrow d$ <br> $S \rightarrow \varepsilon$ | $S \rightarrow \varepsilon$ |
| S' | $S' \rightarrow \varepsilon$ | $S' \rightarrow (S)SS'$ <br> $S' \rightarrow \varepsilon$ | $S' \rightarrow \varepsilon$ |

$\rightarrow$ multiple production

The grammar is not LL(1), since it has a multiple transaction productions

  v) parse the i/p string: $(()())$

| stack | Input | Action |
|-------|-------|--------|
| S\$ | $(()())$\$ | |
| S'\$ | $(()())$ \$ | $S \rightarrow S'$  $\rightarrow$ ambiguity |
| (S)SS'\$ | $(()())$ \$ | $S' \rightarrow (S)SS'$ |
| S)SS'\$ | $()())$\$ | match '(' |
| S')SS'\$ | $()())$\$ | push $S \rightarrow S'$ |

7) $S \rightarrow (L) | a$
$L \rightarrow L, S | S \Rightarrow (a, a)$

Step i) Remove left Recursion
$$L \rightarrow \underset{A}{\cancel{L}}, \underset{\alpha}{S | S} \underset{\beta}{}$$

$L \rightarrow S L'$
$L' \rightarrow , S L' | \varepsilon$
$S \rightarrow (L) | a$

ii) Remove left Recursion factoring
not required

iii) write FIRST and FOLLOW set

|  | S | L | L' |
|---|---|---|---|
| FIRST | ( | ( | , |
|  | a | a | $\varepsilon$ |
| FOLLOW | \$ | ) | ) |
|  | , |  |  |
|  | ) |  |  |

iv) write the productive parsing table

| | ( | a | ) | , | $ |
|---|---|---|---|---|---|
| S | S→(L) | S→a | | | |
| L | L→SL' | L→SL' | | | |
| L' | | | L'→ε | L'→,SL' | |

v)

| Stack | Input | action |
|---|---|---|
| S$ | (a,a) $ | |
| (L)$ | (a,a) $ | S→(L) |
| L)$ | a,a) $ | match ( |
| SL')$ | a,a) $ | L→SL' |
| aL')$ | a,a)$ | S→a |
| L')$ | ,a)$ | match a |
| ,SL')$ | ,a)$ | L'→,SL' |
| SL')$ | a)$ | match , |
| aL')$ | a)$ | S→a |
| L')$ | )$ | match a |
| )$ | )$ | L'→ε |
| $ | $ | match ) |

The grammar i/p is successfully parsed

8) S→ S+S| SS| (S) |S* |a ⟹ (a+a) *a

i) Remove left recursion
S→ S+S |SS| (S) |S* | a
S→(S)S'/a s'
S' → ~~(s)s'~~ ~~+a s'~~ +SS' |SS' |* S' |ε

ii) remove left factoring
not required

iii) Find FIRST and FOLLOW set

| FIRST | S | S' |
|-------|---|-----|
| | ( | + |
| | a | a |
| | | * |
| | | ε |

| FOLLOW | S | S' |
|--------|---|-----|
| | $ | $ |
| | ) | ) |
| | + | + |
| | ( | ( |
| | a | a |
| | * | * |

iv) write predictive parse table

| | ( | a | ) | + | * | $ |
|---|---|---|---|---|---|---|
| S | S→(S)S' | S→aS' | | | | |
| S' | S'→+SS' | | | S'→+SS' | | |
| | S→SS' | | | | | |
| | S→*S' | | | | | |
| | S→ε | | | | | |

∴ The given grammar is not LL(1)

g) S→aSbS |bSaS| ε  ⇒ aabbab

i) Remove left factoring — not required
ii) Remove left Recursion —not required
iii) write FIRST and FOLLOW set

| FIRST | S |
|-------|---|
| | a |
| | b |
| | ε |

| FOLLOW | S |
|--------|---|
| | b |
| | a |
| | $ |

iv) write a predictive parsing table

| | a | b | $ |
|---|---|---|---|
| S | S→aSbS | S→bSaS | S→ε |
| | S→ε | S→ε | |

V.

| Stack | Input | Action |
|---|---|---|
| S$ | aabbab$ | |
| asbs$ | aabbab$ | s→asbs |
| sbs$ | abbab$ | match a |
| asbsbs$ | abbab$ | s→asbs |
| sbsbs$ | bbab$ | match a |
| bsasbs.bs$ | bbab$ | s→bsas |
| Sasbsbs$ | bab$ | match b |
| bSa.Sa.sbsbs$ | bab$ | s→bsas |
| Sasasbsbs$ | | |

  └→ ambigous

10) bexpn → bexpn on bterm / bterm
    bterm → bterm and bfactor / bfactor
    bfactor → not bfactor | (bexpn) / true / false
    i/p: not (true on false)

→ i) Remove left Recursion
    bexpn → ~~bterm~~ bexpnon bterm / bterm
    bexpn' → ~~on~~ bterm bexpn' / ~~ε~~ → bterm→bterm and bfactor/bfactor
    bterm → ~~bfactor bterm~~ bfactor bterm'
    bterm' → and bfactor bterm' / ε
    bfactor → not bfactor | (bexpn) / true / false

ii) No left recursion factoring

iii) FIND the FIRST and FOLLOW set

|        | bexpr | bexpr' | bterm | bterm' | bfactor |
|--------|-------|--------|-------|--------|---------|
| FIRST  | not ( true false | or ε | not ( true false | and ε | not ( true false |
| FOLLOW | $ ) | $ ) | or ) $ | or $ ) | and or ) $ |

predictive parsing table:

bexpr → bterm bexpr' | bterm bexpr'

|         | (              | not           | )  | or            | and            | true          | false          | $ |
|---------|----------------|---------------|----|---------------|----------------|---------------|----------------|---|
| bexpr   | bterm bexpr'   | bterm bexpr'  |    |               |                | bterm bexpr'  | bterm bexpr'   |   |
| bexpr'  |                |               | ε  | or bterm bexpr' |              |               |                | ε |
| bterm   | bfactor bterm' | bfactor bterm'|    |               | ●              | bfactor bterm'| bfactor bterm' |   |
| bterm'  |                |               | ε  | ε             | and bfactor bterm' |           |                | ε |
| bfactor | (bexpr)        | not bfactor   |    |               |                | true          | false          |   |

| Stack | Input | Action |
|---|---|---|
| bexpr $ | not (true or false) $ | bexpr → btermbexpr' |
| btermbexpr' $ | not (true or false) $ | bterm → bfactor bterm' |
| bfactorbterm'bexpr'$ | not (true or false) $ | bfactor → not bfactor |
| not bfactorbterm' | not (true or false) $ | bfactor → (bexpr) |
| bexpr' $ | | |
| bfactorbterm' | (true or false) $ | match ( |
| bexpr' $ | | |
| bexpr) bterm' | true or false) $ | match ( |
| bexpr' $ | | |
| btermbexpr') bterm' | true or false) $ | bexpr → btermbexpr' |
| bexpr' $ | | |
| bfactorbterm' bexpr') | true or false) $ | bterm → bfactorbterm' |
| bterm' bexpr' $ | | |
| bterm' bexpr') | or false) $ | bfactor → bterm' |
| bterm' bexpr' $ | | |
| bexpr') bterm' | or false) $ | bterm' → ε |
| bexpr' $ | | |
| bterm bexpr') | false ) $ | bexpr' → btermbexpr |
| bterm' bexpr' $ | | |
| bfactorbterm' | false ) $ | bterm → bfactorbterm' |
| bexpr) bterm'bexpr'$ | | |
| | | bfactor → ε |
| bterm' bexpr) bterm' | ) $ | |
| bexpr' $ | | bterm' → ε & bexpr' → |
| bterm' bexpr'$ | $ | & match) |
| $ | $ | |

Error recovery in predictive parser:

1. panic mode Recovery: In the blank entries of the follow set of all NT, place "synch"

| Stack | i/p | Table Entry | Action |
|-------|-----|-------------|--------|
| NT | T | blank | skip the terminal from the i/p |
| NT | T | synch | Remove NT from the stack except start symbol |
| T | T | match | pop the terminal from stack & input |

Example:

1. $E \rightarrow E + T | T$         $E \rightarrow T E'$
   $T \rightarrow T * F / F$          $E' \rightarrow + T E' | \varepsilon$
   $F \rightarrow (E) | id$           $T \rightarrow F T'$
   i/p: $) id * + id$                 $T' \rightarrow * F T' | \varepsilon$
                                      $F \rightarrow (E) | id$

|        | E | E' | T | T' | F |
|--------|---|-----|---|-----|---|
| FIRST  | ( | + | ( | * | ( |
|        | id | $\varepsilon$ | id | $\varepsilon$ | id |
| Follow | $ | $ | + | + | * |
|        | ) | ) | $ | $ | $ |
|        |   |   | ) | ) | ) |
|        |   |   |   |   | + |

predictive parsing table

|    | ( | id | ) | + | * | $ |
|----|---|-----|---|---|---|---|
| E  | $E \rightarrow TE'$ | $E \rightarrow TE'$ | synch | | | |
| E' | | | $E' \rightarrow \varepsilon$ | $E' \rightarrow + TE'$ | | $E' \rightarrow \varepsilon$ |
| T  | $T \rightarrow FT'$ | $T \rightarrow FT'$ | synch | synch | | synch |
| T' | | | $T' \rightarrow \varepsilon$ | $T' \rightarrow \varepsilon$ | $T' \rightarrow * FT'$ | $T' \rightarrow \varepsilon$ |
| F  | $F \rightarrow (E)$ | $F \rightarrow id$ | synch | synch | synch | synch |

| stack | input | Action |
|---|---|---|
| E $ | )id * +id $ | Since It Is the start Symbol Skip the i/p [E, )] = Synch |
| E $ | id * +id $ | E → TE' |
| TE' $ | id* + id $ | T → FT' |
| FT'E' $ | id* +id $ | F → id |
| id T'E' $ | id * +id $ | match id |
| T'E' $ | * +id $ | T' → *FT' |
| *FT'E' $ | * +id $ | match * |
| FT'E' $ | +id $ | [F, +] = Synch, remove NT from Stack |
| T'E' $ | +id $ | T' → ε |
| E' $ | +id $ | E' → +TE' |
| +TE' $ | +id$ | match + |
| TE' $ | id $ | T → FT' |
| FT'E' $ | id $ | F → id |
| id T'E' $ | id $ | match id |
| T'E' $ | $ | T' → ε |
| E' $ | $ | E' → ε |
| $ | $ | match $ |

8) Show that the following grammar is ambigous.
E → E+E | E-E | E*E | E/E | (E) | id    and i/p: id + id * id
Give an unambigous grammar such that precedence
order from lowest on highest are +, -, *, /, ( ), id
and all are left to right associative.

## Bottom up parses:

### Introduction:

→ A bottom-up parse Corresponds to the Construction of a parse tree for an input string beginning at the leaves (the bottom) and working up towards the root (the top)

→ eg: A bottom-up parses for id * id



### Handle pruning:

Bottom up parsing during a left to right Scan of the input Constructs a right most derivation in reverse. Informally a 'handle' is a substring that matches the body of the production, and whose reduction represents one step along the reverse of the right most derivation

### Example:

adding subscripts to the tokens id for clarity, The handles during the parse of $id_1 * id_2$ a/c to the expressions are shown in the figure.

grammar → $\left\{ \begin{array}{l} E \rightarrow E+T \mid T \\ T \rightarrow T*F \mid F \\ F \rightarrow (E) \mid id \end{array} \right\}$

Although T is the body of the production E→T, the symbol T is not a handle in the sentential form T*id₂. If T were indeed replaced by E, we would get a string E*id₂, which cannot be derived from the start symbol E. Thus, the leftmost substring that matches the body of some production need not be a handle.

formally, if $S \xRightarrow{*}_{rm}$

| RIGHT SENTENTIAL FORM | HANDLE | REDUCTION PRODUCTION |
|---|---|---|
| $id_1 * id_2$ | $id_1$ | $F \to id$ |
| $F * id_2$ | $F$ | $T \to F$ |
| $T * id_2$ | $id_2$ | $F \to id$ |
| $F$ | $T * F$ | $E \to T * F$ |

Handles during a parse of $id_1 * id_2$ (a)

Formally, If $S \xRightarrow{*}_{rm} \alpha A w \xRightarrow{}_{rm} \alpha \beta w$, as in figure (b), then production A→β, in the position following α is a handle of αβw. Alternatively, a handle of a right sentential form γ is a production A→β and a position of γ where the string β may be found, such that replacing β at that position by A produces the previous right sentential form in a rightmost derivation of γ.

Notice that the string w to the right of the handle must contain only terminal symbols. For convenience, we prefer to the body β neither than A→β as a handle. Note we "a handle" neither than "the handle", because

the grammar could be ambigous, with moore than one rightmost derivation of $\alpha \beta w$. If a grammar is unambigous, then every right-sentenhal form of the grammar has exactly one handle.

A right-most derivation in reverse can be obtained by "handle prunning". That is, we start with a string of terminals $w$ to be parsed. If $w$ is a sentence of a grammar at hand, then let $w = \gamma_n$, where $\gamma_n$ is the $n^{th}$ right sentenhal form of some as yet unknown rightmost derivation

$$S = \gamma_0 \underset{rm}{\Rightarrow} \gamma_1 \underset{rm}{\Rightarrow} \gamma_2 - \cdots \underset{rm}{\Rightarrow} \gamma_{n-1} \underset{rm}{\Rightarrow} \gamma_n = w$$



figure (b). A handle $A \rightarrow \beta$ in the parse tree for $\alpha \beta w$

To reconstruct this derivation in reverse order, we locate the handle $\beta_n$ in $\gamma_n$ and replace $\beta_n$ by the head of relevant production $A_n \rightarrow \beta_n$ to obtain the previous right sentential form $\gamma_{n-1}$. Note that we do not know how handles are to be found, but we shall see methods of doing so shortly

We the repeat this process. That is we locate the handle $\beta_{n-1}$ in $\gamma_{n-1}$ and reduce this handle to obtain the right sentenhal form $\gamma_{n-2}$. If by continuing this process we produce a right sentential form consisting only of the start symbol $S$. the we halt & that's successfull.

## Shift-Reduce parsing:

Shift-reduce parsing is a form of bottom up parsing in which a stack holds grammar symbols and an input buffer holds the rest of the string to be parsed.

→ we use $ to mark the bottom of the stack and also the right end of the input. Conventionally, when discussing bottom-up parsing, we show the top of the stack on the right, rather on the left as we did for top-down parsing.

Initially, the stack is empty, and string ω is on the input, as follows:

| STACK | INPUT |
|-------|-------|
| $ | ω$ |

During left to right scan of the input string, The parses. shifts zero/more input symbols onto the stack untill it is ready to reduce a string β of the grammar symbols on top of the stack. It then reduces β to the head of the appropriate production. The parses repeats this cycle untill it has detected an error or untill the stack contains the start symbol and input is empty.

## Actions in shift-reduce parsing:

1. Shift: shift the next-input symbol onto the top of stack

2. Reduce: The right end of the string to be reduced must be at the top of the stack. Locate the left end of the string within the stack and decide with what non terminal to replace the string

3. **Accept:** Announce successfull completion of parsing

4. **Error:** Discover a syntax error and can an error recovery routine

Find the handles for the given RSF and Construct shift-reduce parser:

1. $E \to E+T|T$
   $T \to T*F|F$
   $F \to (E)|id$

   ilp: id + id
   id + id * id

$\to$ RMD
$E \Rightarrow E+T$
$\Rightarrow E+F$
$\Rightarrow E+id$
$\Rightarrow T+id$
$\Rightarrow F+id$
$\Rightarrow id+id$

| RSF | Handle | Action |
|---|---|---|
| $id_1 + id_2$ | $id_1$ | $F \to id$ |
| $F + id_2$ | $F$ | $T \to F$ |
| $T + id_2$ | $T$ | $E \to T$ |
| $E + id_2$ | $id_2$ | $F \to id$ |
| $E + F$ | $F$ | $T - F$ |
| $E + T$ | $E + T$ | $E \to E+T$ |
| $E$ | | |

| Stack | RSF | Action |
|---|---|---|
| $\$$ | $id_1 + id_2 \$$ | shift $id_1$ |
| $\$ id_1$ | $+ id_2 \$$ | reduce $F \to id$ |
| $\$ F$ | $+ id_2 \$$ | reduce $T \to F$ |
| $\$ T$ | $+ id_2 \$$ | reduce $E \to T$ |
| $\$ E$ | $+ id_2 \$$ | shift $+$ |
| $\$ E +$ | $+ id_2 \$$ | shift $id_2$ |
| $\$ E + id_2$ | $\$$ | reduce $F \to id_2$ |
| $\$ E + F$ | $\$$ | reduce $T \to F$ |
| $\$ E + T$ | $\$$ | reduce $E \to T$ |
| $\$ E$ | | Success |

i/p: id+id*id

RMD:

$E \rightarrow E+T$

$\Rightarrow E+T*F$

$\Rightarrow E+T*id$

$\Rightarrow E+F*id$

$\Rightarrow E+id*id$

$\Rightarrow T+id*id$

$\Rightarrow F+id*id$

$\Rightarrow id+id*id$

| RSF | Handle | Action |
|---|---|---|
| $id_1+id_2*id_3$ | $id_1$ | $F \rightarrow id$ |
| $F+id_2*id_3$ | $F$ | $T \rightarrow F$ |
| $T+id_2*id_3$ | $T$ | $E \rightarrow F$ |
| $E+id_2*id_3$ | $id_2$ | $F \rightarrow id_2$ |
| $F+F*id_3$ | $F$ | $T \rightarrow F$ |
| $E+T*id_3$ | $id_3$ | $F \rightarrow id_3$ |
| $E+T*F$ | $T*F$ | $T \rightarrow F*F$ |
| $E+T$ | $E+T$ | $E \rightarrow E+T$ |
| $E$ | | |

| Stack | RSF | Action |
|---|---|---|
| $\$$ | $id_1+id_2*id_3\$$ | shift $id_1$ |
| $\$id_1$ | $+id_2*id_3\$$ | $F \rightarrow id_1$ |
| $\$F$ | $+id_2*id_3\$$ | $T \rightarrow F$ |
| $\$T$ | $+id_2*id_3\$$ | $E \rightarrow T$ |
| $\$E$ | $+id_2*id_3\$$ | Shift $+$ |
| $\$E+$ | $id_2*id_3\$$ | shift $id_2$ |
| $\$E+id_2$ | $*id_3\$$ | reduce $F \rightarrow id$ |
| | | $T \rightarrow F$ |
| | | shift $*$ |
| | | shift $id_3$ |
| $\$E+T*id_3$ | $\$$ | reduce $F \rightarrow id_3$ |
| $\$E+T*F$ | $\$$ | reduce $T \rightarrow T*F$ |
| | | reduce $F \rightarrow E+T$ |
| $\$E$ | $\$$ | Success |

2) $S \to 0S1 \mid 01$   i/p: 000111

$S \xrightarrow{rm} 0S1$

$\Rightarrow 00S11b$

$\Rightarrow 000111$

| RSF | Handle | Action |
|---|---|---|
| 000111 | 01 | S→0S1 |
| 00S11 | 0S1 | S→0S1 |
| 0S1 | 0S1 | S→0S1 |

| Stack | RSF | Action |
|---|---|---|
| $ | 000111$ | Shift 0 |
| $0 | 00111$ | shift 0 |
| $00 | 0111$ | shift 0 |
| $000 | 111$ | shift 1 |
| $0001 | 11$ | reduce S→01 |
| $00S | 11$ | shift 1 |
| $00S1 | 1$ | reduce S→0S1 |
| $0S | 1$ | shift 1 |
| $0S1 | $ | reduce S→0S1 |
| $S | $ | success |

3) $S \to SS+ \mid SS* \mid a$   i/p: aaa*a++

$S \to SS+$

$\Rightarrow SSS++$

$\Rightarrow SSa++$

$\Rightarrow SSS*a++$

$\Rightarrow SSa*a++$

$\Rightarrow Saa*a++$

$\Rightarrow aaa*a++$

| RSF | Handle | Action |
|---|---|---|
| aaa*a++ | a | S→a |
| Saa*a++ | a | S→a |
| SSa*a++ | a | S→a |
| SSS*a++ | SS* | S→SS* |
| SSa++ | a | S→a |
| SSS++ | SS+ | S→SS+ |
| SS+* | SS+ | S→SS+ |
| S | | |

| Stack | RSF | Action |
|---|---|---|
| $ | aaa*a++$ | shift a |
| $a | aa*a++$ | reduce s→a |
| $s | aa*a++$ | Shift a |
| $Sa | a*a++$ | reduce s→a |
| $SS | a*a++$ | Shift a |
| $SSa | *a++$ | reduce s→a |
| $SSS | *a++$ | Shift * |
| $SS* | a++$ | reduce s→ss* |
| $SS | a++$ | shift a |
| $SSa | a++$ | reduce s→a |
| $SSa | s++$ | shift + |
| $SS+ | +$ | reduce s→ss+ |
| $SS | +$ | Shift + |
| $SS+ | $ | reduce s→ss+ |
| $S | $ | Success |

Types of conflicts in shift-reduce parsers

Conflicts during shifts-Reduce parsing.

There are CFG's for which shift-reduce parsing cannot be used. Every shift-reduce parser for such a grammar can reach a configuration in which the parser. Knowing the entire stack contents & the next input symbol

Types:

i) shift/reduce conflict:
→ Cannot decide whether to shift on to reduce called shift-reduce conflict

ii) reduce/reduce conflict:

→ Cannot decide whether which of several reductions to make called reduce-reduce conflicts

eg: Consider the grammar
$$E \longrightarrow E+E$$
$$| E-E \qquad i/p: 2+3*4$$
$$| NUM$$
$$| ID$$

| No. | Stack | operation/grammar |
|-----|-------|-------------------|
| 1 | 2 NUM | shift 2 |
| 2 | E | reduce E→NUM |
| 3 | E+ | shift + |
| 4 | E+3 | shift 3 |
| 5 | E+E | reduce E→NUM |
| 6 | E | reduce E→E+E on shift * ie Shift reduce conflict |

$\frac{58}{11}$  $E \to T$

$T \to id$

$F \to id$

$idt.$

eg 2) An ambigous grammar can never be LR. for eg:

Consider the dangling-else grammar

$$stmt \to if\ expr\ then\ stmt$$
$$| \quad If\ expr\ then\ stmt\ else\ stmt$$
$$| \quad other$$

If we have a Shift-reduce parser in Configuration

| STACK | INPUT |
|---|---|
| ...If expr then stmt | else... $ |

→ we cannot tell whether If expr then stmt is the handle, no matter what happens below it on the stack. here There is a shift-/reduce Conflict. Depending on what follows the __else__ on the input, it might be consuct to reduce __If expr then stmt__ to __stmt__, or it might be consuct to shift __else__ then to look for another __stmt__ to Complete the alternative __If expr Then stmt else stmt__

eg 3)  (1)   stmt $\to$ id (parameter_list)

(2)   stmt $\to$ expr: = expr

(3)   parameter_list $\to$ parameter_list, parameter

(4)   parameter_list $\to$ parameter

(5)   parameter $\to$ id

(6)   expr $\to$ id (expr_list)

(7)   expr $\to$ id

(8)   expr_list $\to$ expr_list, expr

(9)   expr_list $\to$ expr

```
        STACK                    INPUT
        ...id (id                   ,id)...
```

In the problem id on the top of the stack must be
reduced, but by which production? The correct
choice is production (5) if P is a procedure, but
production (7) if P is an array. The stack does not
tell which information in the symbol table obtained
from the declaration of P must be used

So, In this case we have the conflict that reduce id
by parameter or exper. It is called as

reduce-reduce conflict.

or in the configuration above. In the former case, we choose reduction by production (5); in the latter case by production (7). Notice how the symbol third from the top of the stack determines the reduction to be made, even though it is not involved in the reduction. Shift-reduce parsing can utilize information far down in the stack to guide the parse.                                      □

## 4.6 OPERATOR-PRECEDENCE PARSING

The largest class of grammars for which shift-reduce parsers can be built successfully – the LR grammars – will be discussed in Section 4.7. However, for a small but important class of grammars we can easily construct efficient shift-reduce parsers by hand. These grammars have the property (among other essential requirements) that no production right side is $\epsilon$ or has two adjacent nonterminals. A grammar with the latter property is called an *operator grammar*.

**Example 4.27.** The following grammar for expressions

$$E \rightarrow EAE \mid (E) \mid -E \mid \text{id}$$
$$A \rightarrow + \mid - \mid * \mid / \mid \uparrow$$

is not an operator grammar, because the right side *EAE* has two (in fact three) consecutive nonterminals. However, if we substitute for *A* each of its alternatives, we obtain the following operator grammar:

$$E \rightarrow E+E \mid E-E \mid E*E \mid E/E \mid E \uparrow E \mid (E) \mid -E \mid \text{id} \qquad (4.17)$$

We now describe an easy-to-implement parsing technique called operator-precedence parsing. Historically, the technique was first described as a manipulation on tokens without any reference to an underlying grammar. In fact, once we finish building an operator-precedence parser from a grammar, we may effectively ignore the grammar, using the nonterminals on the stack only as placeholders for attributes associated with the nonterminals.

As a general parsing technique, operator-precedence parsing has a number of disadvantages. For example, it is hard to handle tokens like the minus sign, which has two different precedences (depending on whether it is unary or binary). Worse, since the relationship between a grammar for the language being parsed and the operator-precedence parser itself is tenuous, one cannot always be sure the parser accepts exactly the desired language. Finally, only a small class of grammars can be parsed using operator-precedence techniques.

Nevertheless, because of its simplicity, numerous compilers using operator-precedence parsing techniques for expressions have been built successfully. Often these parsers use recursive descent, described in Section 4.4, for statements and higher-level constructs. Operator-precedence parsers have even been built for entire languages.

In operator-precedence parsing, we define three disjoint *precedence relations*, $<\cdot$, $\doteq$, and $\cdot>$, between certain pairs of terminals. These precedence relations guide the selection of handles and have the following meanings:

| RELATION | MEANING |
|----------|---------|
| $a \lessdot b$ | $a$ "yields precedence to" $b$ |
| $a \doteq b$ | $a$ "has the same precedence as" $b$ |
| $a \gtrdot b$ | $a$ "takes precedence over" $b$ |

We should caution the reader that while these relations may appear similar to the arithmetic relations "less than," "equal to," and "greater than," the precedence relations have quite different properties. For example, we could have $a \lessdot b$ and $a \gtrdot b$ for the same language, or we might have none of $a \lessdot b$, $a \doteq b$, and $a \gtrdot b$ holding for some terminals $a$ and $b$.

There are two common ways of determining what precedence relations should hold between a pair of terminals. The first method we discuss is intuitive and is based on the traditional notions of associativity and precedence of operators. For example, if $*$ is to have higher precedence than $+$, we make $+ \lessdot *$ and $* \gtrdot +$. This approach will be seen to resolve the ambiguities of grammar (4.17), and it enables us to write an operator-precedence parser for it (although the unary minus sign causes problems).

The second method of selecting operator-precedence relations is first to construct an unambiguous grammar for the language, a grammar that reflects the correct associativity and precedence in its parse trees. This job is not difficult for expressions; the syntax of expressions in Section 2.2 provides the paradigm. For the other common source of ambiguity, the dangling else, grammar (4.9) is a useful model. Having obtained an unambiguous grammar, there is a mechanical method for constructing operator-precedence relations from it. These relations may not be disjoint, and they may parse a language other than that generated by the grammar, but with the standard sorts of arithmetic expressions, few problems are encountered in practice. We shall not discuss this construction here; see Aho and Ullman [1972b].

## Using Operator-Precedence Relations

The intention of the precedence relations is to delimit the handle of a right-sentential form, with $\lessdot$ marking the left end, $\doteq$ appearing in the interior of the handle, and $\gtrdot$ marking the right end. To be more precise, suppose we have a right-sentential form of an operator grammar. The fact that no adjacent nonterminals appear on the right sides of productions implies that no right-sentential form will have two adjacent nonterminals either. Thus, we may write the right-sentential form as $\beta_0 a_1 \beta_1 \cdots a_n \beta_n$, where each $\beta_i$ is either $\epsilon$ (the empty string) or a single nonterminal, and each $a_i$ is a single terminal.

Suppose that between $a_i$ and $a_{i+1}$ exactly one of the relations $\lessdot$, $\doteq$, and $\gtrdot$ holds. Further, let us use $\$$ to mark each end of the string, and define $\$ \lessdot b$ and $b \gtrdot \$$ for all terminals $b$. Now suppose we remove the nonterminals from the string and place the correct relation $\lessdot$, $\doteq$, or $\gtrdot$, between each

pair of terminals and between the endmost terminals and the $'s marking the ends of the string. For example, suppose we initially have the right-sentential form **id + id \* id** and the precedence relations are those given in Fig. 4.23. These relations are some of those that we would choose to parse according to grammar (4.17).

|      | id  | +   | *   | $   |
|------|-----|-----|-----|-----|
| id   |     | ·>  | ·>  | ·>  |
| +    | <·  | ·>  | <·  | ·>  |
| *    | <·  | ·>  | ·>  | ·>  |
| $    | <·  | <·  | <·  |     |

**Fig. 4.23.** Operator-precedence relations.

Then the string with the precedence relations inserted is:

$$\$ <\cdot \text{ id } \cdot> + <\cdot \text{ id } \cdot> * <\cdot \text{ id } \cdot> \$ \qquad (4.18)$$

For example, $<\cdot$ is inserted between the leftmost $ and **id** since $<\cdot$ is the entry in row $ and column **id**. The handle can be found by the following process.

1. Scan the string from the left end until the first $\cdot>$ is encountered. In (4.18) above, this occurs between the first **id** and $+$.

2. Then scan backwards (to the left) over any $\doteq$'s until a $<\cdot$ is encountered. In (4.18), we scan backwards to $.

3. The handle contains everything to the left of the first $\cdot>$ and to the right of the $<\cdot$ encountered in step (2), including any intervening or surrounding nonterminals. (The inclusion of surrounding nonterminals is necessary so that two adjacent nonterminals do not appear in a right-sentential form.) In (4.18), the handle is the first **id**.

If we are dealing with grammar (4.17), we then reduce **id** to $E$. At this point we have the right-sentential form $E + \text{id} * \text{id}$. After reducing the two remaining **id**'s to $E$ by the same steps, we obtain the right-sentential form $E + E * E$. Consider now the string $\$ + * \$$ obtained by deleting the nonterminals. Inserting the precedence relations, we get

$$\$ <\cdot + <\cdot * \cdot> \$$$

indicating that the left end of the handle lies between $+$ and $*$ and the right end between $*$ and $. These precedence relations indicate that, in the right-sentential form $E + E * E$, the handle is $E * E$. Note how the $E$'s surrounding the $*$ become part of the handle.

Since the nonterminals do not influence the parse, we need not worry about distinguishing among them. A single marker "nonterminal" can be kept on

the stack of a shift-reduce parser to indicate placeholders for attribute values.

It may appear from the discussion above that the entire right-sentential form must be scanned at each step to find the handle. Such is not the case if we use a stack to store the input symbols already seen and if the precedence relations are used to guide the actions of a shift-reduce parser. If the precedence relation $<\cdot$ or $\doteq$ holds between the topmost terminal symbol on the stack and the next input symbol, the parser shifts; it has not yet found the right end of the handle. If the relation $\cdot>$ holds, a reduction is called for. At this point the parser has found the right end of the handle, and the precedence relations can be used to find the left end of the handle in the stack.

If no precedence relation holds between a pair of terminals (indicated by a blank entry in Fig. 4.23), then a syntactic error has been detected and an error recovery routine must be invoked, as discussed later in this section. The above ideas can be formalized by the following algorithm.

**Algorithm 4.5.** Operator-precedence parsing algorithm.

*Input.* An input string $w$ and a table of precedence relations.

*Output.* If $w$ is well formed, a *skeletal* parse tree, with a placeholder nonterminal $E$ labeling all interior nodes; otherwise, an error indication.

*Method.* Initially, the stack contains $ and the input buffer the string $w$. To parse, we execute the program of Fig. 4.24.      □

```
(1)  set ip to point to the first symbol of w$;
(2)  repeat forever
(3)      if $ is on top of the stack and ip points to $ then
(4)          return
     else begin
(5)          let a be the topmost terminal symbol on the stack
                 and let b be the symbol pointed to by ip;
(6)          if a <· b or a ≐ b then begin
(7)              push b onto the stack;
(8)              advance ip to the next input symbol;
             end;
(9)          else if a ·> b then            /* reduce */
(10)             repeat
(11)                 pop the stack
(12)             until the top stack terminal is related by <·
                     to the terminal most recently popped
(13)         else error()
     end
```

**Fig. 4.24.** Operator-precedence parsing algorithm.

### Operator-Precedence Relations from Associativity and Precedence

We are always free to create operator-precedence relations any way we see fit and hope that the operator-precedence parsing algorithm will work correctly when guided by them. For a language of arithmetic expressions such as that generated by grammar (4.17) we can use the following heuristic to produce a proper set of precedence relations. Note that grammar (4.17) is ambiguous, and right-sentential forms could have many handles. Our rules are designed to select the "proper" handles to reflect a given set of associativity and precedence rules for binary operators.

1.  If operator $\theta_1$ has higher precedence than operator $\theta_2$, make $\theta_1 \cdot> \theta_2$ and $\theta_2 <\cdot \theta_1$. For example, if $*$ has higher precedence than $+$, make $* \cdot> +$ and $+ <\cdot *$. These relations ensure that, in an expression of the form $E + E * E + E$, the central $E * E$ is the handle that will be reduced first.

2.  If $\theta_1$ and $\theta_2$ are operators of equal precedence (they may in fact be the same operator), then make $\theta_1 \cdot> \theta_2$ and $\theta_2 \cdot> \theta_1$ if the operators are left-associative, or make $\theta_1 <\cdot \theta_2$ and $\theta_2 <\cdot \theta_1$ if they are right-associative. For example, if $+$ and $-$ are left-associative, then make $+ \cdot> +$, $+ \cdot> -$, $- \cdot> -$, and $- \cdot> +$. If $\uparrow$ is right associative, then make $\uparrow <\cdot \uparrow$. These relations ensure that $E - E + E$ will have handle $E - E$ selected and $E \uparrow E \uparrow E$ will have the last $E \uparrow E$ selected.

3.  Make $\theta <\cdot$ **id**, **id** $\cdot> \theta$, $\theta <\cdot ($, $( <\cdot \theta$, $) \cdot> \theta$, $\theta \cdot> )$, $\theta \cdot> \$$, and $\$ <\cdot \theta$ for all operators $\theta$. Also, let

|            |              |               |
|------------|--------------|---------------|
| $( \doteq )$  | $\$ <\cdot ($   | $\$ <\cdot$ **id** |
| $( <\cdot ($  | **id** $\cdot> \$$  | $) \cdot> \$$    |
| $( <\cdot$ **id** | **id** $\cdot> )$   | $) \cdot> )$     |

These rules ensure that both **id** and $(E)$ will be reduced to $E$. Also, $\$$ serves as both the left and right endmarker, causing handles to be found between $\$$'s wherever possible.

**Example 4.28.** Figure 4.25 contains the operator-precedence relations for grammar (4.17), assuming

1.  $\uparrow$ is of highest precedence and right-associative,

2.  $*$ and $/$ are of next highest precedence and left-associative, and

3.  $+$ and $-$ are of lowest precedence and left-associative,

(Blanks denote error entries.) The reader should try out the table to see that it works correctly, ignoring problems with unary minus for the moment. Try the table on the input **id** $* ($ **id** $\uparrow$ **id** $) -$ **id** $/$ **id**, for example.                                   □

| | + | – | * | / | ↑ | id | ( | ) | $ |
|---|---|---|---|---|---|---|---|---|---|
| + | ·> | ·> | <· | <· | <· | <· | <· | ·> | ·> |
| – | ·> | ·> | <· | <· | <· | <· | <· | ·> | ·> |
| * | ·> | ·> | ·> | ·> | <· | <· | <· | ·> | ·> |
| / | ·> | ·> | ·> | ·> | <· | <· | <· | ·> | ·> |
| ↑ | ·> | ·> | ·> | ·> | <· | <· | <· | ·> | ·> |
| id | ·> | ·> | ·> | ·> | ·> | | | ·> | ·> |
| ( | <· | <· | <· | <· | <· | <· | <· | ≐ | |
| ) | ·> | ·> | ·> | ·> | ·> | | | ·> | ·> |
| $ | <· | <· | <· | <· | <· | <· | <· | | |

**Fig. 4.25.** Operator-precedence relations.

## Handling Unary Operators

If we have a unary operator such as ¬ (logical negation), which is not also a binary operator, we can incorporate it into the above scheme for creating operator-precedence relations. Supposing ¬ to be a unary prefix operator, we make $\theta <\cdot \neg$ for any operator $\theta$, whether unary or binary. We make $\neg \cdot> \theta$ if ¬ has higher precedence than $\theta$ and $\neg <\cdot \theta$ if not. For example, if ¬ has higher precedence than &, and & is left-associative, we would group $E \& \neg E \& E$ as $(E \& (\neg E)) \& E$, by these rules. The rule for unary postfix operators is analogous.

The situation changes when we have an operator like the minus sign – that is both unary prefix and binary infix. Even if we give unary and binary minus the same precedence, the table of Fig. 4.25 will fail to parse strings like **id**∗–**id** correctly. The best approach in this case is to use the lexical analyzer to distinguish between unary and binary minus, by having it return a different token when it sees unary minus. Unfortunately, the lexical analyzer cannot use lookahead to distinguish the two; it must remember the previous token. In Fortran, for example, a minus sign is unary if the previous token was an operator, a left parenthesis, a comma, or an assignment symbol.

## Precedence Functions

Compilers using operator-precedence parsers need not store the table of precedence relations. In most cases, the table can be encoded by two *precedence functions* $f$ and $g$ that map terminal symbols to integers. We attempt to select $f$ and $g$ so that, for symbols $a$ and $b$,

1. $f(a) < g(b)$ whenever $a <\cdot b$,
2. $f(a) = g(b)$ whenever $a \doteq b$, and
3. $f(a) > g(b)$ whenever $a \cdot> b$.

Thus the precedence relation between $a$ and $b$ can be determined by a

numerical comparison between $f(a)$ and $g(b)$. Note, however, that error entries in the precedence matrix are obscured, since one of (1), (2), or (3) holds no matter what $f(a)$ and $g(b)$ are. The loss of error detection capability is generally not considered serious enough to prevent the using of precedence functions where possible; errors can still be caught when a reduction is called for and no handle can be found.

Not every table of precedence relations has precedence functions to encode it, but in practical cases the functions usually exist.

**Example 4.29.** The precedence table of Fig. 4.25 has the following pair of precedence functions,

|   | + | − | * | / | ↑ | ( | ) | id | $ |
|---|---|---|---|---|---|---|---|----|---|
| $f$ | 2 | 2 | 4 | 4 | 4 | 0 | 6 | 6 | 0 |
| $g$ | 1 | 1 | 3 | 3 | 5 | 5 | 0 | 5 | 0 |

For example, $* <\cdot$ id, and $f(*) < g(\text{id})$. Note that $f(\text{id}) > g(\text{id})$ suggests that id $\cdot>$ id; but, in fact, no precedence relation holds between id and id. Other error entries in Fig. 4.25 are similarly replaced by one or another precedence relation.                                                                      □

A simple method for finding precedence functions for a table, if such functions exist, is the following.

**Algorithm 4.6.** Constructing precedence functions.

*Input.* An operator precedence matrix.

*Output.* Precedence functions representing the input matrix, or an indication that none exist.

*Method.*

1.  Create symbols $f_a$ and $g_a$ for each $a$ that is a terminal or $. 

2.  Partition the created symbols into as many groups as possible, in such a way that if $a \doteq b$, then $f_a$ and $g_b$ are in the same group. Note that we may have to put symbols in the same group even if they are not related by $\doteq$. For example, if $a \doteq b$ and $c \doteq b$, then $f_a$ and $f_c$ must be in the same group, since they are both in the same group as $g_b$. If, in addition, $c \doteq d$, then $f_a$ and $g_d$ are in the same group even though $a \doteq d$ may not hold.

3.  Create a directed graph whose nodes are the groups found in (2). For any $a$ and $b$, if $a <\cdot b$, place an edge from the group of $g_b$ to the group of $f_a$. If $a \cdot> b$, place an edge from the group of $f_a$ to that of $g_b$. Note that an edge or path from $f_a$ to $g_b$ means that $f(a)$ must exceed $g(b)$; a path from $g_b$ to $f_a$ means that $g(b)$ must exceed $f(a)$.

4.  If the graph constructed in (3) has a cycle, then no precedence functions exist. If there are no cycles, let $f(a)$ be the length of the longest path

beginning at the group of $f_a$; let $g(a)$ be the length of the longest path from the group of $g_a$.                                                    □

**Example 4.30.** Consider the matrix of Fig. 4.23. There are no $\doteq$ relationships, so each symbol is in a group by itself. Figure 4.26 shows the graph constructed using Algorithm 4.6.



**Fig. 4.26.** Graph representing precedence functions.

There are no cycles, so precedence functions exist. As $f_S$ and $g_S$ have no out-edges, $f(\$) = g(\$) = 0$. The longest path from $g_+$ has length 1, so $g(+) = 1$. There is a path from $g_{id}$ to $f_*$ to $g_*$ to $f_+$ to $g_+$ to $f_S$, so $g(id) = 5$. The resulting precedence functions are:

|   | + | * | id | \$ |
|---|---|---|----|----|
| $f$ | 2 | 4 | 4 | 0 |
| $g$ | 1 | 3 | 5 | 0 |

□

### Error Recovery in Operator-Precedence Parsing

There are two points in the parsing process at which an operator-precedence parser can discover syntactic errors:

1.  If no precedence relation holds between the terminal on top of the stack and the current input.[1]
2.  If a handle has been found, but there is no production with this handle as a right side.

Recall that the operator-precedence parsing algorithm (Algorithm 4.5) appears to reduce handles composed of terminals only. However, while nonterminals

---

[1] In compilers using precedence functions to represent the precedence tables, this source of error detection may be unavailable.

are treated anonymously, they still have places held for them on the parsing stack. Thus when we talk in (2) above about a handle matching a production's right side, we mean that the terminals are the same and the positions occupied by nonterminals are the same.

We should observe that, besides (1) and (2) above, there are no other points at which errors could be detected. When scanning down the stack to find the left end of the handle in steps (10-12) of Fig. 4.24, the operator-precedence parsing algorithm, we are sure to find a $<\cdot$ relation, since \$ marks the bottom of stack and is related by $<\cdot$ to any symbol that could appear immediately above it on the stack. Note also that we never allow adjacent symbols on the stack in Fig. 4.24 unless they are related by $<\cdot$ or $\doteq$. Thus steps (10-12) must succeed in making a reduction.

Just because we find a sequence of symbols $a <\cdot b_1 \doteq b_2 \doteq \cdots \doteq b_k$ on the stack, however, does not mean that $b_1 b_2 \cdots b_k$ is the string of terminal symbols on the right side of some production. We did not check for this condition in Fig. 4.24, but we clearly can do so, and in fact we must do so if we wish to associate semantic rules with reductions. Thus we have an opportunity to detect errors in Fig. 4.24, modified at steps (10-12) to determine what production is the handle in a reduction.

## Handling Errors During Reductions

We may divide the error detection and recovery routine into several pieces. One piece handles errors of type (2). For example, this routine might pop symbols off the stack just as in steps (10-12) of Fig. 4.24. However, as there is no production to reduce by, no semantic actions are taken; a diagnostic message is printed instead. To determine what the diagnostic should say, the routine handling case (2) must decide what production the right side being popped "looks like." For example, suppose $abc$ is popped, and there is no production right side consisting of $a$, $b$ and $c$ together with zero or more nonterminals. Then we might consider if deletion of one of $a$, $b$, and $c$ yields a legal right side (nonterminals omitted). For example, if there were a right side $aEcE$, we might issue the diagnostic

illegal $b$ on line (line containing $b$)

We might also consider changing or inserting a terminal. Thus if $abEdc$ were a right side, we might issue a diagnostic

missing $d$ on line (line containing $c$)

We may also find that there is a right side with the proper sequence of terminals, but the wrong pattern of nonterminals. For example, if $abc$ is popped off the stack with no intervening or surrounding nonterminals, and $abc$ is not a right side but $aEbc$ is, we might issue a diagnostic

missing $E$ on line (line containing $b$)

Here $E$ stands for an appropriate syntactic category represented by nontermi-nal $E$. For example, if $a$, $b$, or $c$ is an operator, we might say "expression;" if $a$ is a keyword like if, we might say "conditional."

In general, the difficulty of determining appropriate diagnostics when no legal right side is found depends upon whether there are a finite or infinite number of possible strings that could be popped in lines (10-12) of Fig. 4.24. Any such string $b_1 b_2 \cdots b_k$ must have $\doteq$ relations holding between adjacent symbols, so $b_1 \doteq b_2 \doteq \cdots \doteq b_k$. If an operator precedence table tells us that there are only a finite number of sequences of terminals related by $\doteq$, then we can handle these strings on a case-by-case basis. For each such string $x$ we can determine in advance a minimum-distance legal right side $y$ and issue a diagnostic implying that $x$ was found when $y$ was intended.

It is easy to determine all strings that could be popped from the stack in steps (10-12) of Fig. 4.24. These are evident in the directed graph whose nodes represent the terminals, with an edge from $a$ to $b$ if and only if $a \doteq b$. Then the possible strings are the labels of the nodes along paths in this graph. Paths consisting of a single node are possible. However, in order for a path $b_1 b_2 \cdots b_k$ to be "poppable" on some input, there must be a symbol $a$ (pos-sibly $) such that $a \lessdot b_1$. Call such a $b_1$ *initial*. Also, there must be a sym-bol $c$ (possibly $) such that $b_k \gtrdot c$. Call $b_k$ *final*. Only then could a reduction be called for and $b_1 b_2 \cdots b_k$ be the sequence of symbols popped. If the graph has a path from an initial to a final node containing a cycle, then there are an infinity of strings that might be popped; otherwise, there are only a fin-ite number.



Fig. 4.27. Graph for precedence matrix of Fig. 4.25.

**Example 4.31.** Let us reconsider grammar (4.17):

$$E \rightarrow E+E \mid E-E \mid E*E \mid E/E \mid E \uparrow E \mid (E) \mid -E \mid \text{id}$$

The precedence matrix for this grammar was shown in Fig. 4.25, and its graph is given in Fig. 4.27. There is only one edge, because the only pair related by $\doteq$ is the left and right parenthesis. All but the right parenthesis are initial, and all but the left parenthesis are final. Thus the only paths from an initial to a final node are the paths $+$, $-$, $*$, $/$, id, and $\uparrow$ of length one, and the path from ( to ) of length two. There are but a finite number, and each corresponds to the terminals of some production's right side in the grammar. Thus the error checker for reductions need only check that the proper set of

nonterminal markers appears among the terminal strings being reduced. Specifically, the checker does the following:

1.  If $+$, $-$, $*$, $/$, or $\uparrow$ is reduced, it checks that nonterminals appear on both sides. If not, it issues the diagnostic

                          missing operand

2.  If **id** is reduced, it checks that there is no nonterminal to the right or left. If there is, it can warn

                          missing operator

3.  If ( ) is reduced, it checks that there is a nonterminal between the parentheses. If not, it can say

              no expression between parentheses

Also it must check that no nonterminal appears on either side of the parentheses. If one does, it issues the same diagnostic as in (2).         □

If there are an infinity of strings that may be popped, error messages cannot be tabulated on a case-by-case basis. We might use a general routine to determine whether some production right side is close (say distance 1 or 2, where distance is measured in terms of tokens, rather than characters, inserted, deleted, or changed) to the popped string and if so, issue a specific diagnostic on the assumption that that production was intended. If no production is close to the popped string, we can issue a general diagnostic to the effect that "something is wrong in the current line."

### Handling Shift/Reduce Errors

We must now discuss the other way in which the operator-precedence parser detects errors. When consulting the precedence matrix to decide whether to shift or reduce (lines (6) and (9) of Fig. 4.24), we may find that no relation holds between the top stack symbol and the first input symbol. For example, suppose $a$ and $b$ are the two top stack symbols ($b$ is at the top), $c$ and $d$ are the next two input symbols, and there is no precedence relation between $b$ and $c$. To recover, we must modify the stack, input or both. We may change symbols, insert symbols onto the input or stack, or delete symbols from the input or stack. If we insert or change, we must be careful that we do not get into an infinite loop, where, for example, we perpetually insert symbols at the beginning of the input without being able to reduce or to shift any of the inserted symbols.

One approach that will assure us no infinite loops is to guarantee that after recovery the current input symbol can be shifted (if the current input is $, guarantee that no symbol is placed on the input, and the stack is eventually shortened). For example, given $ab$ on the stack and $cd$ on the input, if $a \lessdot c^2$

---

[2] We use $\lessdot$ to mean $<\cdot$ or $\doteq$.

we might pop $b$ from the stack. Another choice is to delete $c$ from the input if $b \leq \cdot d$. A third choice is to find a symbol $e$ such that $b \leq \cdot e \leq \cdot c$ and insert $e$ in front of $c$ on the input. More generally, we might insert a string of symbols such that

$$b \leq \cdot e_1 \leq \cdot e_2 \leq \cdots \leq e_n \leq \cdot c$$

if a single symbol for insertion could not be found. The exact action chosen should reflect the compiler designer's intuition regarding what error is likely in each case.

For each blank entry in the precedence matrix we must specify an error-recovery routine; the same routine could be used in several places. Then when the parser consults the entry for $a$ and $b$ in step (6) of Fig. 4.24, and no precedence relation holds between $a$ and $b$, it finds a pointer to the error-recovery routine for this error.

**Example 4.32.** Consider the precedence matrix of Fig. 4.25 again. In Fig. 4.28, we show the rows and columns of this matrix that have one or more blank entries, and we have filled in these blanks with the names of error handling routines.

|    | id | ( | ) | $ |
|----|----|----|----|----|
| id | e3 | e3 | ·> | ·> |
| (  | <· | <· | ≐ | e4 |
| )  | e3 | e3 | ·> | ·> |
| $  | <· | <· | e2 | e1 |

**Fig. 4.28.** Operator-precedence matrix with error entries.

The substance of these error handling routines is as follows:

e1: /* called when whole expression is missing */
    insert id onto the input
    issue diagnostic: "missing operand"

e2: /* called when expression begins with a right parenthesis */
    delete ) from the input
    issue diagnostic: "unbalanced right parenthesis"

e3: /* called when id or ) is followed by id or ( */
    insert + onto the input
    issue diagnostic: "missing operator"

e4: /* called when expression ends with a left parenthesis */
    pop ( from the stack
    issue diagnostic: "missing right parenthesis"

Let us consider how this error-handling mechanism would treat the

erroneous input **id** + ). The first actions taken by the parser are to shift **id**, reduce it to $E$ (we again use $E$ for anonymous nonterminals on the stack), and then to shift the + . We now have configuration

| STACK | INPUT |
|-------|-------|
| $\$E +$ | )$ |

Since + $\cdot>$ ) a reduction is called for, and the handle is + . The error checker for reductions is required to inspect for $E$'s to left and right. Finding one missing, it issues the diagnostic

<center>missing operand</center>

and does the reduction anyway.

Our configuration is now

| $\$E$ | )$ |
|-------|-----|

There is no precedence relation between $ and ), and the entry in Fig. 4.28 for this pair of symbols is e2. Routine e2 causes diagnostic

<center>unbalanced right parenthesis</center>

to be printed and removes the right parenthesis from the input. We are now left with the final configuration for the parser.

| $\$E$ | $ | □ |
|-------|---|---|

## 4.7 LR PARSERS

This section presents an efficient, bottom-up syntax analysis technique that can be used to parse a large class of context-free grammars. The technique is called $LR(k)$ parsing; the "L" is for left-to-right scanning of the input, the "R" for constructing a rightmost derivation in reverse, and the $k$ for the number of input symbols of lookahead that are used in making parsing decisions. When $(k)$ is omitted, $k$ is assumed to be 1. LR parsing is attractive for a variety of reasons.

- LR parsers can be constructed to recognize virtually all programming-language constructs for which context-free grammars can be written.

- The LR parsing method is the most general nonbacktracking shift-reduce parsing method known, yet it can be implemented as efficiently as other shift-reduce methods.

- The class of grammars that can be parsed using LR methods is a proper superset of the class of grammars that can be parsed with predictive parsers.

- An LR parser can detect a syntactic error as soon as it is possible to do so on a left-to-right scan of the input.

handle - substring that match
the right side of the
CLASSMATE
Date
Page
prod

**Shift Reduce Process:**

| stack | RSF | Action |
|-------|-----|--------|
| $ | 000 111 | shift 01 |
| $ 01 | 0011 | reduce $S \to 01$ |
| $ s | 0011 | |

**24/04/18:** Conflict in shift reduce parsing:

1. Shift Reduce conflict     — example:
2. Reduce-Reduce conflict

| Stack | RSF | Action |
|-------|-----|--------|
| $ | iE+ses$ | |
| $iE+s | es$ | → shift els |
| s | | reduce |
| | | By |

↓

2 productions with same
production on the
right

$S \to id$    { Parser doesn't know either id shou
$P \to id$    be reduced to $S$ or $P$.

---

**Operator precedence Parser:**

Grammer G is operator grammer iff:

i) No E production

ii) No two adjacent non-terminals.

operator
——→ Grammer

↓ input

Operator precedence Pa

↓
parse tree (postfix
expression

Ex: $E \to EAE | id$   } not a OG
    $A \to * | +$

$E \to E+E | E*E | id$   ✓ OG.

---

Steps for operator precedence parsing:

~~Problem~~:

1. Check whether the given grammer is operator grammer or not.
possible try to convert.

2. Generate operator relation table

3. Parse the input string

4. Construct the parse tree.

PROBLEM:

1. Construct the operator precedence parser for the given grammar and parse the given input string.

$$E \rightarrow EAE \mid id$$
$$A \rightarrow + \mid *$$

(i) Converting to operator grammar

$$E \rightarrow E+E \mid E*E \mid id.$$

(ii) Operator Relation Table:

Assumption: identifier - highest precedence          Right associative
     * - left associative   $(\cdot >)$                    $(< \cdot)$
     + - left associative.
     $ - least precedence

|     | id  | +       | *       | $       |
|-----|-----|---------|---------|---------|
| id  | –   | $\cdot >$ | $\cdot >$ | $\cdot >$ |
| +   | $< \cdot$ | $\cdot >$ | $< \cdot$ | $\cdot >$ |
| *   | $< \cdot$ | $\cdot >$ | $\cdot >$ | $\cdot >$ |
| $   | $< \cdot$ | $< \cdot$ | $< \cdot$ | – Accept |

(iv)



(iii) Parse the input string.

| Stack     | input        | Relation    | Action   |
|-----------|--------------|-------------|----------|
| $         | id+id*id$    | $< \cdot$   | push id  |
| $id       | +id*id$      | $\cdot >$   | pop id   |
| $         | +id*id$      | $< \cdot$   | push +   |
| $+        | id*id$       | $< \cdot$   | push id  |
| $+id      | *id$         | $\cdot >$   | pop id   |
| $+        | *id$         | $< \cdot$   | push *   |
| $+*       | id$          | $< \cdot$   | push id  |
| $+*id     | $            | $\cdot >$   | pop id   |
| $+*       | $            | $\cdot >$   | pop *    |
| $+        | $            | $\cdot >$   | pop +    |
| $         | $            | –           | Accept   |

$\$ \rightarrow T$
$T \rightarrow T$
$\downarrow$
$T \rightarrow$
$\downarrow T$

28/04/18

∴  E → E+E | E-E | E*E | E/E | E↑E | (E) | id.

input: id * (id ↑ id) - id / id

i) It is an OG

ii) Generate relation table

    id - highest precedence      + - → left associative
    () - right associative      $ - least precedence
    ↑ - right associative
    */ - left associative

|     | id | + | - | * | / | ↑ | ( | ) | $ |
|-----|----|----|----|----|----|----|----|----|----|
| id  | -  | ⋗ | ⋗ | ⋗ | ⋗ | ⋗ | - | ⋗ | ⋗ |
| +   | ⋖ | ⋗ | ⋗ | ⋖ | ⋖ | ⋖ | ⋖ | ⋗ | ⋗ |
| -   | ⋖ | ⋗ | ⋗ | ⋖ | ⋖ | ⋖ | ⋖ | ⋗ | ⋗ |
| *   | ⋖ | ⋗ | ⋗ | ⋗ | ⋗ | ⋖ | ⋖ | ⋗ | ⋗ |
| /   | ⋖ | ⋗ | ⋗ | ⋗ | ⋗ | ⋖ | ⋖ | ⋗ | ⋗ |
| ↑   | ⋖ | ⋗ | ⋗ | ⋗ | ⋗ | ⋖ | ⋖ | ⋗ | ⋗ |
| (   | ⋖ | ⋖ | ⋖ | ⋖ | ⋖ | ⋖ | ⋖ | ≐ | - |
| )   | -  | ⋗ | ⋗ | ⋗ | ⋗ | ⋗ | - | ⋗ | ⋗ |
| $   | ⋖ | ⋖ | ⋖ | ⋖ | ⋖ | ⋖ | ⋖ | - | Accept |

iii) Parse input:

| Stack | input | Relation | Action |
|-------|-------|----------|--------|
| $ | id*(id↑id)-id/id$ | ⋖ | push id |
| $id | *(id↑id)-id/id$ | ⋗ | pop id |
| $ | *(id↑id)-id/id$ | ⋖ | push * |
| $* | (id↑id)-id/id$ | ⋖ | push ( |
| $*( | id↑id)-id/id$ | ⋖ | push id |
| $*(id | ↑id)-id/id$ | ⋗ | pop id |
| $*( | ↑id)-id/id$ | ⋖ | push ↑ |
| $*(↑ | id)-id/id$ | ⋖ | push id |
| $*(↑id | )-id/id$ | ⋗ | pop id |
| $*(↑ | )-id/id$ | ⋗ | pop ↑ |
| $*( | )-id/id$ | ≐ | push ) |
| $*() | -id/id$ | ⋗ | pop ),( |
| $* | -id/id$ | ⋗ | pop * |

| | | | |
|---|---|---|---|
| $ - | id/id $ | < | push id |
| $-id | /id $ | > | pop id |
| $ - | /id $ | < | push / |
| $-/ | id $ | < | push id |
| $-/ id | $ | > | pop id |
| $-/ | $ | > | pop / |
| $- | $ | > | pop - |
| $ | $ | - | Accept |

**Algorithm: Operator precedence parsing algorithm:**

Input: An input string w and a table of precedence relations.

Output: If w is well formed, a skeletal parse tree, with a placeholder non-terminal E labelling all interior nodes otherwise, an error indication.

Method: Initially the stack contains $ and the input buffer the input buffer the string w$. To parse, we execute the program.

1. Set input to point to the first symbol of w$.

2. repeat forever:

3. if $ is on top of stack and ip points to $ then

4. return

5. Let a be the topmost terminal symbol on the stack and let b be the symbol pointed to by ip.

6. If a < b or a ≐ b then begin

7.     push b onto the stack

8. advance ip to next input symbol
   end;

9. else if a -> b then

10. repeat

11. pop the stack over any of ≐

12. until the top stack terminal is related by < to the terminal most recently popped

13. else error()

3. $S \rightarrow (L) \mid a$

$L \rightarrow L, S \mid S$

i/p : $(a, (a,a))$

a - highest     ( - left

$ least      , - left

i) It is an OG

iiy Relation table:

Note: Do remember to put the relation operator based on associativity i.e which one is evaluated first. After parsing we get the parse tree

|     | id  | ,   | (   | )   | $   |
|-----|-----|-----|-----|-----|-----|
| id  | —   | ·>  | ·>  | ·>  | ·>  |
| ,   | <·  | —   | <·  | ·>  | ·>  |
| (   | <·  | <·  | <·  | ≐   | —   |
| )   | —   | ·>  | —   | ·>  | ·>  |
| $   | <·  | <·  | <·  | —   | Accept |

iii) Input is $(a, (a,a))$

| Stack | Input | Relation | Action. |
|-------|-------|----------|---------|
| $ | $(a,(a,a))$ | $<$ | Push ( |
| $( | $a,(a,a))$ | $<$ | Push a |
| $(a | $,(a,a))$ | $>$ | Pop a |
| $( | $,(a,a))$ | $<·$ | push , |
| $(, | $(a,a))$ | $<·$ | push ( |
| $(,( | $a,a))$ | $<·$ | push a |
| $(,(a | $,a))$ | $>$ | pop a |
| $(,( | $,a))$ | $<·$ | pop, push , |
| $(,(, | $a))$ | $<·$ | push, push a |
| $(,(,a | $))$ | $>$ | pop, pop |
| $(,(, | $))$ | $>$ | pop, pop |
| $(,( | $))$ | $≐$ | push, push ) |
| $(,() | $)$ | $>$ | pop, pop |
| $(, | $)$ | $>$ | pop , |
| $( | $)$ | $≐$ | psh ) |
| $() | $ | $>$ | pop () |
| $ | $ | Accept | |

iv)



Drawbacks of operator relation table:

→ It is very difficult to handle tokens like '—' which has two precedence functions based on whether it is unary operator or binary operator.

→ Only small class of grammars can be parsed

→ If we ever have 4 operators, then the no. of entries in the table are $4 \times 4 = 16$ entries i.e in general, if the number of operators are $n$, we need $O(n^2)$ entries. To overcome this we go for operator precedence functions.

Operator precedence functions:

- The parsers doesnot store relation table instead they make use of precedence functions which map the terminal symbols to integers.

- It uses two functions i.e $f_a$ and $g_b$ for the symbols $a$ and $b$.
   (edge)

(i) if $a > b$, then there is an arrow from function $f_a$ to function $g_b$

(ii) if $a < b$, then there is an edge from $g_b$ to $f_a$

(iii) If $a = b$ then $f_a = g_b$ are in the same group. Note that even if they are not related by $=$ directly we group them together for example if $a = b$ and $c = b$ then $f_a$ and $f_c$ are in the same group, since they are both in the same group as $g_b$

(iv) If the graph constructed has no cycle, then the precedence functions exists.

(v) Find the longest path in the function starting from terminal to $ i.e $f(a)$ to $ and $g(a)$ to $. Using these

Example: $E \to E+E \mid E*E \mid id$

|    | id  | +   | *   | $   |
|----|-----|-----|-----|-----|
| id | -   | ·>  | ·>  | ·>  |
| +  | <·  | ·>  | <·  | ·>  |
| *  | <·  | ·>  | ·>  | ·>  |
| $  | <·  | <·  | <·  | -   |



The longest path is

$fid \to g* \to f+ \to q+· \to f\$$

$gid \to f* \to g* \to f+ \to q+ \to$

|   | id | + | * | $ |
|---|----|----|----|----|
| f | 4  | 2  | 4  | 0 |
| g | 5  | 1  | 3  | 0 |

Advantages: lesser entries

disadvantages: For blank entries of relation tables we g
non-blank entries in junction table i.e we cant mak
out the errors during passing

H·W Construct an operator precedence parser for the given gram
and parse an input string

✓ (i) $E \to E+E \mid E*E \mid (E) \mid id$    ip: $(id + id * id)$: it is an OG

(ii) $E \to E+T \mid T$

$T \to T*v \mid v$    input: $a + b * c * d$.

$v \to a \mid b \mid c \mid d$

(ii) i) It is an OG

ii) Relation table

|    | id | +  | *  | $  |
|----|----|----|----|----|
| id | -  | ·> | ·> | ·> |
| +  | <· | ·> | <· | ·> |

input     a+b*c*d$

| Stack | Input | Relation | Action |
|---|---|---|---|
| $ | a+b*c*d$ | <· | push a |
| $a | +b*c*d$ | ·> | pop a |
| $ | +b*c*d$ | <· | push + |
| $+ | b*c*d$ | <· | push b |
| $+b | *c*d$ | ·> | pop b |
| $+ | *c*d$ | <· | push * |
| $+* | c*d$ | <· | push c |
| $+*c | *d$ | ·> | pop c |
| $+* | *d$ | ·> | pop * |
| $+ | *d$ | <· | push * |
| $+* | d$ | <· | push d |
| $+*d | $ | ·> | pop d |
| $+* | $ | ·> | pop * |
| $+ | $ | ·> | pop + |
| $ | $ | Accept | |

(i)

| Relation table | id | + | * | ( | ) | $ |
|---|---|---|---|---|---|---|
| id | - | ·> | ·> | - | ·> | ·> |
| + | <· | ·> | <· | <· | ·> | ·> |
| * | <· | ·> | ·> | <· | ·> | ·> |
| ( | <· | <· | <· | <· | ≐ | - |
| ) | - | ·> | ·> | - | ·> | ·> |
| $ | <· | <· | <· | <· | - | Accept |

| Stack | input | Relation | Action |
|---|---|---|---|
| $ | (id+id*id)$ | <· | push ( |
| $( | id+id*id)$ | <· | push id |
| $(id | +id*id)$ | ·> | pop id |
| $( | +id*id)$ | <· | push + |
| $(+ | id*id)$ | <· | push id |
| $(+id | *id)$ | ·> | pop id |
| $(+ | *id)$ | <· | push * |
| $(+* | id)$ | <· | push id |
| $(+*id | )$ | ·> | pop id |
| $(+* | )$ | ·> | pop * |
| $(+ | )$ | ·> | pop + |
| $( | )$ | ≐ | push ) |
| $() | $ | ·> | pop |
| $ | $ | Accept | |

# UNIT-5

# SYNTAX - DIRECTED TRANSLATION

## CONTENTS

→ Syntax directed definations
→ Evaluat^n orders for SDD's
→ Applicat^n of SDT
→ SDT schemes.

## Syntax Directed Definations:

A syntax directed defination in a context free grammers with attributes & rules. Attributes are associated with grammer symbols & rules with productions. If 'x' is a symbol, 'a' is one of attributes then we write X.a to denote value of a at a particular parse tree Node labelled X.

* Attributes may be of many kinds: numbers, types, table references, strings, etc...

* 2 types of attributes: ① Synthesized attr
                         ②> Inherited attr

① Synthesized attr: A synth attr for a nonterminal A at a parse tree node N is defined by a semantic rules associated with the production at N. Note that the production must have A as its head. A synthesized attr at node N is defined only in terms of attribute values at the children of N & at N itself.

② Inherited attr   A inherited attr for a nonterminal B at a parse tree node N is defined by a semantic rule associated with the product^n at the parent of N. Note that the product^n must have B as a symbol in its body. An inherited attr at node N is defined only

in terms of attr values at N's parent, N itself & N's sibling

① Synthesised ⟨N → Nodes undr considration⟩ ⟨C → child⟩



Case(i) Single child    Case(ii) Rightmost child    Case(iii) Nochild
⇒ itself.

② Inherited ⟨P → parent⟩ ⟨N → node undr consideratn⟩ ⟨S → Sibling⟩ ⟨O → operator.⟩



Case(i) Sibling    Case (2): inherited from both parent & sibling.

③ Terminals can have Synthesized attributes but not inherited attributes.

★ Attr of terminals have lexical value that are supplied by lexical analysn.

★ Types of SDD: ① S Attributed SDD
                ② L Attributed SDD

① S-Attributed SDD

★ A SDD that involves only synthesised attr then it is called S'-attributed SDD

★ In s-attributed SDD, each rule computes

an attribute for the terminal at the head of a production from attribute taken from body of production.

* S-attributed SDD can be implemented naturally in conjunction with an LR parser / bottom up parser.

* **Annotated parse tree**

A parse tree showing the value(s) of attribute(s) is called an annotated parse tree.

* It is used in bottom-up parser
* Order of evaluation is postorder traversal.

* Example of S-attributed SDD.

| Production | Semantic rules |
|---|---|
| 1) $L \to E n$ | $L.val = E.val$ |
| 2) $E \to E1 + T$ | $E.val = E1.val + T.val$ |
| 3) $E \to T$ | $E.val = T.val$ |
| 4) $T \to T_1 * F$ | $T.val = T_1.val * F.val$ |
| 5) $T \to F$ | $T.val = F.val$ |
| 6) $F \to (E)$ | $F.val = E.val$ |
| 7) $F \to digit$ | $F.val = digit.lexval$ |

② **L-Attributed SDD**

* Example of mixed attributes / L-attributed SDD

| Production | Semantic rules |
|---|---|
| 1) $T \to F T'$ | $T'.inh = F.val$ |
| | $T.val = T'.syn$ |
| 2) $T' \to * F T'_1$ | $T'_1.inh = T'.inh * F.val$ |
| | $T.syn = T'_1.syn$ |
| 3) $T' \to \epsilon$ | $T'.syn = T'.inh$ |
| 4) $F \to digit$ | $F.val = digit.lexval$ |

* A SDD which has both synthesised & inheritid attributes
is called as L-attributed SDD.
* It is used in top down parsing.
* Order of evaluation is topological sorting.

## Evaluating Orders For SDD's

* A dependency graph is used to determine the order
of computation of attributes.

* While an annotated parse trees shows the values
of attributes, a dependency graph helps us to
determine how those values can be computed.

## DEPENDENCY GRAPHS

A dependency graph predicts the flow of
information among the attribute instances in a particular
parse tree. An edge from one attribute instance to
on other means that the value of first is needed to
compute the second.

> For each parse tree node, say node X, the
dependency graph has a node for each attribute
associated with X.

2) If a semantic rule (defines) associated with a
product^n 'p' defines the value of synthized attribute
A.b in terms of value of X.c then the
dependency graph has an edge from X.c to A.b

3) If a semantic rule associated with a product^n
'p' defines the value of inherited attribute
B.c in terms of value of X.a then the

dependency graph has edge from X.C to B.C

Eg1: **production**              **Semantic Rule**

$$E \rightarrow E_1 + T \qquad E.val = E_1.val + T.val.$$

fig: E.val is synthesized from $E_1$.val of T.val



Eg2: **Production**                    **Semantic Rule.**

$$T \rightarrow F T'$$          $T'.inh = F.val$

$$T' \rightarrow * F T'$$        $T.val = T'.syn$

$$T' \rightarrow \epsilon$$        $P'.inh = P'.inh + F.val$

$$F \rightarrow digit$$          $T'.syn = T'.syn$

                             $T'.syn = T'.inh$

                             F.val = digit.lexval

fig: Dependency graph for above production.



## Ordering the evaluation of attributes

* If the dependency graph has an edge from node M to N, then the attribute corresponding to M must be evaluated before attribute of N.

* Thus the only allowable orders of evaluation are those sequences of nodes $N_1, N_2, \ldots N_k$ $\exists$ if there is an edge of the dependency graph from $N_i$ to $N_j$ then $i < j$.

* Such an ordering is called a topological sorting of a graph

* If there is any cycle then no topological sorts, i.e., evaluation of SDD not possible.

* Eg: For dependency graph for (Eq 2) in previous page topological sorting:    1, 2, 3, 4, 5, 6, 7, 8, 9

(or)

1, 3, 5, 2, 4, 6, 7, 8, 9

## S-Attributed Definations
<u>_____</u>

→ An SDD is S-attributed if every attribute is synthesized.

→ When SDD is S-attributed, it attributes evaluated in any bottom-up order of nodes of parse tree.

→ we can have post-order traversal of parse tree to evaluate attributes in S-attributed definat?.

```
Postorder (N) {
        for(each child C of N, from the left)
            postorder (c);
        evaluate the attributes associated with
        Node N;
    }
```

→ S-Attributed definations can be implemented during bottom up parsing without the need to explicitly create parse trees.

# - Attributed Definations

* A SDD is L-attributed if the edges in dependency graph goes from Left to Right but not from Right to left.

* More precisely, each attribute must be either
  → Synthesized.

  → Inherited, but if there is a production $A \to X_1 X_2 \dots X_n$, if there is an inherited attribute $X_i \cdot a$ Computed by a rule associated with this product then the rule many only use:

  (a) Inherited attributes associated with the head A.

  (b) Either inherited or synthesized attr associated with the occurance of symbols $X_1, X_2 \dots, X_{i-1}$ located to the left of $X_i$.

  (c) Inherited/Synthesized attr associated with the, occurance of $X_i$ itself but only in such a way that there is no cycle in the graph.

# PROBLEMS

(5.1) Write a SDD for simple disk calculator

SDL: SDD defination    $i/p = 3 * 5 + 4n$

| PRODUCTION | SEMANTIC RULES |
|---|---|
| 1) $L \to En$ | $L.val = E.val$ |
| 2) $E \to E_1 + T$ | $E.val = E_1.val + T.val$ |
| 3) $E \to T$ | $E.val = T.val$ |
| 4) $T \to T_1 * F$ | $T.val = T_1.val * F.val$ |
| 5) $T \to F$ | $T.val = F.val$ |
| 6) $F \to (E)$ | $F.val = E.val$ |
| 7) $F \to digit$ | $F.val = digit.lexval$ |

L.val = 19

E.val = 19        n

E.val = 15   +        T.val = 4

T.val = 15              F.val = 4

T.val = 3  *            digit.lexval = 4

F.val = 3     F.val = 5

digit.lexval = 3   digit.lexval = 5

Step3: Dependency graph.



L.val ⑫

E.val ⑪        n

E.val ⑦   +        T.val ⑩

T.val ⑥              F.val ⑨

T.val ③  *      F.val ⑤    digit lexval ⑧

F.val ②        digit lexval ④

digit lexval ①

Step4: Topological Ordering : ① ② ③ ④ ⑤ ⑥ ⑦ ⑧ ⑨ ⑩ ⑪ ⑫

Write the SDD & construct annotated parse tree, dependency graph.

a) $(3 + 4) * (5 + 6)n$

b) $1 * 2 * 3 * (4 + 5)n$

c) $(9 + 8 * (7 + 6) + 5) * 4n$

d) $(3 + 4) * (5 + 6)n$

Step1 : SDD Defination

| PRODUCTION | SEMANTIC RULES |
|---|---|
| $L \to En$ | $L.val = E.val$ |
| $E \to E + T$ | $E.val = E.val + T.val$ |
| $E \to T$ | $E.val = T.val$ |
| $T \to T * F$ | $T.val = T.val * F.val$ |
| $T \to F$ | $T.val = F.val$ |
| $F \to (E)$ | $F.val = E.val$ |
| $F \to digit$ | $F.val = digit.lexval$ |

Step2 : Annotated Parse Tree



$L.val = 77$

$E.val = 77 \quad n$

$T.val = 77$

$T.val = 7 \quad *$

$F.val = 11$

$F.val = 7$

$( \quad E.val = 11 )$

$( \quad E.val = 7 )$

$E.val + \quad T.val = 5 \quad T.val = 6$

$E.val = 3 \quad + \quad T.val = 4$

$T.val = 5 \quad F.val = 6$

$T.val = 3$

$F.val = 4$

$F.val = 5 \quad digit.lexval = 6$

$F.val = 3$

$digit.lexval = 3 \quad digit.lexval = 4 \quad digit.lexval = 5$

Step3: Dependency Graph

b) $1 * 2 * 3 * (4 + 5) n$

Step1

| production | Semantic rules |
|---|---|
| $L \to E n$ | $L.val = E.val$ |
| $E \to E + T$ | $E.val = E.val + T.val$ |
| $E \to T$ | $E.val = T.val$ |
| $T \to T * F$ | $T.val = T.val * F.val$ |
| $T \to F$ | $T.val = F.val$ |
| $F \to (E)$ | $F.val = F.val$ |
| $F \to digit$ | $F.val = digit.lexval$ |

R: Anotated parse tree

$d.val = 54$

$E.val = 54$     $n.$

$T.val = 54$

$T.val = 6$  *     $F.val = 9$

$T.val = 2$ *   $F.val = 3$   ( digit.le   $E = 9$ )

$T.val = 1$ *   $F.val = 2$   digit. lexval = 3   $E.val + = 4$   $T.val = 5$

$F.val = 1$         digit. lexval = 2       $T.val = 4$       $F.val = 5$

digit. lexval = 1                           $F.val = 4$       digit. lexval = 5

                                            digit. lexval = 4

## Stp3 : Dependency graph



Stp4 : Topological sortig : 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21

c) $(9 * 8 * (7+6) + 5) * 4n$
         $+$

Step1: SDD definition

| Production | Semantic rules |
|---|---|
| $\alpha \rightarrow En$ | $\alpha.val = E.val$ |
| $E \rightarrow E + T$ | $E.val = E.val + T.val$ |
| $E \rightarrow T$ | $E.val = T.val$ |
| $T \rightarrow T * F$ | $T.val = T.val * F.val$ |
| $T \rightarrow F$ | $T.val = F.val$ |
| $F \rightarrow (E)$ | $F.val = E.val$ |
| $F \rightarrow digit$ | $F.val = digit.lexval$ |

Step2: Annotated parse tree

$\alpha.val = 472$

$E.val = 472$      $n.$

$T = 472$

$T.val = \frac{3*4}{118}$      $F.val = 4.$

$F.val = 118$      $digit.lexval = 4$

$($  $E.val\ 118$  $)$

$E.val\ 113 +$      $T.val = 5$

$E.val +$  $T.val = 104$  $F.val = 5$
$= 9$

$T.val = 9$  $T.val = 8 *$  $F.val = 13$      $digit.lexval = 5$

$F.val = 9$  $F.val = 8$  $($  $E.val) = 13$

$digit.lexval = 9$  $digit.lexval = 8$  $E.val + $  $T.val$

$T.val$  $F.val$

$F.val$  $digit.lexval = 6$

$digit.lexval = 7$

# Step3: Dependency graph



## Topological sorting:
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18
19 20 21 22 23 24 25 26 27 28 29

Write L-attributed SDD (01) write a SDD for top down parser & construct annotated parse tree, dependency graph for given i/p.

① 3 * 5

Stp1 : SDD Dyination

| Production | Semantic rules. |
|---|---|
| $E \rightarrow TE'$ | $E'. inh = T. val$ <br> $E. val = E'. syn$ |

| | |
|---|---|
| $E' \rightarrow +TE'_1$ | $E'_1.inh = E'.inh + P.inh$ |
| | $E'.syn = E'_1 syn$ |
| $E' \rightarrow \epsilon$ | $E'.syn = E'.inh$ |
| $T \rightarrow FT'$ | $T'.inh = F.val$ |
| | $T.syn = T'.syn$ |
| $T' \rightarrow * FT'_1$ | $T'_1.inh = T'.inh * F.val$ |
| | $T'.syn = T'_1.syn$ |
| $T' \rightarrow \epsilon$ | $T'.syn = T'.inh$ |
| $F \rightarrow (E)$ | $F.val = E.val$ |
| $F \rightarrow digit$ | $F.val = digit.lexval$ |

## NOTE

① The original grammar symbols [i.e., E, T, F] will only have synthesized attributes [i.e., val]

② The augmented grammar symbols [ ' ] will have both inherited & synthesised attributes i.e., inh & syn

Step 2: Annotated parse tree

$E.val = 15$

$T.val = 15$   $E'.inh = 15$
              $.syn = 15$

$F.val = 3$   $T'.inh = 3$   $\epsilon$
              $.syn = 15$

digit lexval   *   $F.val = 15$   $T'.inh = 15$
= 3                              $.syn = 15$

              digit.lexval       $\epsilon$
              = 5

## Q.3 : Dependency graph



Topological sorting : 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18
19 20 21 22 23 24 25 26 27 28 29

Write L-attributed SDD (02) Write a SDD for top down parser & construct annotated parse tree, dependency graph for given i/p.

① 3 * 5

Step 1 : SDD Definition

| Productn | Semantic rules. |
|---|---|
| E → TE' | E'.inh = T. Val<br>E. syn = E'. syn |

| | |
|---|---|
| $E' \rightarrow +TE'_1$ | $E'_1.inh = E'.inh + T.inh$ <br> $E'.syn = E'_1.syn$ |
| $E' \rightarrow \epsilon$ | $E'.syn = E'.inh$ |
| $T \rightarrow FT'$ | $T'.inh = F.val$ <br> $T.syn = T'.syn$ |
| $T' \rightarrow * FT'_1$ | $T'_1.inh = T'.inh * F.val$ <br> $T'.syn = T'_1.syn$ |
| $T' \rightarrow \epsilon$ | $T'.syn = T'.inh$ |
| $F \rightarrow (E)$ | $F.val = E.val$ |
| $F \rightarrow digit$ | $F.val = digit.lexval$ |

## NOTE

① The original grammar symbols [i.e., E, T, F] will only have synthesized attributes $\boxed{i.e., val}$

② The augmented grammar symbols $\boxed{[\,']\,}$ will have both inherited & synthesized attributes i.e.,

$$inh \And syn$$

Step2: Annotated parse tree

Step3: Dependency graph.

E.val (12)

(11) inh E'. syn (10)

T.val (9)

F.val (2) (3) inh T' syn (8)

(1) E

digit lexval *

F.val (5)(6) inh T' syn (7)

(4)

digit lexval E

Step4: Topological Order → 1 2 3 4 5 6 7 8 9 10 11 12.

② 3 + 5

Step1: SDD definition

// Same as previous problem.

Step2: Annotated parse tree

E.val = 8

T.val = 3     E'. inh = 3
              . syn = 8

F.val = 3  T'. inh = 3
           . syn = 3   +   T.val = 5   E'. inh = 5+3 = 8
                                       . syn = 8

digit. lexval = 3     E

                                F.val    T'. inh = 5
                                = 5      . syn = 5

                                         E

                              digit. lexval = 5

Step3: Dependency graph

E.val (15)

(6) inh E'. syn (14)

T.val (5)

(13)

F.val (2) (3) inh T' syn (4)    +    T.val (11)  inh E' syn (12)

(1)                                              E

digit lexval     E

                              F.val (8) inh T' syn (10)

                              (7)        (9)

                              digit lexval    E

Sol<sup>n</sup>: Topological Ordering = 1 2 3 4 .. .. .. .. 14 15

③ 3 * 5 + 4

Stp1: SDD dyfination

// Same as problem II → ①

Stp2: Annotated parse tree



E.val=19

T.val=15

E'. inh =15
E'. syn=19

F.val=3

T'. inh=3
T'. syn=15

+

T.val=4 E'. inh =15+4=19
E'. syn=19

digit.lexval
=3

*

F.val
=5

T'. inh=3*5=15
T'. syn=15

F.val
=4

T'. inh=4
T'. syn=4

digit.lexval
=5

ε

digit.lexval = 4

ε

Stp3: Dependency graph



E.val ⑲

T.val⑨ ⑩ inh E' inh ⑱

F.val ② ③ inh T' syn ⑧ + T.val⑮ ⑯ inh E' syn ⑰

① * digit.lexval ⑤ F.val inh T' syn ⑬ inh T' syn ⑭ ε

④ digit.lexval ε ⑫ F.val ⑪ digit.lexval ε

Topological Order: 1 R 3 4 ... 19

$3 * 5 + 4.$  $(3+4)*(5+6)$

Stp1: SDD definatⁿ
            // same as SDD of II → ①

Stp2: Annotated parse tree



E.val = 77
T.val = 77   E' inh = 77
                syn = 77
                 E
F.val = 7   7 = inh. T'
            77 = syn.
(   E.val 7   )        F.val = 11   T' .inh = 7 * 11 = 77
                                         .syn = 77
                                          E
T.val = 3   E' .inh = 3        (   E.val = 11   )
            .syn = 7
F.val   T' .inh = +     T.val = 7   E' .inh = 3 + 7   T.val = 5   E' .inh = 5   = 5+6 = 11
= 3     .syn = 3                    .syn = 7                      .syn = 11
|        E              F.val T' .inh = +                  F.val T' .inh = +     T.val = 6   E' .inh
digit.lexval            = 7   .syn = 7        digit.lexval   = 5  .syn = 5                    .syn = 11
= 3                      |                    = 5                  |                           E
                      digit.lexval                            E.               F.val T' .inh = 6
                         = 7                 digit.lexval      digit.lexval    = 6   .syn = 6
                                                = 5                                   E
                                                                                 digit.lexval
                                                                                 = 6

Stp3: Dependency graph



Eval  ㊳
Tval ㊵        ㊱ inh E' syn  ㊲
                      E
Fval ⑯   ⑰ inh T' syn  ㉞
⑮                                        ㉜  inh T' syn ㉓
Eval                         *         Fval ㉛
⑤                                      Eval ㉚
Tval   ⑥ inh E' syn  ⑭                 Tval ㉔ ㉓ inh E' syn  ⑲
Fval ② ③ inh T syn +  ⑪ ⑫ E' syn ⑬    Tval inh E' syn ⑳
①                         E             ⑩                      ㉕
digit.lexval   E    Fval ⑧ ⑨ inh T' syn    ⑲ Fval inh T syn + ㉑  Tval inh E syn ㉒
                        ⑦                   ⑱                        Fval
                     digit.lexval           digit.lexval              ㉔
                                                                    digit.lexval

⑤ $1 * 2 * 3 * (4 + 5)$

Step1 : SDD definaton

// Same as II → ①

Step2 : Annotated parse True

E.val = 54

T.val = 54          E'.inh = 54
                    .syn = 54
                    ϵ

F.val = 1       T'.inh = 1
                .syn = 54

digit.lexval      *     F.val = 2   T'.inh = 1*2 = 2
= 1                                 .syn = 54

                digit.lexval *   F.val = 3   T'.inh = 2*3 = 6.
                = 2                           .syn = 54

                        digit.lexval  *  F.val = 9   T'.inh = 6*9 = 54
                        = 3                           .syn = 54
                                                      ϵ

                                          (    E.val = 9

                                        T.val = 4   E'.inh = 4
                                                    .syn = 9

                                  F.val = 4  T'.inh = 4   T.val E'.inh
                                             .syn = 4   = 5  .syn
                                  digit.lexval  ϵ
                                  = 4        F.val  T'.inh
                                             = 5   .syn

                                             digit.lexval ϵ
                                             = 5

Step3 : Dependency graph ㉞

E val
T val ㉛              ㉜ inh E' syn ㉝
                     ϵ
F val ②  ③ inh T' syn �30
      ①                ㉙
digit lexval * ⑤ Fval inh ⑥ T' syn
           ④      digit.lexval
              digit.lexval   ⑧ F val ⑨ inh T' syn ㉘
              ⑦ digit lexval  *        ㉕ ㉖ inh T syn ㉗
                                    Fval  ↑
                                    E val ㉔
                                 T val ⑭ ⑮ inh E' syn ㉝
                              ⑪ F val ⑫ inh T syn  +  ⑳ ㉑ inh E' syn ㉒
                                ⑩ digit lexval ϵ    Tval
                                                    ⑲ F val ⑱ inh T syn ϵ
                                                    ⑯ digit lexval ϵ

eg1 : The following dyfination is $L$ attributed. Here the inherited attribute of $T'$ gets its value from its left sibling F. Similarly $T_1'$ gets its value from its parent $T'$ & left sibling F

| Production | Semantic Rules |
|---|---|
| $T \rightarrow FT'$ | $T'.inh = F.val$ |
| $T' \rightarrow *FT_1'$ | $T_1'.inh = T'.inh * F.val$ |

eg2 : The dyfinations below are not $L$-attributed as B.i depends on its right sibling C's attribute.

| Production | Semantic Rules |
|---|---|
| $A \rightarrow BC$ | $A.s = B.b$ |
|  | $B.i = f(C.c, A.s)$ |

## SIDE EFFECTS

Evaluation of semantic rules may generate intermediate codes, Eg : A disk calculator might print a result, a code generator might enter the type of an identifier into a symbol table, may perform type checking & may issue error msgs. These are known as side effects.

## SEMANTIC RULES WITH CONTROLLED SIDE EFFECTS

In practise translation involves side effects. Attribute grammers has no side effects & allow any evaluation order consistent with dependency graph wheras translation schemes impose left to right evaluation & allow scheme actions to contain any program fragment.

# Ways to Control Side Effects

1. permit incidential side effects that do not constrain attribute evaluation.

   In other words, permit side effects when attr evaluat? based on any topological sort of the dependency graph produces a <u>correct</u> translation.

2. Impose constraints on allowable evaluation orders so that the same translation is produced for any allowable order.

Write an SDD for simple type declaration

a) i/p: inta

Step1:

| Production | Semantic rules |
|---|---|
| $D \rightarrow T \alpha$ | $\alpha.inh = T.type$ |
| $T \rightarrow int$ | $T.type = int$ |
| $T \rightarrow float$ | $T.type = float$ |
| $\alpha \rightarrow \alpha_1 , id$ | $\alpha_1.inh = \alpha.inh$  addtype(id.entry, $\alpha.inh$) |
| $\alpha \rightarrow id$ | addtype(id.entry, $\alpha.inh$) |

Step2: Annotated parse tree

```
            D
          /   \
   T type=int   α  inh = int
       |           entry = a
      int           |
                 id. entry = a
```

Step3: Dependency graph

```
          D ⑥
         /    \
   T type②   ③inh α entry ⑤
      |              |
    int①        id.entry ④
```

## explanation:

Non terminal D represents a declaration, which from production 1, consists of a type T followed by a list L of identifiers. T has one attribute: T.type, which is the type in the declaration D. Nonterminal L has one attribute, which call inh to emphasize that it is an inherited attribute. The purpose of L.inh is to pass the declared type down the list of identifiers, so that it can be the appropriate symbol table entries. Production ② & ③ each evaluate the synthesized attribute T.type giving it the appropriate value, integer or float. This type is passed to the attribute L.inh in the rule of production 1. Production 4 passes L.inh down the parser tree i.e., the value of L1.inh is compared at a parse tree node by copying the value of L1.inh from the parent of that node, the parent corresponds to the head of production. Production ④ & ⑤ also have a rule in which a function addtype is called with 2 arguments:

1) id.entry a lexical value that points to a symbol table object.

2) L.inh, the type being assigned to every identifier on the list.

The function addType properly installs the type L.inh as the type of the represented identifier. Note that the side effect, adding the type info to the table, doesnot affect the evaluation order.

**b)** int a,b

**Step1**

| Production | Semantic rules |
|---|---|
| $D \rightarrow T \alpha$ | $\alpha.inh = T.type$ |
| $T \rightarrow int$ | $T.type = int$ |
| $T \rightarrow float$ | $T.type = float$ |
| $\alpha \rightarrow \alpha_1, id$ | $addtype(id.entry, \alpha_1.in)$ $\alpha_1.inh = \alpha.inh$ |
| $\alpha \rightarrow id$ | $addtype(id.entry, \alpha.inh)$ |

**Step2 : Annotated parse tree**



$$T.type = int$$
$$int = a$$
$$\alpha.inh = int \quad \alpha.entry = b$$
$$\alpha.inh = int \quad \alpha.entry = b$$
$$id.entry = b$$
$$id.entry = a$$

**Step3 : Dependency graph**



$T$ type ③, int ①, inh $\alpha$ entry ④, inh $\alpha$ entry ⑤ , id entry ⑦, id entry ⑥, ⑧, ⑨, ② id entry

**Step4 : Topological order** ① ② ③ ④ ⑤ ⑥ ⑦ ⑧ ⑨

(c) | float  a, b, c. |    or  int  a, b, c          ⟨Exercise 5.2.2⟩

Stp1 : SAD defination

| Production | SAD |
|---|---|
| $D \to T\alpha$ | $\alpha.inh = T.type$ |
| $T \to int$ | $T.type = int$ |
| $T \to float$ | $T.type = float$ |
| $\alpha \to \alpha, id$ | $\alpha_1.inh = \alpha.inh$<br>addtype (id.entry, $\alpha.inh$) |
| $\alpha \to id$ | addtype (id.entry, $\alpha.inh$) |

Stp2 : Annotated parse tree



Stp3 : Dependency graph



Stp4 : Topological Ordering : 1  2  3 . . . . . 10

(d) float ω, x, y, z

Step1 SDD definition

// Same as previous problem.

Step2 : Annotated parse tree



D

T.type
|
float

α · inh = float
α · entry = z

α · inh = float
α · entry = y ,

id · entry = z

α · inh = float
α · entry = x ,

id · entry = y

α · inh = float
· entry = ω ,

id · entry = x

id · entry = ω

Step3 : Dependency graph



D

T type ⑤

⑥ inh α entry ⑦

float

⑧

inh α entry

⑩

id entry ⑭

⑨

id entry ⑤

⑩

inh α entry

⑪ ,

id entry ③

⑫

⑬

inh α entry ,

id entry ①

Exercies 5.2.4 This grammar generates binary no with a
decimal point

$$S \rightarrow \alpha . \alpha \mid \alpha$$
$$\alpha \rightarrow \alpha B \mid B$$
$$B \rightarrow 0 \mid 0$$

Assign an $\alpha$ attributed SDD to compute S.val, the decimal no value of i/f string. For eg, translating string 101.101 should be the decimal no 5.625.

Solo

| Production | SAD |
|---|---|
| 1) $S \rightarrow \alpha . \alpha_1$ | 1. $\alpha.inh = 0$ <br> 2. $\alpha_1.inh = -1$ <br> 3. $S.val = \alpha.syn + \alpha_1.syn$ |
| 2) $S \rightarrow \alpha$ | 1. $\alpha.inh = 0$ <br> 2. $S.val = \alpha.syn$ |
| 3) $\alpha \rightarrow \alpha_1 B$ | 1. $\alpha_1.inh = \alpha.inh + 1$ <br> 2. $B.inh = \alpha.inh$ <br> 3. $\alpha.syn = \alpha_1.syn * B.syn$ |
| 4) $\alpha \rightarrow B$ | 1. $\alpha.syn = B.syn * R \wedge \alpha.inh$ |
| 5) $B \rightarrow 0 \mid 1$ | 1. $B.syn = digit.lexval$ |

Exercise 5.2.5 Design an S-attributed SDD for grammar of translation described in 5.2.4

| Production | Semantic Rule. |
|---|---|
| 1) $S \rightarrow \alpha . \alpha_1$ | $S.val = L.lhs + \alpha_1.rhs$ |
| 2) $S \rightarrow \alpha$ | $S.val = \alpha.lhs$ |

| | |
|---|---|
| 3) $A \rightarrow \alpha, B$ | 1) $\alpha . lhs = \alpha_1 . lhs + (2 \times lhs\_exponent * B.val)$ |
| | 2) $\alpha . rhs = \alpha_1 . rhs + (2 ^ \wedge \alpha . rhs\_exponent * B.val)$ |
| | 3) $\alpha . lhs\_exponent = \alpha_1 . lhs\_exponent + 1$ |
| | 4) $\alpha . rhs\_exponent = \alpha_1 . rhs\_exponent + 1$ |
| 4) $\alpha \rightarrow B$ | 1) $\alpha . lhs = 2 ^ \wedge \alpha . lhs exponent * B.val$ |
| | 2) $\alpha . rhs = 2 ^ \wedge \alpha . rhs exponent * B.val$ |
| | 3) $\alpha . lhs\_exponent = 0$ |
| | 4) $\alpha . rhs\_exponent = -1$ |
| 5) $B \rightarrow 0 | 1$ | $B.val = digit . lexval$ |

# Application Of Syntax - Directed Translation

## 1. Construction Of Syntax Tree

SDD's are useful for construction of syntax tree.
A syntax tree is condensed form of parse tree



Parse tree                                    Syntax tree

& Syntax trees are useful for representing programming language constructs like expression & statements.

* They help computer design by decoupling parsing from translation.

* Each node of a syntax tree represent a construct; the children of the node represent the meaningful components of a construct

Eg: A syntax tree node representing an expression $E1 + ER$ has label + & R children representing the sub expression $E1$ & $ER$

* Each node is implemented by object with suitable no of fields; each object will have an Op field that is the label of node with additional fields as follows:

a) If the node is a leaf, an addition field holds the lexical value for the leaf. This is created by function Leaf (op, val).

b) If the node is an interior node, there are as many fields as the node has children in syntax tree. This is created by function Node(op, c1, c2, ... ck)

Example: The S-attributed definition in fig below constructs Syntax trees for a simple expr grammar involving only binary operators + & -. As usual these operators are at the same precedence level & are jointly left associative. All nonterminals have one synthesized attr node, which represents a node of the syntax tree.

$$\boxed{a-h+c}$$

| | PRODUCTION | SEMANTIC RULES |
|---|---|---|
| 1. | $E \rightarrow E_1 + T$ | E.node = new Node ('+', $E_1$.node, T.node) |
| 2. | $E \rightarrow E_1 - T$ | E.node = new Node('-', $E_1$.node, T.node) |
| 3. | $E \rightarrow T$ | E.node = T.node |

| | |
|---|---|
| $T \rightarrow (E)$ | $T.node = E.node$ |
| $T \rightarrow id$ | $T.node = newleaf(id, id.entry)$ |
| $T \rightarrow num$ | $T.node = newleaf(num, num.val)$ |

Stp2 : Parse tree

Syntax tree



Stp3: Syntax tree for a-4+C using above SDD :

Steps in construction of the syntax tree for a-4+c

If the rules are evaluated during a post order traversal of the parse tree, or with reduction during a bottom up parse, then the sequence of steps shown below ends with p5 pointing to the root of the constructed syntax tree.

1) $P_1$ = new Leaf (id, entry_a)

2) $P_2$ = new Leaf (num, 4)

3) $P_3$ = new Node ('-', $P_1$, $P_2$)

4) $P_4$ = new Leaf (id, entry - c)

5) $P_5$ = new Node ('+', $P_3$, $P_4$)

## Constructing Syntax Trees during Top down parsing

With a grammar designed for top-down parsing, the syntax trees are constructed, using the same sequence of steps, even though the structure of the parse trees differs significantly from that of syntax trees. The L-attributed definition below performs the same translation as the S attributed definition shown before.

| | | |
|---|---|---|
| 1) | $E \to TE^1$ | $E.node = E^1.syn$ <br> $E^1.inh = T.node$ |
| 2) | $E^1 \to +TE_1^1$ | $E_1^1.inh = \cancel{E^1.inh + P_1}$ <br> new Node ('+', $E^1.inh$, $T.node$) <br> $E^1.syn = E_1^1.syn$ |
| 3) | $E^1 \to -TE_1^1$ | $E_1^1.inh = $ new Node ('-', $E^1.inh$, $T.node$) <br> $E^1.syn = E_1^1.syn$ |
| 4) | $E^1 \to \epsilon$ | $E^1.syn = E^1.inh$ |
| 5) | $T \to (E)$ | $T.node = E.node$ |

| | | |
|---|---|---|
| 6) | $T \rightarrow id$ | T.node = new Leaf (id, id. entry) |
| 7) | $T \rightarrow num$ | T.node = new Leaf (num, num, val) |

## Dependency Graph for $a-H+C$ with L attributed SDD



## STRUCTURE OF A TYPE

This is an example of how inherited attribute can be used to carry info from one part of the parse tree to another. In C the type $int[2][3]$ can be read as $\boxed{\text{"array of 2 arrays of 3 integer"}}$. The corresponding type expression array $(2, \text{array} (3, \text{integer}))$ is represented by the tree as shown below.



| | production | Semantic rules |
|---|---|---|
| 1) | $T \rightarrow BC$ | $T.t = C.t$<br>$C.b = B.t$ |
| 2) | $B \rightarrow int$ | $B.t = integer$ |
| 3) | $B \rightarrow float$ | $B.t = float$ |
| 4) | $C \rightarrow [num]C_1$ | $C_1.t = array (num.val, C_1)$<br>$C_1.b = C.b$ |
| 5) | $C \rightarrow \epsilon$ | $C.t = C.b$ |

→ The non terminals S & T have a synthesized attribute t representing a type

→ The non terminal C has 2 attributes: an inherited attr (b) & a synthesized attr (t).

→ The inherited attribute, b pass a basic type down the tree

↪ They synthesized attribute, t accumulate the result.

→ An anotated parse tree for e/p : int [2]

(i) int [2]



T.t

B.t =integer

int

C.b = integer
C.t = array (2, integer.

[ num.val
    |
    2

]

C.t = integer.
C.b = integer

ε

Dependency graph



(ii) int [2][3]

Annotated parse tree:



T.t = array (2, array (3, integer))

B.t =integer
|
int

C.b = integer
C.t = array (2, array (3, integer))

[ num.val =2 ]
    |
    2

C.b =integer
C.t =array (3, integer)

[ num.val
    |  =3
    3
]

C.b=integer
C.t
integer
|
ε

<u>Dependency graph</u>



# SYNTAX DIRECTED TRANSLATION SCHEME:

SDT is a complementary notation to SDD.

* All application of SDD can be implemented using SDT.

* SDT is a CFG with program fragments called <u>Seman</u> actions embedded with production bodies.

* Any SDT can be implemented by first building a parse tree & then performing the actions in a left to right, depth first order i.e., during preorder traversal.

* Typically SDT's are implemented during parsing without building parse tree. During parsing, an action in a production body is executed as soon as all the grammar symbols to the left of action have been matched.

* SDT's that can be implemented during parsing can be characterized by introducing distinct marker non terminals in places of each embedded action.

* Each marker M has only one production M→ε.
* If grammar with marker non terminals can be parsed by a given method, then SDT can be implemented

# UNIT-6
# INTERMEDIATE CODE
# GENERATION

Intermediate Code generation:

In the analysis-synthesis model of a compiler, the front end analyzes a source program and creates an intermediate representation, from which the back end generates target code.



Logical structure of a compiler front end

Parsing, static checking and intermediate code generation are done sequentially; sometimes they can be combined and folded into parsing.

Static checking includes type checking, which ensures that operators are applied to compatible operands. It also includes any syntactic checks that remain after parsing.

Ex: It ensures assures that a break-statement in C is enclosed within a while-, for- or switch-statement; an error is reported if such an enclosing statement does not exist.

In the process of translating a program in a given source language into code for a given target machine, a compiler may construct a sequence of intermediate representation as

Source program $\longrightarrow$ High level Intermediate Representation $\rightarrow \ldots \rightarrow$ Low level Intermediate Representation $\longrightarrow$ Target code

High level representations are close to the source language and are well suited to tasks like static type checking.

Ex: Syntax tree

Low level representations are close to the target machine & are suitable for machine dependent tasks like register allocation and instruction selection.

An intermediate representation may either be an actual Language or it may consist of internal data structure that are shared by phases of the compiler.

Variants of Syntax Trees

Nodes in a syntax tree represent constructs in the source program; the children of a node represent the meaningful components of a construct.

A directed acyclic graph (DAG) for an expression identifies the common subexpressions of the expression.

# ① Directed Acyclic Graphs for Expressions.

On DAG leaves represents the atomic operands and interior nodes represents the operators. as in the syntax tree.

A node N in a DAG has more than one parent if N represents a common subexpression; But in the syntax tree, the tree for the common subexpression would be duplicated as many times as the subexpression appears in the original expression.

DAG gives the compiler important clues regarding the generation of efficient code to evaluate the expressions.

Ex: DAG for the expression

$$a + a * (b - c) + (b - c) * d$$



⟶ The leaf for 'a' has 2 parents, because 'a' appears twice in the expression

⟶ The 2 occurance of the common subexpression b-c are represented by one node, the node labeled '−'

## SDD to produce DAG.

| PRODUCTION | SEMANTIC RULES |
|---|---|
| (i) $E \rightarrow E_1 + T$ | $E.node = new\ Node\ ('+', E_1.node, T.node)$ |
| (ii) $E \rightarrow E_1 - T$ | $E.node = new\ Node\ ('-', E_1.node, T.node)$ |
| (iii) $E \rightarrow T$ | $E.node = T.node$ |
| (iv) $T \rightarrow (E)$ | $T.node = E.node$ |
| (v) $T \rightarrow id$ | $T.node = new\ Leaf\ (id, id.entry)$ |
| (vi) $T \rightarrow num$ | $T.node = new\ Leaf\ (num, num.val)$ |

It will construct a DAG, before creating a new node, these functions first check whether an identical node already exists.

If a previously created identical ~~ex~~ node exists, the existing node is returned.

Steps for constructing the DAG for above example.

(i) $P_1 = Leaf\ (id, entry - a)$

(ii) $P_2 = Leaf\ (id, entry - a) = P_1$

(iii) $P_3 = Leaf\ (id, entry - b)$

(iv) $P_4 = Leaf\ (id, entry - c)$

(v) $P_5 = Node\ ('-', P_3, P_4)$

(vi) $P_6 = Node\ ('x', P_1, P_5)$

(vii) $P_7 = Node\ ('+', P_1, P_6)$

(viii) $P_8 = Leaf\ (id, entry - b) = P_3$

(ix) $P_9 = Leaf\ (id, entry - c) = P_4$

(x) $P_{10} = Node\ ('-', P_3, P_4) = P_5$

(xi) $P_{11} = Leaf\ (id, entry - d)$

(xII) $P_{12}$ = Node ('*', $P_5$, $P_{11}$)

(xIII) $P_{13}$ = Node ('+', $P_7$, $P_{12}$)

When the call to Leaf (id, entry- a) is repeated at step 2, the node created by the previous call is returned, so $P_2 = P_1$.

The Value-Number Method for Constructing DAG's

* The nodes of a DAG are stored in an array of records

* Each row of array represents one record & therefore one node.

* In each record, the first field is an operation code, indicating the label of the node. Leaves have one additional field which holds the lexical value and interior nodes have 2 additional fields indicating the left and right children.

Ex:   DAG for $i = i + 10$ allocated in an Array



| | | | |
|---|---|---|---|
| 1 | id | | to entry for i |
| 2 | num | 10 | |
| 3 | + | 1 | 2 |
| 4 | = | 1 | 3 |
| 5 | | | |

DAG                              Array

* In this array, we refer to nodes by giving the integer index of the record for that node within the array.

* This integer historically has been called the "value number" for the node or for the expression represented by the node

* For above example node labeled + has value number 3 & its left & right children have value numbers 1 & 2 respectively

Suppose that nodes are stored in an array & each node is referred to by its value numbers. Let the signature of an Interior node be the triple <op, l, r> where op is the label, l is its left child's value number & r its right child's value number. A unary operator may be assumed to have r = 0.

ALGORITM: To construct the nodes of a DAG using value number method.

INPUT: Label op, node l and node r

OUTPUT: The value number of a node in the array with signature <op, l, r>

MECHOD: Search the array for a node M with label op, left child l and right child r. If there is such a node, return the value number of M. If not, create in the array a new node N with label op, left child l and right child r & return its value number.

Above algorithm yields the desired output, but searching the entire array every time we are asked to locate one node is expensive.

A more efficient approach is to use a hash table, in which the nodes are put into "buckets" each of which typically will have only a few nodes. It supports dictionaries. which is an abstract data type that allows us to insert & delete elements of a set & to determine whether a given element is currently in the set.

To construct a hash table for the nodes of a DAG, we need a hash function $R$ that computes the index of the bucket for a signature $\langle op, l, r \rangle$.

The bucket index $R(op, l, r)$ is computed deterministically from $op, l$ & $r$ so that we may repeat the calculation & always get to the same bucket index for node $\langle op, l, r \rangle$

The buckets can be implemented as linked list as,



Array of buckets q
Readers Indexed
by hash values

List elements
representing nodes

An array indexed by hash value, holds the bucket Readers, each of which points to the first cell of a list. Within the linked list for a bucket, each cell holds the value number of one of the nodes that hash to that bucket. That is, node $\langle op, l, r \rangle$ can be found on the list whose Reader is at index $R(op, l, r)$ of the array.

Thus, given the output input node $op, l$ & $r$ we compute the bucket index $R(op, l, r)$ & search the list of cells in this bucket for the given input node.

For each value number 'v' found in a cell, we must

check whether the signature $\langle op, l, r \rangle$ of the input node matches the node with value number $v$ in the list of the cells. If we find a match, we return $v$. If we find no match, we know no such node can exist in any other bucket, so we create a new cell, add it to the list of cells for bucket-index $h(op, l, r)$ & return the value number in that new cell.

Problems.

Construct the DAG for the expression.

$$((x+y) - ((x+y) * (x-y))) + ((x+y) * (x-y))$$



Construct the DAG & identify the value number for the sub expressions of the following expressions, assuming + associa from the left.

(a) $a + b + (a + b)$

| 1 | ßd | a | |
|---|----|---|---|
| 2 | ßd | 6 | |
| 3 | + | 1 | 2 |
| 4 | + | 3 | 3 |

(b) $a + 6 + a + 6$



| 1 | ßd | a | |
|---|----|---|---|
| 2 | ßd | 6 | |
| 3 | + | 1 | 2 |
| 4 | + | 3 | 1 |
| 5 | + | 4 | 2 |

(c) $a + a + (a + a + a + (a + a + a + a))$



| 1 | ßd | a | |
|---|----|---|---|
| 2 | + | 1 | 1 |
| 3 | + | 2 | 1 |
| 4 | + | 3 | 1 |
| 5 | + | 3 | 4 |
| 6 | + | 2 | 5 |

# Three - Address Code

In 3-address code, there is atmost one operator on the right R side of an instruction, i.e, no built up arithmetic expression are permitted.

Thus a source-language expression like $x + y * z$ might be translated into the sequence of 3 address instructions

$$t_1 = y * z$$

$$t_2 = x + t_1$$

where $t_1$ & $t_2$ are compiler generated names.

3 address code is a linearized representation of a DAG in which explicit names correspond to the interior nodes of the graph.

Ex: Write DAG & its corresponding 3 address code for the expression $a + a * (b-c) + (b-c) * d$.



$$t_1 = b - c$$
$$t_2 = a * t_1$$
$$t_3 = a + t_2$$
$$t_4 = t_1 * d$$
$$t_5 = t_3 + t_4 .$$

DAG                                3 -address code

## Addresses and Instructions

3-address code is built from 2 concepts: address & instructions.

An address can be one of the following.

↳ A name :- For convenience, we allow source program names to appear as addresses in 3-address code. In an implementation, a source name is replaced by a pointer to its symbol table entry, where all information about the name is kept.

↳ A constant : A compiler must deal with many different types of constants and variables.

↳ A compiler generated temporary : It is useful, especially in optimizing compilers, to create a distinct name each time a temporary is needed.

Symbolic labels will be used by instructions that alter the flow of control. A symbolic label represents the index of a 3-address instruction in the sequence of instructions. Actual indexes can be substituted for the labels, either by making a separate pass or by "backpatching"

Here is a list of the common 3-address instruction forms

(i) Assignment instructions of the form x = y op z where op is a binary arithmetic or logical operation & x, y & z are addresses.

(ii) Assignments of the form x = op y, where op is a unary operation

(iii) Copy instructions of the form x = y, where x is assigned the value of y.

(iv) An unconditional jump goto L. The 3-address instruction with label L is the next to be executed.

(v) Conditional jumps of the form if x goto L and if False x goto L. These instructions execute the instruction with label L next if x is true and false respectively.

(vi) Conditional jumps such as if x relop y goto L which apply a relational operator (<, ==, >= etc) to x & y & execute the instruction with label L next if x stands in relation relop to y. If not, the 3 address instruction following if x relop y goto L is executed next, in sequence.

(vii) Procedures calls & returns are implemented using the following instructions: param x for parameters; call p, n & y = call p, n for procedure & function calls respectively & return y, where y representing a returned value, is optional

$$\text{param } x_1$$
$$\text{param } x_2$$
$$\ldots$$
$$\text{param } x_n$$
$$\text{call } p, n$$

The integer n indicating the no. of actual parameters in call p, n is not redundant because calls can be nested.

(viii) Indexed copy instructions of the form x = y[i] and x[i] = y. The instruction x = y[i] sets x to the value in the location i memory units beyond location y. The instruction x[i] = y sets the contents of the location i units beyond x to the value of y.

(Ex) Address & pointer assignments of the form $x = \&y$, $x = *y$
and $*x = y$.

Ex: Consider the statement

$$do \; i = i + 1; \; while \; (a[i] < v);$$

2 possible translations of this statement are

L:  $t_1 = i + 1$
    $i = t_1$
    $t_2 = i * 8$
    $t_3 = a[t_2]$
    if $t_3 < v$ goto L

(a) Symbolic labels

100:  $t_1 = i + 1$
101:  $i = t_1$
102:  $t_2 = i * 8$
103:  $t_3 = a[t_2]$
104:  if $t_3 < v$ goto 100

(b) position number.

The translation in (a) uses a symbolic label L, attached to the first instruction. The translation in (b) shows position number for the instructions, starting arbitrarily at position 100. In both translations, the last instruction is a conditional jump to the first instruction. The multiplication $i * 8$ is appropriate for an array of elements that each take 8 units of space.

## Quadruples

A quadruples has 4 fields, which we call op, arg1, arg2, & result. The op field contains an internal code for the operator.

Ex:  $x = y + z$ is represented by placing $+$ in op, y in arg1, z in arg2 & x in result.

The following are some exceptions to this rule.

(i) Instructions with unary operators like $x = $ minus $y$ or $x = y$ do not use arg2. Note that for a copy statement like $x = y$, op is $=$, while for most other operations, the assignment operator is implied.

(ii) Operators like param use neither arg2 nor result.

(iii) Conditional & unconditional jumps put the target label in result.

Ex: Write quadruples for $a = b * -c + b * -c$.

| | op | arg1 | arg2 | result |
|---|---|---|---|---|
| 0 | minus | c | | $t_1$ |
| 1 | * | b | $t_1$ | $t_2$ |
| 2 | minus | c | | $t_3$ |
| 3 | * | b | $t_3$ | $t_4$ |
| 4 | + | $t_2$ | $t_4$ | $t_5$ |
| 5 | = | $t_5$ | | a |

$t_1 = $ minus $c$

$t_2 = b * t_1$

$t_3 = $ minus $c$

$t_4 = b * t_3$

$t_5 = t_2 + t_4$

$a = t_5$

(a) 3-address code

(b) Quadruples.

The special operator minus is used to distinguish the unary minus operator.

## Triples

A triples has only 3 fields, op, arg1 & arg2. Note that the result field in quadruples is used primarily for temporary names. Using triples, we refer to the result of the operation $x$ op $y$ by its position, rather than by an explicit temporary names.

# Ex: Write triples for $a = b \ast -c + b \ast -c$



(a) Syntax tree

| | op | arg1 | arg2 |
|---|---|---|---|
| 0 | minus | c | |
| 1 | * | b | (0) |
| 2 | minus | c | |
| 3 | * | b | (2) |
| 4 | + | (1) | (3) |
| 5 | = | a | (4) |

(b) Triples.

A ternary operator like $x[i] = y$ requires 2 entries in the triple structure, for example, we can put $x$ & $i$ in one triple & $y$ in the next.

A benifits of quadruples over triples can be seen in an optimizing compiler, where instructions are often moved around. With quadruples, if we move an instruction that computes a temporary $t$, then the instructions that use $t$ require no change. With triples, the result of an operation is referred to by its position, so moving an instruction may require us to change all references to that result.

## Indirect triples.

Indirect triples consist of a listing of pointers to triples, rather than a listing of triples themselves.

With indirect triples, an optimizing compiler can move an

instruction by reordering the instruction list without affecting the triples themselves.

Ex: Write indirect triples for $a = 6 * -C + 6 * -C$:

Instruction

| | |
|----|-----|
| 35 | (0) |
| 36 | (1) |
| 37 | (2) |
| 38 | (3) |
| 39 | (4) |
| 40 | (5) |
| | ... |

| | op | arg1 | arg2 |
|----|-------|------|------|
| 0 | minus | c | |
| 1 | * | 6 | (0) |
| 2 | minus | c | |
| 3 | * | 6 | (2) |
| 4 | + | (1) | (3) |
| 5. | = | a | (4) |
| | | | |

... ...

## Static Single-Assignment Form.

Static single assignment form (SSA) is an intermediate representation that facilitates certain code optimization.

2 distinctive aspects of SSA that distinguish SSA from 3-address code

(i) All assignments in SSA are to variables with distinct names.

Ex:

$$p = a + 6 \qquad\qquad p_1 = a + 6$$

$$q = p - c \qquad\qquad q_1 = p_1 - c$$

$$p = q * d \qquad\qquad p_2 = q_1 * d$$

$$p = e - p \qquad\qquad p_3 = e - p_2$$

$$q = p + q \qquad\qquad q_2 = p_3 + q_1$$

3-address code            Static single assignment form

The same variable may be defined in 2 different control flow paths in a program. For example, the source program

  if (flag) x = -1; else x = 1;

  y = x * a;

If we use different names for x in the true part & false then conflict arises which name should use in $y = x*a$.

(ii) SSA uses a notational convention called $\phi$-function to combine the 2 definitions of $x$

  if (flag) $x_1 = -1$; else $x_2 = 1$;

  $x_3 = \phi(x_1, x_2)$;

Here $\phi(x_1, x_2)$ has the value $x_1$ if the control flow passes through the true part of the conditional & the value $x_2$ if the control flow passes through the false part.
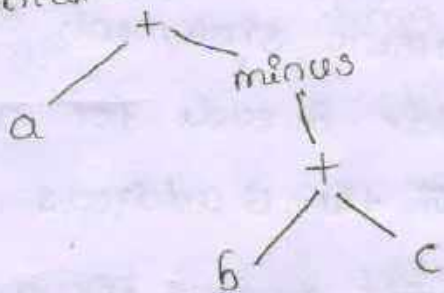
Translate the arithmetic expression $a + -(b+c)$ into

(a) A syntax tree

(b) Quadruples

(c) Triples

(d) Indirect triples.

(a) Syntax tree.



$t_1 = b + c$

$t_2 = minus \; t_1$

$t_3 = a + t_2$

## (b) Quadruples

| | op | arg 1 | arg 2 | result |
|---|---|---|---|---|
| 0 | + | b | c | $t_1$ |
| 1 | minus | $t_1$ | | $t_2$ |
| 2 | = | a | | $t_3$ |

| | op | arg1 | arg2 |
|---|---|---|---|
| 0 | + | b | b c |
| 1 | minus | (0) | |
| 2 | = | a | (1) |

### (d) Indirect triples.

Instructions

| | |
|---|---|
| 35 | (0) |
| 36 | (1) |
| 37 | (2) |

| | op | arg1 | arg2 |
|---|---|---|---|
| 0 | + | b | c |
| 1 | minus | (0) | |
| 2 | = | a | (1) |

## Translation of Expressions.

An expression with more than one operator, like $a + b * c$, well translate into instructions with almost one operator per instruction. An array reference $A[i][j]$ well expand into a sequence of 3-address instructions that calculate an address for the reference.

## Operations within Expressions

The following syntax-directed definition builds up the 3-address code for an assignment statement S using attribute code for S & attributes addr & code for an expression E. Attributes S.code & E.code denotes the 3 address code for S & E respectively. Attribute E.addr denotes the address

# QUESTIONS

1. Define quadruples, triples and static single assignment form.

A quadruples has 4 fields, op, arg1, arg2 & re...

The op field contains an incremental code for the operato...

Ex: the quadruples for $a = (6 * - c) + (6 * - c)$

$t_1 = minus \ c$

$t_2 = 6 * t_1$

$t_3 = minus \ c$

$t_4 = b * t_3$

$t_5 = t_2 + t_4$

$a = t_5$

| | op | arg1 | arg2 | result |
|---|---|---|---|---|
| 0 | minus | c | | $t_1$ |
| 1 | * | 6 | $t_1$ | $t_2$ |
| 2 | minus | c | | $t_3$ |
| 3 | * | 6 | $t_3$ | $t_4$ |
| 4 | + | $t_2$ | $t_4$ | $t_5$ |
| 5 | = | $t_5$ | | a |

A triples has only 3 fields op, arg1 & arg2

Ex: The triples for $a = 6 * - c + 6 * - c$.

| | op | arg1 | arg2 |
|---|---|---|---|
| 0 | -minus | c | |
| 1 | * | 6 | (0) |
| 2 | minus | c | |
| 3 | * | 6 | (2) |
| 4 | + | (1) | (3) |
| 5 | = | a | (4) |

Static single-statement assignment form is an intermediate representation that facilitates certain code optimizations

Ex:

$$p = a + b$$
$$q = p - c$$
$$p = q * d$$
$$p = e - p$$
$$q = p + q$$

(a) 3-address code

$$p_1 = a + b$$
$$q_1 = p_1 - c$$
$$p_2 = q_1 * d$$
$$p_3 = e - p_2$$
$$q_2 = p_3 + q_1$$

(b) Static single assignment form.

2. Develop SDD to produce directed acyclic graph for an expression show the steps for constructing the DAG for the expression $a + a * (b-c) + (b-c) * d$.

Syntax directed definition is,

| $E \rightarrow$ PRODUCTION | SEMANTIC RULES |
|---|---|
| (i) $E \rightarrow E_1 + T$ | $E.node = new\ Node\ ('+', E_1.node, T.node)$ |
| (ii) $E \rightarrow E_1 - T$ | $E.node = new\ Node\ ('-', E_1.node, T.node)$ |
| (iii) $E \rightarrow T$ | $E.node = T.node$ |
| (iv) $T \rightarrow (E)$ | $T.node = E.node$ |
| (v) $T \rightarrow id$ | $T.node = new\ Leaf\ (id, id.entry)$ |
| (vi) $T \rightarrow num$ | $T.node = new\ Leaf\ (num, num.val)$ |

Steps for constructing the DAG

(i) $P_1$ = Leaf (id, entry-a)

(ii) $P_2$ = Leaf (id, entry-a) = $P_1$

(iii) $P_3$ = Leaf (id, entry-b)

(iv) $P_4$ = Leaf (id, entry-c)

(v) $P_5$ = Node ('−', $P_3$, $P_4$)

(vi) $P_6$ = Node ('×', $P_1$, $P_5$)

(vii) $P_7$ = Node ('+', $P_1$, $P_6$)

(ix) (viii) $P_8$ = Leaf (id, entry-b) = $P_3$

(ix) $P_9$ = Leaf (id, entry-c) = $P_4$

(x) $P_{10}$ = Node ('−', $P_3$, $P_4$) = $P_5$

(xi) $P_{11}$ = Leaf (id, entry-d)

(xii) $P_{12}$ = Node ('×', $P_5$, $P_{11}$)

(xiii) $P_{13}$ = Node ('+', $P_7$, $P_{12}$)

DAG

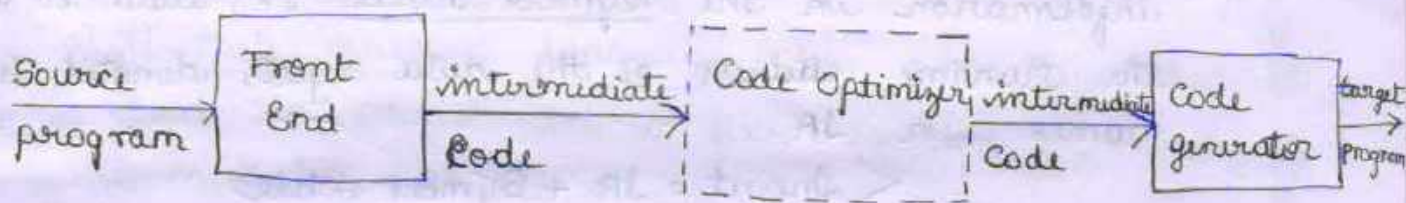# UNIT - 8
# CODE GENERATION

## INTRODUCTION :

* Code generation is the final phase in the compiler design.

* The code optimizer accepts intermediate code represent-ation which is generated from the front end of the compiler & produces another intermediate code representation which is Optimized.

* Code generator takes intermediate representation produced by code optimizer along with supplementary information in symbol table of the source program & produce as output an equivalent target program.



Source program → [Front End] → intermediate code → [Code Optimizer] → intermediate Code → [Code generator] → target program

* Code generator has 3 main tasks:
  1) Instruction selection
  2) Register allocation & assignment
  3) Instruction Ordering

## 1) INSTRUCTION SELECTION :

Choose a appropriate target machine instructions to implement the IR [intermediate represent] statements

## 2) REGISTER ALLOCATION & ASSIGNMENT :

Decide what values to keep in which registers

## 3) INSTRUCTION ORDERING

Decide in what order, to schedule the execution of instructions.

## 8.1 ISSUES IN THE DESIGN OF CODE GENERATOR:

1) Input to the code generator
2) The Target program
3) Instruction selection
4) Register Allocation
5) Evaluation Order

### 1) Input to the Code generator

* Input to the code generator is the intermediate represen of the source program produced by the front end along with information in the symbol table i.e., used to determine the runtime address of the data objects denoted by the names in IR.

$$< Input = IR + Symbol\ table >$$

* IR has several choices
  (a) 3-address representation: quadruples, triples, indirect triples
  (b) Virtual machine representation: byte codes & stack machine codes
  (c) Linear representation such as postfix notation
  (d) Graphical representation such as syntax trees & DAG's

* Assumptions made are
  (i) Front end produces low-level IR, i.e., values of names in it can be directly manipulated by the machine instruction.
  (ii) Syntatic & semantic errors have been already detected

2) The Target Program:

* The Output of code generator is <u>target program</u>.
* The <u>inst<sup>n</sup> set architecture</u> of the target machine has a significant impact on the design of code generator.
* Most common architecture are:

(a) <u>CISC</u>: It has few registers, has maximum of 2 operands & variety of addressing mode, variable length instructions & instruct<sup>n</sup> with side effects.

(b) <u>RISC</u>: It has many registers, has maximum of 3 operands with simple addressing modes, & relatively simple instruct<sup>n</sup> set architecture.

* Output may take variety of forms.

a) Absolute machine language [Executable Code]
b) Relocatable machine language [object files for linker]
c) Assembly language [facilitates debugging]

a) <u>Absolute machine language</u> has advantage that it can be placed in a fixed location in memory & immediately executed.

b) <u>Relocatable machine language</u> program allows subprograms to be compiled separately.

c) Producing <u>Assembly language</u> program as output makes the process of code generat<sup>n</sup> somewhat easier.

3) Instruction Selection

The code generator must map the IR program into a code sequence that can be executed by the target machine.

* The complexity of performing this mapping is determined by the factors such as:

   (i) the level of the IR
   (ii) the nature of the instruction set architectures.
   (iii) the desired quality of the generated code.

(i) the levels of the IR :

   > If the IR is high level , use code templates To translate each IR statements into a sequence of machine instruction.

   > produces poor code, needs further optimizat".

   > If the IR is low level , use code (this) low level informat" to generate more efficient code sequence.

(ii) the nature of the instruction set architectures has strong effect on difficulty of instruct" select".

   > Uniformity & completeness of the instruct" set are imp factors.

   > If we do not care about the efficiency of the target program, instruct" select" is straightforward.

   > For eg:

   $$x = y + z \Rightarrow \quad \begin{array}{ll} LD & R_0, y \\ ADD & R_0, R_0, z \\ ST & x, R_0 \end{array}$$

   ∴ produces redundant LD & store

   eg2:

   $$\begin{array}{l} a = b + c \\ d = a + b \end{array} \Rightarrow \quad \begin{array}{ll} LD & R_0, b \\ ADD & R_0, R_0, c \\ \boxed{\begin{array}{ll} ST & a, R_0 \\ LD & R_0, a \end{array}} \quad \rightarrow REDUNDANT \\ ADD & R_0, R_0, c \\ ST & d, R_0 \end{array}$$

(iii) the quality of the generated code is determined by its speed of size.

> For eg:

$$a = a + 1 \Rightarrow \begin{array}{l} LD\ R_0, a \\ ADD\ R_0, R_0, \#1 \\ ST\ a, R_0 \end{array} \left. \begin{array}{l} \\ \text{replaced} \\ \xrightarrow{\hspace{1cm}} \\ \text{by} \end{array} \right\} INC\ a$$

4) Register Allocation:

* Instruction involving register operands are usually shorter & faster than those involving operands in memory.

* 2 subproblems:

(i) Register allocation: Select the set of variables that will reside in registers at each point in the program.

(ii) Register assignment: Select specific register that a variable will reside in.

* Complications imposed by the hardware architecture
   Eg: Register pairs for multiplication & division.

* Multiplication instruction is of the form

$$\boxed{M \quad x, y}$$

where $x \longrightarrow$ Multiplicand, is the odd register of an even/odd register pair.

$y \longrightarrow$ Multiplier, is a single register.

$\Rightarrow$ Product $\longrightarrow$ occupies the entire even/odd register pair.

* Division instruction is of the form

$$\boxed{D \quad x, y}$$

where $x \longrightarrow$ dividend, occupies even register

$y \longrightarrow$ divisor, occupies odd/even register

$\Rightarrow$ quotient $\longrightarrow$ stored in odd register
   remainder $\longrightarrow$ stored in even register

Eg: two-3 address code sequences

$$t = a + b$$
$$t = t * c$$
$$t = t / d$$

$$t = a + b$$
$$t = t + c$$
$$t = t / d$$

Optimal machine - Code sequences

```
L   R1, a
A   R1, b
M   R0, C
D   R0, d
ST  R1, t
```

```
L    R0, a
A    R0, b
A    R0, c
SRDA R0, 32
D    R0, d
ST   R1, t
```

## 5) Evaluation Order:

* The order in which computations are performed can effect the efficiency of the Target code.

* when instruct^n are independent their evaluation order can be changed.

* Some computat^n orders require fewer registers to hold intermediate results than others.

* However picking a best order in the general case is a difficult NP-complete problem.

---

ADDITIONAL INFORMATION: Eg

$$a + b - (c+d)*e \Rightarrow$$

```
t1 = a+b
t2 = c+d
t3 = e*t2
t4 = t1-t3
```

$$\Rightarrow$$

```
MOV R0, a
ADD R0, b
MOV R1, R0
mov R1, c
ADD R1, d
mov R0, e
MUL R0, R1
mov R1, t1
SUB R1, R0
MOV t4, t4, R1
```

Reorder ↓

```
t2 = c+d
t3 = e*t2
t1 = a+b
t4 = t1-t3
```

```
mov R0, c
ADD R0, d
mov R1, e
MUL R1, R0
mov R0, a
ADD R0, b
SUB R0, R1
mov t4, R0
```

THE

8.2 THE TARGET LANGUAGE:

For designing a good code generator, we need to have familiarity with target machine & its instruction set. Instead of generating code on a specific target machine, a general machine consisting of many registers are considered.

A SIMPLE TARGET MACHINE MODEL:

The characteristics of target machine mode with instruction format & instruction set are shown below:

* Our hypothetical machine:

(i) It is a 3-address machine with the following format

$$\boxed{OP \quad destination, Source1, Source2}$$

NOTE:

A 3 address instruction can have 2 operands or 1 operands also but it can have max of 3 operands

(ii) The target machine is byte addressable i.e., it can access 8 bit of info from specific address

(iii) It has n no of registers denoted by $R_0, R_1, R_2, \ldots, R_{n-1}$

* Various types of instruction that are used by target m/c:

(i) Load Instruction

(ii) Store Instruction

(iii) Computational Instruction

(iv) Unconditional Instruction

(v) Conditional Instruction

(i) Load Instruction: Used to copy the data into destination operand which must be a register.

SYNTAX: LD dst, addr

where addr operand $\longrightarrow$ register or memory locat

(ii) Store instruction: Used to copy the data into memor location specified in the destination operand.

SYNTAX: ST dst, sr
where dst $\longrightarrow$ destination of st is a memo location
or $\longrightarrow$ register.

Computational operation.

(iii) Arithmetic instruction: They are performed using these instruction.

SYNTAX: OP dst, Src1, Src2.

where 1st operand, dst $\longrightarrow$ destination
2nd & 3rd operand $\longrightarrow$ Operands where R values fetched for Operat'' to be p

Eg1: ADD R0, R1, R2    // R0 = R1 + R2
Eg2: SUB R0, R0, R1    // R0 = R0 - R1
Eg3: MUL R2, R0, R1    // R2 = R0 * R1

(iv) Unconditional Jumps: The branch instruct'' without any condit'' are called unconditional jumps.

SYNTAX: BR label

where BR $\longrightarrow$ BRanch instruct''

(v) Conditional Jumps: Based on the value stored in a register i.e., whether it is true or zero or -ve, if branching takes place, then the branch inst'' are called Conditional jumps.

SYNTAX: Bcond or, label
where B stands from Branch,
Cond can be LT, GT, LTE, GTE
less than  greater than  less than or equal  greater than or equal

$\pi \rightarrow$ register, contains value such as 0, true or -ve.

eg1: B $\pi$    R0, T1      // Branch to T1, if R0 contains
                                                 -ve value
eg2: B $\pi$ TZ    R1, TR      // Branch to TR, if R1 contains
                                                 either 0 or -ve value

* Different addressing modes supported by generalized target machine:

     1) Direct addressing mode
     2) Indexed —"——————
     3) Integer Indexed —"————
     4) Indirect —"————
     5) Immediate —"————

(i) **Direct A/M :**

     Address of the data to be accessed is directly present in the instruct$^n$, i.e., location is identified by a variable name $x$.

     Eg: ~~LDP~~ LD    R1, $x$        // Load value stored in
                                       memory locat$^n$ $x$ into R1

(ii) **Indexed A/M:** The data can be accessed from a memory locat$^n$ using index. This addressing mode is useful for accessing arrays, where $a$ is the base address of the array & register holds the index value

     Eg: LD    R1, $a$ (R2)      // Accefs the data stored in
                                 $R_1 = $ Contents ($a$ + contents (R2))

(iii) **Indexed A/M** where memory locat$^n$ is integer

     It is same as previous one except that a memory locat$^n$ is identified as integer.

     Eg: LD    R1, 100(R2)      // R1 = contents (100+ contents (R2))

(iv) Indirect A/M : Contents of the data can be accessed by dereferencing using * operators as shown below:

$$\boxed{LD \quad R1, *(R2)}$$  // R2 contains memory locat the data stored in that memory locat$^n$ is copied in register R1

$$\boxed{LD \quad R1, *100(R2)}$$  // R1 = Contents ( contents( 100 + contents (R

(v) Immediate A/M : The data to be manipulated is directly present in the instruction & preceded by

$$\boxed{LD \quad R1, \#100}$$  // R1 ← 100

# EXERCISE :

1. Generate code for 3 address statement for $x = y - z$

| | | |
|---|---|---|
| LD | R1, y | // R1 = y |
| LD | R2, z | // R2 = z |
| ADD | R1, R1, R2 | // R1 = R1 + R2 |
| ST | x, R1 | // x = R1 |

2. Generate code for 3 address statement $x = *p$

| | | |
|---|---|---|
| LD | R1, p | // R1 ← p |
| LD | R2, 0(R1) | // R2 = Contents (0 + conte |
| ST | x, R2 | // x = R2 |

(iv) **Indirect A/M** : Contents of the data can be accessed by
dereferencing using * operators as shown below :

| LD   R1, *(R2) |    // R2 contains memory locat
the data stored in that
memory locat$^n$ is copied in
register R1

| LD   R1, *100 (R2) |    // R1 = Contents ( contents(100 +
contents(R

(iv) **Immediate A/M** : The data to be manipulated is
directly present in the instruction & preceded by

| LD   R1, #100 |    // R1 ← 100

## EXERCISE :

1. Generate code for, 3 address statement for $x = y - z$

$$LD \quad R1, y \qquad\qquad // R1 = y$$
$$LD \quad R2, z \qquad\qquad // R2 = z$$
$$ADD \quad R1, R1, R2 \qquad // R1 = R1 + R2$$
$$ST \quad x, R1 \qquad\qquad // x = R1$$

2. Generate code for, 3 address statement $x = *p$

$$LD \quad R1, p \qquad\qquad // R1 ← p$$
$$LD \quad R2, 0(R1) \qquad\qquad // R2 = Contents (0 + conte$$
$$ST \quad x, R2 \qquad\qquad // x = R2$$

3. Generate code for 3 address statement $*p = y$

   LD   R1, p       // R1 = p

   LD   R2, y       // R2 = y

   ST    0(R1), R2    // contents (0 + contents (R1)) = R2.

4. Generate m/c code for 3 address statement $b = a[i]$

   LD   R1, i        // R1 = i

   MUL   R1, R1, 8     // R1 = R1 * 8

   LD    R2, a[R1)     // R2 = contents (a + contents (R1))

   ST    b, R2       // b = R2

5. Generate m/c code for 3 address statement $a[j] = c$

   LD   R1, j        // R1 = j

   LD   R2, C.       // R2 = C

   MUL   R1, R1, 8     // R1 = R1 * C

   ST   a[R1], R2     // contents (a + contents (R1)) = R2

6. Generate m/c code for 3 address statement

    if $x < y$ goto l

   LD    R1, x       // R1 = x

   LD    R2, y       // R2 = y

   SUB   R1, R1, R2.    // R1 = R1 - R2

   BLTZ   R1, M      // if R1 < 0 jump to M.

# Program & Instruction Cost

* For simplicity we take the cost of an instruction to be one plus the costs associated with the addressing modes of the operands.

* A/M involves registers have zero additional cost.
* A/M involving memory locat^n or constant have additional cost of 1.
* For example:

  a) $\mathcal{L}$D     R0, R1          $\Rightarrow$ Cost = 1
  b) $\mathcal{L}$D     R0, M           $\Rightarrow$ Cost = 2
  c) $\mathcal{L}$D     R1, *100(RR)    $\Rightarrow$ Cost = 3

* Cost of Addressing mode:

| | Mode | Form | Address | Added Cost |
|---|---|---|---|---|
| ① | Absolute direct A/M | M | M | 1 |
| ② | Register direct A/M | R | R | 0 |
| ③ | Indexed A/M | C(R) | C + contents (R) | 1 |
| ④ | Indirect register A/M | *R | contents (R) | 0 |
| ⑤ | Indirect indexed A/M | *C(R) | contents (C + contents (R)) | 1 |
| ⑥ | Immediate A/M | #C | N/A | 1 |

NOTE: Cost of each statement = 1 + cost (Addressing mode)

# EXERCISES (8.2)

1. Determine the costs of the following instruction sequence

$$
\begin{aligned}
&\text{LD} \quad R0, y \qquad\qquad\qquad \text{Cost} = 1 + \text{cost(AM)}\\
&\text{LD} \quad R1, z \qquad\qquad\qquad \text{Cost} = 1 + 1 = 2\\
&\text{ADD} \quad R0, R0, R1 \qquad\;\; \text{Cost} = 1 + 1 = 2\\
&\text{ST} \quad x, R0 \qquad\qquad\qquad \text{Cost} = 1 + 0 = 1\\
&\qquad\qquad\qquad\qquad\qquad\quad\;\; \text{Cost} = 1 + 1 = 2
\end{aligned}
$$

Total Cost = 7

2.
```
LD   R0, i
MUL  R0, R0, 8
LD   R1, a(R0)
ST   b, R1
```

| | | Cost = Cost (AM) + 1. |
|---|---|---|
| LD | R0, i | Cost = 1 + 1 = 2 |
| MUL | R0, R0, 8 | Cost = 1 + 1 = 2 |
| LD | R1, a(R0) | Cost = 1 + 1 = 2. |
| ST | b, R1 | Cost = 1 + 1 = 2 |
| | | Total cost = 8 |

3.
```
LD   R0, C
LD   R1, i
MUL  R1, R1, 8
ST   a(R1), R0
```

| | Cost = cost (A.M) + 1 |
|---|---|
| LD  R0, C | 1 + 1 = 2 |
| LD  R1, i | 1 + 1 = 2 |
| MUL R1, R1, 8 | 1 + 1 = 2 |
| ST  a(R1), R0 | 1 + 1 = 2 |
| | Total cost = 8 |

4. LD    RO, p
   LD    RI, O(RO)
   ST    x, RI

| | Cost = Cost (A.M) + 1 |
|---|---|
| LD  RO, p | 1 + 1 = 2 |
| LD  RI, O(RO) | 1 + 1 = 2.   // ·: of constan |
| ST  x, RI | 1 + 1 = 2 |

Total  cost = 6

5. LD    RO, P
   LD    RI, x
   ST    O(RO), RI

| | Cost = Cost (A.M) + 1 |
|---|---|
| LD  RO, P | 1 + 1 = 2 |
| LD  RI, x | 1 + 1 = 2 |
| ST  O(RI), RI | 1 + 1 = 2 |

Total Cost = 6

6. LD    RO, x
   LD    RI, y
   SUB   RO, RO, RI
   BLTZ  *R3, RO

| | Cost = 1 + Cost (A.M) |
|---|---|
| LD  RO, x | 1 + 1 = 2 |
| LD  RI, y | 1 + 1 = 2 |
| SUB  RO, RO, RI | 1 + 0 = 1 |
| BLTZ  *R3, RO | 1 + 1 = 2   <·: indirect A |

Total  cost = 7