

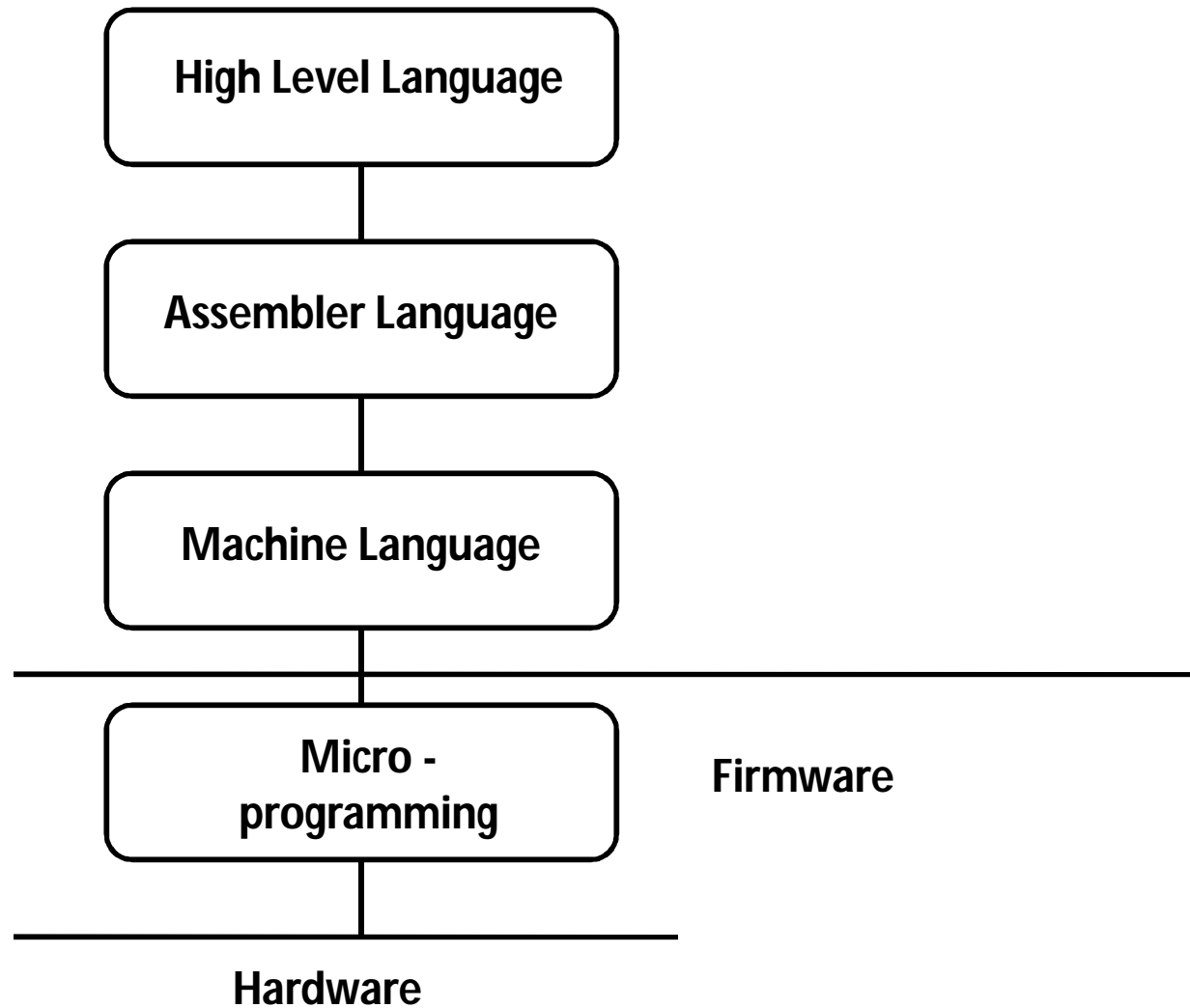
Assemblers

Assembler: Definition

- Translating source code written in assembly language to object code.



Language Levels



Machine code

- Machine code:
 - Set of commands directly executable via CPU
 - Commands in numeric code
 - Lowest semantic level

Machine code language

- Structure:

OpCode	OpAddress
--------	-----------

- Operation code
 - Defining executable operation
- Operand address
 - Specification of operands
 - Constants/register addresses/storage addresses

Elements of the Assembly Language Programming

- An Assembly language is a
 - machine dependent,
 - low level Programming language specific to a certain computer system.

Three features when compared with machine language are

1. Mnemonic Operation Codes
2. Symbolic operands
3. Data declarations

Elements of the Assembly Language Programming

Mnemonic operation codes: eliminates the need to memorize numeric operation codes.

Symbolic operands: Symbolic names can be associated with data or instructions. Symbolic names can be used as operands in assembly statements (need not know details of memory bindings).

Data declarations: Data can be declared in a variety of notations, including the decimal notation (avoids conversion of constants into their internal representation).

Assembly language-structure

<Label>	<Mnemonic>	<Operand>	Comments
----------------------	-------------------------	------------------------	-----------------

- Label
 - symbolic labeling of an assembler address (command address at Machine level)
- Mnemonic
 - Symbolic description of an operation
- Operands
 - Contains of variables or addresse if necessary
- Comments

Statement format

An Assembly language statement has following format:

[Label] <opcode> <operand spec>[,<operand spec>..]

If a label is specified in a statement, it is associated as a symbolic name with the memory word generated for the statement.

<operand spec> has the following syntax:

<symbolic name> [+<displacement>] [(<index register>)]

Eg. AREA, AREA+5, AREA(4), AREA+5(4)

Mnemonic Operation Codes

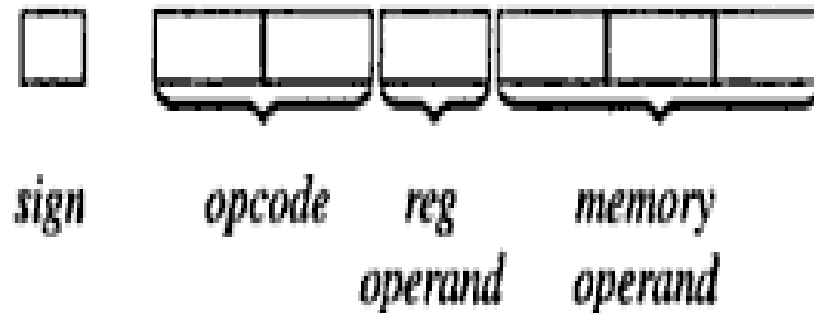
- Each statement has two operands, first operand is always a register and second operand refers to a memory word using a symbolic name and optional displacement.

<i>Instruction opcode</i>	<i>Assembly mnemonic</i>	<i>Remarks</i>
00	STOP	Stop execution
01	ADD	<i>First operand is modified Condition code is set</i>
02	SUB	
03	MULT	
04	MOVER	
05	MOVEM	Memory ← register move
06	COMP	Sets condition code
07	BC	Branch on condition
08	DIV	Analogous to SUB
09	READ	<i>First operand is not used</i>
10	PRINT	

Operation Codes

- *MOVE* instructions move a value between a memory word and a register
- *MOVER* – First operand is target and second operand is source
- *MOVEM* – first operand is source, second is target
- All arithmetic is performed in a register (replaces the contents of a register) and sets *condition code*.
- A Comparison instruction sets *condition code* analogous to arithmetics, i.e. without affecting values of operands.
- *condition code* can be tested by a Branch on Condition (BC) instruction and the format is:
BC <condition code spec> , <memory address>

Machine Instruction Format



- sign is not a part of the instruction
- Opcode: 2 digits, Register Operand: 1 digit, Memory Operand: 3 digits
- Condition code specified in a BC statement is encoded into the first operand using the codes 1- 6 for specifications LT, LE, EQ, GT, GE and ANY respectively
- In a Machine Language Program, all addresses and constants are shown in decimal as shown in the next slide

Example: ALP and its equivalent Machine Language Program

	START	101		
	READ	N	101)	+ 09 0 113
	MOVER	BREG, ONE	102)	+ 04 2 115
	MOVEM	BREG, TERM	103)	+ 05 2 116
AGAIN	MULT	BREG, TERM	104)	+ 03 2 116
	MOVER	CREG, TERM	105)	+ 04 3 116
	ADD	CREG, ONE	106)	+ 01 3 115
	MOVEM	CREG, TERM	107)	+ 05 3 116
	COMP	CREG, N	108)	+ 06 3 113
	BC	LE, AGAIN	109)	+ 07 2 104
	MOVEM	BREG, RESULT	110)	+ 05 2 114
	PRINT	RESULT	111)	+ 10 0 114
	STOP		112)	+ 00 0 000
N	DS	1	113)	
RESULT	DS	1	114)	
ONE	DC	'1'	115)	+ 00 0 001
TERM	DS	1	116)	
	END			

Assembly Language Statements

- An assembly program contains three kinds of statements:
 - 1) Imperative Statements
 - 2) Declaration Statements
 - 3) Assembler Directives

Imperative Statements: They indicate an action to be performed during the execution of an assembled program. Each imperative statement is translated into one machine instruction.

Assembly Language Statements

- **Declaration Statements:** syntax is as follows:

[Label] DS <constant>

[Label] DC '<value>'

- The DS (declare storage) statement reserves memory and associates names with them.
- Ex:
A DS 1 ; reserves a memory area of 1 word, associating the name A to it
G DS 200 ; reserves a block of 200 words and the name G is associated with the first word of the block (G+6 etc. to access the other words)
- The DC (declare constant) statement constructs memory words containing constants.
- Ex:
ONE DC '1' ; associates name one with a memory word containing value 1

Assembly Language Statements

Use of Constants

- The DC statement does not really implement constants
- it just initializes memory words to given values.
- The values are not protected by the assembler and can be changed by moving a new value into the memory word.
- In the above example, the value of ONE can be changed by executing an instruction

MOVEM BREG, ONE

Assembly Language Statements

Use of Constants

- An Assembly Program can use constants just like HLL, in two ways – as immediate operands, and as literals.
- 1) Immediate operands can be used in an assembly statement only if the architecture of the target machine includes the necessary features.
 - Ex: **ADD AREG,5**
 - This is translated into an instruction from two operands – AREG and the value '5' as an immediate operand

Assembly Language Statements

Use of Constants

- 2) A *literal* is an operand with the syntax = '<value>'.
- It differs from a constant because its location cannot be specified in the assembly program.
- Its value does not change during the execution of the program.
- It differs from an immediate operand because no architectural provision is needed to support its use.

ADD AREG, ='5'	➔	ADD AREG, FIVE
		FIVE DC '5'
Use of literals	vs.	Use of DC

Assembly Language Statements

Assembler Directive

- Assembler directives instruct the assembler to perform certain actions during the assembly of a program.
- Some assembler directives are described in the following:
 - 1) **START** *<constant>*
- This directive indicates that the first word of the target program generated by the assembler should be placed in the memory word having address *<constant>*.
- 2) **END** [*<operand spec>*]
- This directive indicates the end of the of the source program. The optional *<operand spec>* indicates the address of the instruction where the execution of the program should begin.

Advantages of Assembly Language

- The primary advantages of assembly language programming over machine language programming are due to the **use of symbolic operand specifications**.
(in comparison to machine language program)
- Assembly language programming holds an edge over HLL programming in situations where it is desirable to use architectural features of a computer.
(in comparison to high level language program)

Fundamentals of LP

- Language processing = analysis of source program + synthesis of target program
- **Analysis of source program** is specification of the source program
 - Lexical rules: formation of valid lexical units(tokens) in the source language
 - Syntax rules : formation of valid statements in the source language
 - Semantic rules: associate meaning with valid statements of the language

Fundamentals of LP

- **Synthesis of target program** is construction of target language statements
 - Memory allocation : generation of data structures in the target program
 - Code generation

A simple Assembly Scheme

- There are two phases in specifying an assembler:
 1. Analysis Phase
 2. Synthesis Phase(the fundamental information requirements will arise in this phase)

A simple Assembly Scheme

Design Specification of an assembler

There are four steps involved to design the specification of an assembler:

- Identify information necessary to perform a task.
- Design a suitable data structure to record info.
- Determine processing necessary to obtain and maintain the info.
- Determine processing necessary to perform the task

Synthesis Phase: Example

Consider the following statement:

MOVER BREG, ONE

The following info is needed to synthesize machine instruction for this stmt:

1. **Address of the memory word with which name **ONE** is associated** [depends on the source program, hence made available by the Analysis phase].
2. **Machine operation code corresponding to **MOVER**** [does not depend on the source program but depends on the assembly language, hence synthesis phase can determine this information for itself]

Note: Based on above discussion, the two data structures required during the synthesis phase are described next

Data structures in synthesis phase

Symbol Table --built by the analysis phase

- The two primary fields are name and address of the symbol used to specify a value.

Mnemonics Table --already present

- The two primary fields are *mnemonic* and *opcode*, along with *length*.

Synthesis phase uses these tables to obtain

- The machine address with which a name is associated.
 - The machine op code corresponding to a mnemonic.
- The tables have to be searched with the
 - **Symbol name and the mnemonic as keys**

Analysis Phase

- Primary function of the Analysis phase is to build the symbol table.
 - It must determine the addresses with which the symbolic names used in a program are associated
 - It is possible to determine some addresses directly like the address of first instruction in the program (ie.,start)
 - Other addresses must be inferred
 - To determine the addresses of the symbolic names we need to fix the addresses of all program elements preceding it through *Memory Allocation*.
- To implement *memory allocation* a data structure called *location counter* is introduced.

Analysis Phase – Implementing memory allocation

- LC(location counter) :
 - is always made to contain the address of the next memory word in the target program.
 - It is initialized to the constant specified at the START statement.
- When a LABEL is encountered,
 - it enters the LABEL and the contents of LC in a new entry of the symbol table.
 LABEL – e.g. N, AGAIN, SUM etc
 - It then finds the number of memory words required by the assembly statement and updates the LC contents
- To update the contents of the LC, analysis phase needs to know lengths of the different instructions
 - This information is available in the Mnemonics table and is extended with a field called length
- We refer the processing involved in maintaining the LC as LC Processing

Example

START 100

MOVER BREG, N

LC = 100 (1 byte)

MULT BREG, N

LC = 101 (1 byte)

STOP

LC = 102 (1 byte)

N DS 5

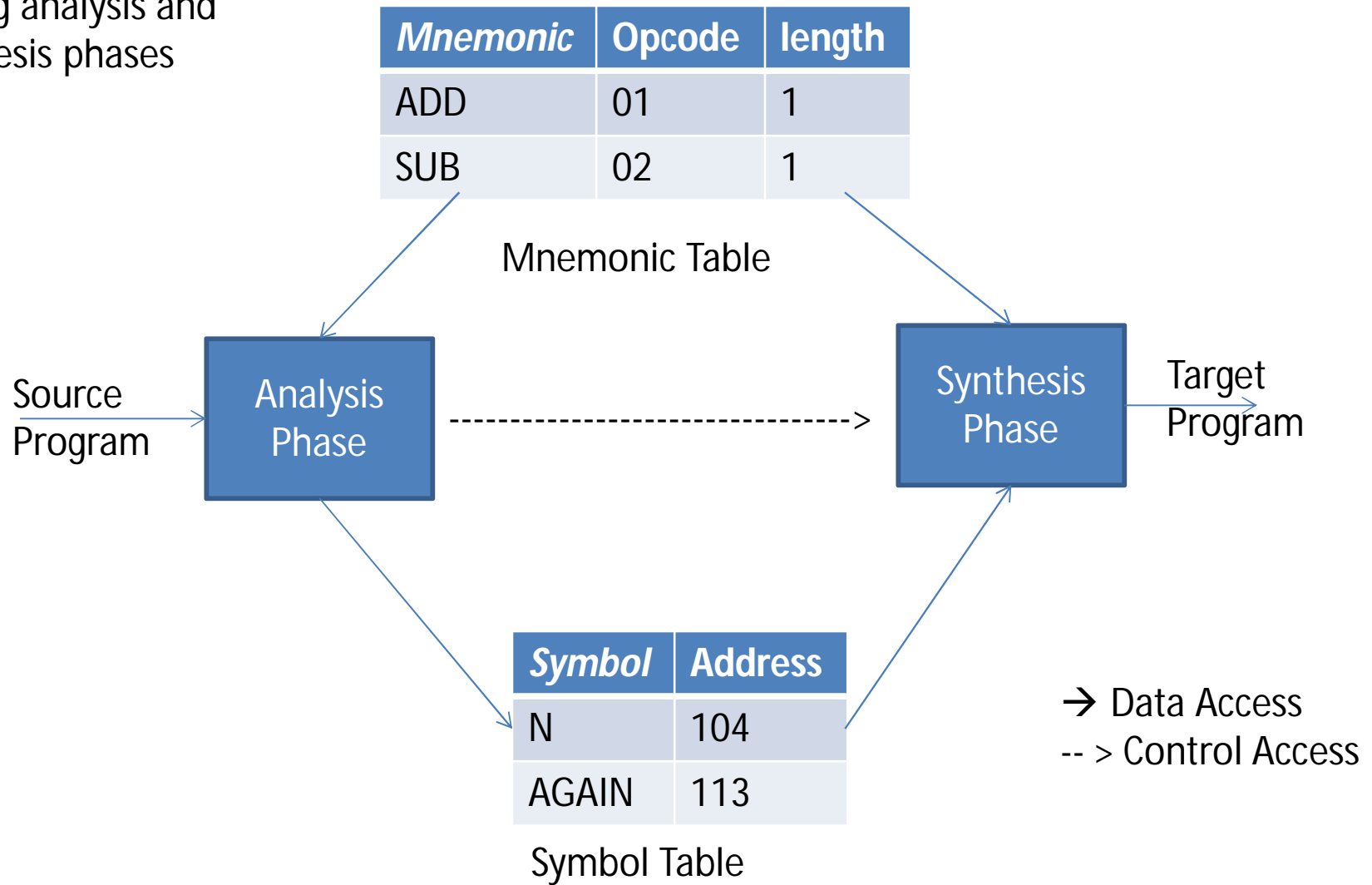
LC = 103

Symbol	Address
N	103

- Since there the instructions take different amount of memory, it is also stored in the mnemonic table in the “length” field

Mnemonic	Opcode	Length
MOVER	04	1
MULT	03	1

Data structures of an assembler
During analysis and
Synthesis phases



Data structures

- Mnemonics table is a fixed table which is merely accessed by the analysis and synthesis phases
- Symbol table is constructed during analysis and used during synthesis

Tasks Performed : Analysis Phase

- Isolate the labels, mnemonic, opcode and operand fields of a statement.
- If a label is present, enter (symbol, <LC>) into the symbol table.
- Check validity of the mnemonic opcode using mnemonics table.
- Update value of LC.

Tasks Performed : Synthesis Phase

- Obtain machine opcode corresponding to the mnemonic from the mnemonic table.
- obtain address of the memory operand from symbol table.
- Synthesize a machine instruction or machine form of a constant, depending on the instruction.

Assembler's functions

- Convert mnemonic operation codes to their machine language equivalents
- Convert symbolic operands to their equivalent machine addresses
- Build the machine instructions in the proper format
- Convert the data constants to internal machine representations
- Write the object program and the assembly listing

Assembler:Design

- The design of assembler can be of:
 - Scanning (tokenizing)
 - Parsing (validating the instructions)
 - Creating the symbol table
 - Resolving the forward references
 - Converting into the machine language


Assembler Design

- Pass of a language processor – one complete scan of the source program
- Assembler Design can be done in:
 - Single pass
 - Two pass
- Single Pass Assembler:
 - Does everything in single pass
 - Cannot resolve the forward referencing
- Two pass assembler:
 - Does the work in two pass
 - Resolves the forward references

Difficulties: Forward Reference

- Forward reference: reference to a label that is defined later in the program.

```
START 100  
MOVER BREG, N  
MULT BREG, N  
STOP  
N DS 5
```



- Problem of forward reference in single pass assembler is resolved by backpatching.

Backpatching

- The problem of forward references is handled using a process called backpatching
 - Initially, the operand field of an instruction containing a forward reference is left blank
 - Ex: MOVER BREG, ONE can be only partially synthesized since **ONE is a forward reference**
 - The instruction opcode and address of BREG will be assembled to reside in location 101
 - To insert the second operand's address later, an entry is added in Table of Incomplete Instructions (TII)
 - The entry TII is a pair (<instruction address>, <symbol>) which is (101, ONE) here

Backpatching

- The problem of forward references is handled using a process called backpatching
 - When END statement is processed, the symbol table would contain the addresses of all symbols defined in the source program
 - So TII would contain information of all forward references
 - Now each entry in TII is processed to complete the instruction
 - Ex: the entry (101, ONE) would be processed by obtaining the address of ONE from symbol table and inserting it in the operand field of the instruction with assembled address 101.
 - Alternatively, when definition of some symbol L is encountered, all forward references to L can be processed

Assembler Design

- First pass:
 - Scan the code by separating the symbol, mnemonic op code and operand fields
 - Perform LC processing
 - Record symbol and its address (LC) in to symbol table
 - Construct intermediate representation
- Second Pass:
 - Converts the code to the machine code

Assembler Directives

- Assembler directives are pseudo instructions.
 - They provide instructions to the assemblers itself.
 - They are not translated into machine operation codes. like START , END.

- Advanced Assembler directives are
 - **ORIGIN <address specification>**

Where <address specification> is <operand specification> or <constant> , this directive instruct the assembler to put <address specification> to location counter. Useful when lack of single contiguous area of memory.

1		START	200		
2		MOVER	AREG, ='5'	200)	+04 1 211
3		MOVEM	AREG, A	201)	+05 1 217
4	LOOP	MOVER	AREG, A	202)	+04 1 217
5		MOVER	CREG, B	203)	+05 3 218
6		ADD	CREG, ='1'	204)	+01 3 212
7		...			
12		BC	ANY, NEXT	210)	+07 6 214
13		LTORG			
			= '5'	211)	+00 0 005
			= '1'	212)	+00 0 001
14		...			
15	NEXT	SUB	AREG, ='1'	214)	+02 1 219
16		BC	LT, BACK	215)	+07 1 202
17	LAST	STOP		216)	+00 0 000
18		ORIGIN	LOOP+2		
19		MULT	CREG, B	204)	+03 3 218
20		ORIGIN	LAST+1		
21	A	DS	1	217)	
22	BACK	EQU	LOOP		
23	B	DS	1	218)	
24		END			

Example 3.4 (The ORIGIN directive) Statement number 18 of Figure 3.8(a), viz. **ORIGIN** LOOP+2, puts the address 204 in the location counter because the symbol LOOP is associated with the address 202. The next statement

MULT CREG, B

is therefore given the address 204. The statement **ORIGIN** LAST+1 puts the address 217 in the location counter. Note that an equivalent effect could have been achieved by using the statements **ORIGIN** 204 and **ORIGIN** 217 at these two places in the program; however, the absolute addresses used in these statements would have to be changed if the address specification in the START statement is changed.

EQU

EQU

The EQU directive has the syntax

<symbol> EQU <address specification>

where *<address specification>* is either a *<constant>* or *<symbolic name> ± <displacement>*. The EQU statement simply associates the name *<symbol>* with the address specified by *<address specification>*. However, the address in the location counter is not affected.

Example 3.5 (The EQU directive) Before processing the 22nd statement of Figure 3.8(a), the assembler had assembled the statement A DS 1 and given the address 217 to A. Hence the location counter contains the address 218 at this time. On encountering the statement BACK EQU LOOP, the assembler associates the symbol BACK with the address of LOOP, i.e., with 202. Note that the address in the location counter is not affected by the EQU statement. That is how the symbol B defined in the next statement B DS 1 is assigned the address 218. In the second pass, the 16th statement, i.e., BC LT, BACK, is assembled as '+ 07 1 202'.

LTORG

The **assembler** should put the values of literals in such a place that control does not reach any of them during execution of the generated program. The **LTORG directive**, which stands for '**origin** for literals', allows a programmer to specify where literals should be placed. The **assembler** uses the following scheme for placement of literals: When the use of a literal is seen in a statement, the **assembler** enters it into a *literal pool* unless a matching literal already exists in the pool. At every **LTORG** statement, as also at the **END** statement, the **assembler** allocates memory to the literals of the literal pool and clears the literal pool. This way, a literal pool would contain all literals used in the program since the start of the program or since the previous **LTORG** statement. Thus, all references to literals are forward references by definition. If a program does not use an **LTORG** statement, the **assembler** would enter all literals used in the program into a single pool and allocate memory to them when it encounters the **END** statement.

Example 3.6 (Memory allocation to literals) In Figure 3.8, the literals =‘5’ and =‘1’ are added to the literal pool in Statements 2 and 6, respectively. The first LTOrg statement (Statement 13) allocates the addresses 211 and 212 to the values ‘5’ and ‘1’. A new literal pool is now started. The value ‘1’ is put into this pool in Statement 15. This value is allocated the address 219 while processing the END statement. The literal =‘1’ used in Statement 15 therefore refers to location 219 of the second pool of literals rather than location 212 of the first pool.

Two Pass Assembler

<i>mnemonic opcode</i>	<i>class</i>	<i>mnemonic info</i>
MOVER	IS	(04,1)
DS	DL	R#7
START	AD	R#11
	:	

OPTAB

<i>symbol</i>	<i>address</i>	<i>length</i>
LOOP	202	1
NEXT	214	1
LAST	216	1
A	217	1
BACK	202	1
B	218	1

SYMTAB

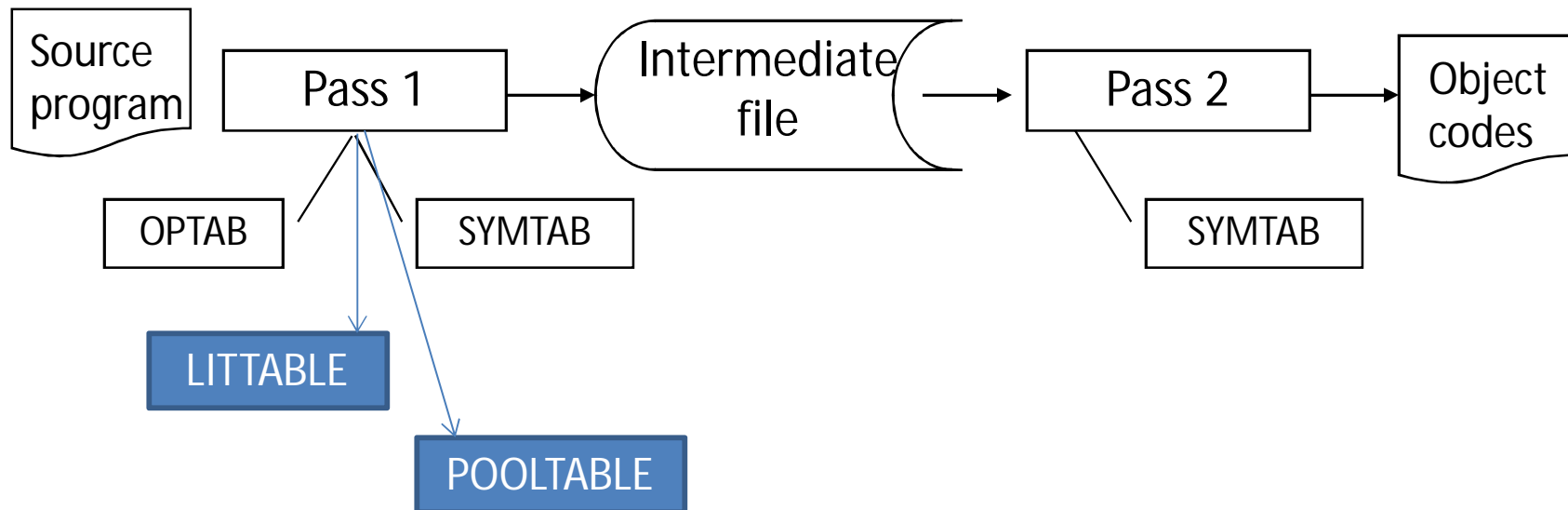
	<i>value</i>	<i>address</i>
1	= '5'	211
2	= '1'	212
3	= '1'	219

LITTAB

	<i>first</i>	<i># literals</i>
1	1	2
2	3	1
3	4	0

POOLTAB

Figure 3.9 Data structures of **assembler** Pass I



Data Structures in Pass I

- OPTAB – a table of mnemonic op codes
 - Contains mnemonic op code, class and mnemonic info
 - Class field indicates whether the op code corresponds to
 - an imperative statement (IS),
 - a declaration statement (DL) or
 - an assembler Directive (AD)
 - For IS, mnemonic info field contains the pair (machine opcode, instruction length)
 - Else, it contains the id of the routine to handle the declaration or a directive statement
 - The routine processes the operand field of the statement to determine the amount of memory required and updates LC and the SYMTAB entry of the symbol defined

Data Structures in Pass I

- SYMTAB - Symbol Table
 - Contains address and length
- LOCCTR - Location Counter
- LITTAB – a table of literals used in the program
 - Contains literal and address
 - Literals are allocated addresses starting with the current value in LC and LC is incremented, appropriately

OPTAB (operation code table)

- Content
 - Menmonic opcode, class and mnemonic info
- Characteristic
 - static table
- Implementation
 - array or hash table, easy for search

SYMTAB (symbol table)

- Content
 - label name, value, flag, (type, length) etc.
- Characteristic
 - dynamic table (insert, delete, search)
- Implementation
 - hash table, non-random keys, hashing function

IC

Variant I

Figure 3.12 shows an assembly program and its intermediate code using Variant I. The first operand in an assembly statement is represented by a single digit number which is either a code in the range 1...4 that represents a CPU register, where 1 represents AREG, 2 represents BREG, etc., or the condition code itself, which is in the range 1...6 and has the meanings described in Section 3.1. The second operand, which is a memory operand, is represented by a pair of the form

(operand class, code)

where *operand class* is one of C, S and L standing for constant, symbol and literal, respectively. For a constant, the *code* field contains the representation of the constant itself. For **example**, in Figure 3.12 the operand descriptor for the statement **START 200** is (C, 200). For a symbol or literal, the *code* field contains the entry number of the operand in SYMTAB or LITTAB. Thus entries for a symbol XYZ and a literal =‘25’ would be of the form (S, 17) and (L, 35), respectively.

	START	200	(AD, 01)	(C, 200)
	READ	A	(IS, 09)	(S, 01)
LOOP	MOVER	AREG, A	(IS, 04)	(1) (S, 01)
	⋮		⋮	
	SUB	AREG, =‘1’	(IS, 02)	(1) (L, 01)
	BC	GT, LOOP	(IS, 07)	(4) (S, 02)
	STOP		(IS, 00)	
A	DS	1	(DL, 02)	(C, 1)
	LTORG		(DL, 05)	
	

Variant II

Figure 3.13 shows an assembly program and its intermediate code using Variant II. This variant differs from Variant I in that the operand field of the intermediate code may be either in the processed form as in Variant I, or in the source form itself. For a declarative statement or an **assembler directive**, the operand field has to be processed in the first pass to support LC processing. Hence the operand field of its intermediate code would contain the processed form of the operand. For imperative statements, the operand field is processed to identify literal references and enter them in the LITTAB. Hence operands that are literals are represented as (L, *m*) in the intermediate code. There is no reason why symbolic references in operand fields of imperative statements should be processed during Pass I, so they are put in the source form itself in the intermediate code.

	START	200	(AD, 01)	(C, 200)
	READ	A	(IS, 09)	A
LOOP	MOVER	AREG, A	(IS, 04)	AREG, A
	⋮		⋮	
	SUB	AREG, = '1'	(IS, 02)	AREG, (L, 01)
	BC	GT, LOOP	(IS, 07)	GT, LOOP
	STOP		(IS, 00)	
A	DS	1	(DL, 02)	(C, 1)
	LTORG		(DL, 05)	
	

Figure 3.13 Intermediate code - Variant II