Mr. Abdul Rehman email id: abdul@northern.edu.pk

Whatsapp# 0308-7792217

(Week 11) Lecture 21-22

Learning objectives:

Review of the last lecture

General Register Organization

Stack Organization

Type of Instruction format

Introduction to instructions execution cycle

(Note: It is important to note that we are considering "Basic computer architecture" processor which we were discussing in the class and it's the continuation of the same concepts. Another important point is that brief part of the last lecture is repeated. Audio clip will be sent separately

for every diagram.)

Resources: Beside these lecture handouts, this lesson will draw from the following

Text Book: Computer System Architecture by Morris Mano (3rd Edition) and

Reference book: Computer Architecture, by William Stallings (4th Edition).

Lecture:

Register Stack Organization

This organization involves stack for performing the operations.

Stack

A useful feature that is included in the CPU of most computers is a stack or last-in, first-out

(UFO) list. A stack is a storage device that stores information in such a manner that the item

stored last is the first item retrieved. The operation of a stack can be compared to a stack of trays.

The last tray placed on top of the stack is the first to be taken off. The stack in digital computers

is essentially a memory unit with an address register that can count only (after an initial value is

loaded into it). The register that holds the address for the stack is called a stack pointer (SP)

1

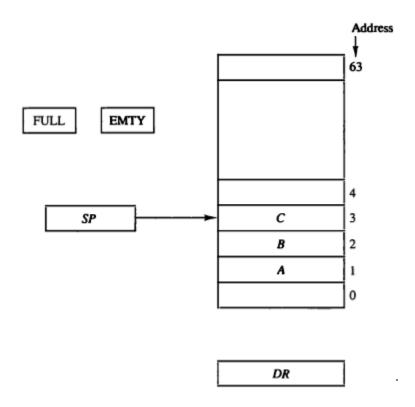
Mr. Abdul Rehman

email id: abdul@northern.edu.pk

Whatsapp# 0308-7792217

because its value always points at the top item in the stack. Contrary to a stack of trays where the tray itself may be taken out or inserted, the physical registers of a stack are always available for reading or writing. It is the content of the word that is inserted or deleted. The two operations of a stack are the insertion and deletion of items. The operation of insertion is called push (or pushdown) because it can be thought of as the result of pushing a new item on top. The operation of deletion is called pop (or pop-up) because it can be thought of as the result of removing one item so that the stack pops up. However, nothing is pushed or popped in a computer stack. These operations are simulated by incrementing or decrementing the stack pointer register.

- Pointer: SP
- Only PUSH and POP operations are applicable



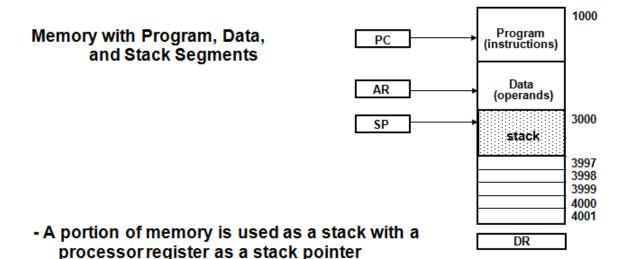
Push and Pop Operations

/* Initially, SP = 0, EMPTY = 1, FULL = 0 */

Mr. Abdul Rehman email id: abdul@northern.edu.pk Whatsapp# 0308-7792217

<u>PUSH</u>	<u>POP</u>
SP ← SP + 1	$DR \leftarrow M[SP]$
$M[SP] \leftarrow DR$	SP ← SP - 1
If (SP = 0) then (FULL \leftarrow 1)	If (SP = 0) then (EMPTY \leftarrow 1)
EMPTY ← 0	FULL ← 0

Register Stack Organization



- PUSH: SP ← SP - 1

M[SP] ← DR

- POP: DR ← M[SP] SP ← SP + 1

 Most computers do not provide hardware to check stack overflow (full stack) or underflow(empty stack)

Expression Evaluation (Reverse Polish Notation)

A stack organization is very effective for evaluating arithmetic expressions. The common mathematical method of writing arithmetic expressions imposes difficulties when evaluated by a computer. The common arithmetic expressions are written in infix notation, with each operator written between the operands. Consider the simple arithmetic expression

$$A*B+C*D$$

The star (denoting multiplication) is placed between two operands A and B or C and D. The plus is between the two products. To evaluate this arithmetic expression it is necessary to compute the

Mr. Abdul Rehman

email id: abdul@northern.edu.pk

Whatsapp# 0308-7792217

product A * B, store this product while computing C * D, and then sum the two products. From this example we see that to evaluate arithmetic expressions in infix notation it is necessary to scan back and forth along the expression to determine the next operation to be performed.

The Polish mathematician Lukasiewicz showed that arithmetic expressions can be represented in prefix notation. This representation, often referred to as Polish notation, places the operator before the operands. The postfix notation, referred to as reverse Polish notation (RPN), places the operator after the operands. The following examples demonstrate the three representations:

- A+ B Infix notation
- + AB Prefix or Polish notation
- AB+ Postfix or reverse Polish notation

The reverse Polish notation is in a form suitable for stack manipulation. The expression

$$A*B+C*D$$

is written in reverse Polish notation as

and is evaluated as follows: Scan the expression from left to right. When an operator is reached, perform the operation with the two operands found on the left side of the operator. Remove the two operands and the operator and replace them by the number obtained from the result of the operation. Continue to scan the expression and repeat the procedure for every operator encountered until there are no more operators.

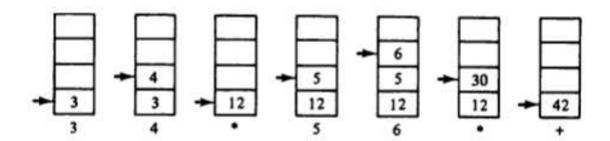
Evaluation of Arithmetic

Expressions Reverse Polish notation, combined with a stack arrangement of registers, is the most efficient way known for evaluating arithmetic expressions. This procedure is employed in some electronic calculators and also in some computers. The stack is particularly useful for handling long, complex problems involving chain calculations. It is based on the fact that any arithmetic expression can be expressed in parentheses-free Polish notation.

Mr. Abdul Rehman email id: abdul@northern.edu.pk Whatsapp# 0308-7792217

The following numerical example may clarify this procedure. Consider the arithmetic expression (3 * 4) + (5 * 6)

In reverse Polish notation, it is expressed as



Instruction formats

The physical and logical structure of computers is normally described in reference manuals provided with the system. Such manuals explain the internal construction of the CPU, including the processor registers available and their logical capabilities. They list all hardware-implemented instructions, specify their binary code format, and provide a precise definition of each instruction. A computer will usually have a variety of instruction code formats. It is the function of the control unit within the CPU to interpret each instruction code and provide the necessary control functions needed to process the instruction.

The format of an instruction is usually depicted in a rectangular box symbolizing the bits of the instruction as they appear in memory words or in a control register. The bits of the instruction are divided into groups called fields. The most common fields found in instruction formats are:

- 1. An operation code field that specifies the operation to be performed.
- 2. An address field that designates a memory address or a processor registers.
- 3. A mode field that specifies the way the operand or the effective address is determined.

Computers may have instructions of several different lengths containing varying number of addresses. The number of address fields in the instruction format of a computer depends on the internal organization of its registers. Most computers fall into one of three types of CPU organizations:

- 1. Single accumulator organization.
- 2. General register organization.
- 3. Stack organization.

Mr. Abdul Rehman email id: abdul@northern.edu.pk Whatsapp# 0308-7792217

To illustrate the influence of the number of addresses on computer programs, we will evaluate the arithmetic statement

$$X = (A + B) * (C + D)$$

Three-Address Instructions

Computers with three-address instruction formats can use each address field to specify either a processor register or a memory operand. The program in assembly language that evaluates $X = (A + B) \cdot (C + D)$ is shown below.

ADD R1, A, B ; $R1 \leftarrow M[A] + M[B]$

ADD R2, C, D ; $R2 \leftarrow M[C] + M[D]$

MUL X, R1, R2 ; $X \leftarrow R1 * R2$

It is assumed that the computer has two processor registers, R 1 and R2. The symbol M[A] denotes the operand at memory address symbolized by A.

The advantage of the three-address format is that it results in short programs when evaluating arithmetic expressions. The disadvantage is that the binary-coded instructions require too many bits to specify three addresses.

Two-Address Instructions

Two-address instructions are the most common in commercial computers. Here again each address field can specify either a processor register or a memory word. The program to evaluate X = (A + B) * (C + D) is as follows:

MOV R1, A; $R1 \leftarrow M[A]$

ADD R1, B ; $R1 \leftarrow R1 + M[B]$

MOV R2, C; $R2 \leftarrow M[C]$

ADD R2, D ; $R2 \leftarrow R2 + M[D]$

MUL R1, R2 ; R1 \leftarrow R1 * R2

MOV X, R1 ; $M[X] \leftarrow R1$

The MOV instruction moves or transfers the operands to and from memory and processor registers. The first symbol listed in an instruction is assumed to be both a source and the destination where the result of the operation is transferred.

One-address instructions

Mr. Abdul Rehman email id: abdul@northern.edu.pk Whatsapp# 0308-7792217

One-address instructions use an implied accumulator (AC) register for all data manipulation. For multiplication and division there is a need for a second register. However, here we will neglect the second register and assume that the AC contains the result of all operations. The program to evaluate X = (A + B) * (C + D) is

LOAD A ; $AC \leftarrow M[A]$

ADD B ; $AC \leftarrow AC + M[B]$

STORE T ; $M[T] \leftarrow AC$

LOAD C ; $AC \leftarrow M[C]$

ADD D ; $AC \leftarrow AC + M[D]$

MUL T ; $AC \leftarrow AC*M[T]$

STORE X ; $M[X] \leftarrow AC$

Zero-Address Instructions

A stack-organized computer does not use an address field for the instructions ADD and MUL. The PUSH and POP instructions, however, need an address field to specify the operand that communicates with the stack. The following program shows how X = (A + B) * (C + D) will be written for a stack organized computer. (TOS stands for top of stack.)

PUSH A ; TOS←A

PUSH B ; TOS←B

ADD ; $TOS \leftarrow (A+B)$

PUSH C ; TOS \leftarrow C

PUSH D ; TOS←D

ADD ; $TOS \leftarrow (C+D)$

MUL ; $TOS \leftarrow (C + D) * (A+B)$

POP X ; $M[X] \leftarrow TOS$

To evaluate arithmetic expressions in a stack computer, it is necessary to convert the expression into reverse Polish notation. The name "zero-address" is given to this type of computer because of the absence of an address field in the computational instructions.

Introduction

Mr. Abdul Rehman email id: abdul@northern.edu.pk

Whatsapp# 0308-7792217

An important aspect of computer architecture is the design of the instruction set for the

processor. The instruction set chosen for a particular computer determines the way that machine

language programs are constructed. Early computers had small and simple instruction sets,

forced mainly by the need to minimize the hardware used to implement them.

In this phase of the lecture we will study two computer architectures

• Complex Instruction Set Computers (CISC)

• Reduced Instruction Set Computer (RISC)

CISC Characteristics:

The design of an instruction set for a computer must take into consideration not only machine

language constructs, but also the requirements imposed on the use of high-level programming

languages. The translation from high-level to machine language programs is done by means of a

compiler program. One reason for the trend to provide a complex instruction set is the desire to

simplify the compilation and improve the overall computer performance. The task of a compiler

is to generate a sequence of machine instructions for each high-level language statement.

The essential goal of CISC architecture is to attempt to provide a single machine instruction for

each statement that is written in a high-level language. The instructions in a typical CISC

processor provide direct manipulation of operands residing in memory. CISC processors have

instructions that use only processor registers, the availability of other modes of operations tend to

simplify high-level language compilation. In summary, the major characteristics of CISC

architecture are:

1. A large number of instructions-typically from 100 to 250 instructions

2. Some instructions that perform specialized tasks and are used infrequently

3. A large variety of addressing modes-typically from 5 to 20 different modes

4. Variable-length instruction formats

5. Instructions that manipulate operands in memory

RISC Characteristics:

8

Mr. Abdul Rehman

email id: abdul@northern.edu.pk

Whatsapp# 0308-7792217

The concept of RISC architecture involves an attempt to reduce execution time by simplifying the instruction set of the computer. The small set of instructions of a typical RISC processor consists mostly of register-to-register operations, with only simple load and store operations for memory access. Thus each operand is brought into a processor register with a load instruction. All computations are done among the data stored in processor registers. Results are transferred to memory by means of store instructions. This architectural feature simplifies the instruction set and encourages the optimization of register manipulation. The use of only a few addressing modes results from the fact that almost all instructions have simple register addressing. Other addressing modes may be included, such as immediate operands and relative mode. The major characteristics of a RISC processor are:

- 1. Relatively few instructions
- 2. Relatively few addressing modes
- 3. Memory access limited to load and store instructions
- 4. All operations done within the registers of the CPU
- 5. Fixed-length, easily decoded instruction format Single-cycle instruction execution
- 6. Hardwired rather than microprogrammed control
- 7. Use of overlapped register windows to speed-up procedure call and return

Overlapped Register:

Windows Procedure call and return occurs quite often in high-level programming languages. When translated into machine language, a procedure call produces a sequence of instructions that save register values, pass parameters needed for the procedure, and then calls a subroutine to execute the body of the procedure. After a procedure return, the program restores the old register values, passes results to the calling program, and returns from the subroutine. Saving and restoring registers and passing of parameters and results involve time consuming operations. Some computers provide multiple-register banks, and each procedure is allocated its own bank of registers. This eliminates the need for saving and restoring register values. Some computers use

Mr. Abdul Rehman

email id: abdul@northern.edu.pk

Whatsapp# 0308-7792217

the memory stack to store the parameters that are needed by the procedure, but this requires a memory access every time the stack is accessed.

The concept of overlapped register windows is illustrated in Figure below. The system has a total of 74 registers. Registers R0 through R9 are global registers that hold parameters shared by all procedures. The other 64 registers are divided into four windows to accommodate procedures A, B, C, and D. Each register window consists of 10 local registers and two sets of six registers common to adjacent windows. Local registers are used for local variables. Common registers are used for exchange of parameters and results between adjacent procedures. The common overlapped registers permit parameters to be passed without the actual movement of data. Only one register window is activated at any given time with a pointer indicating the active window. Each procedure call activates a new register window by incrementing the pointer. The high registers of the calling procedure overlap the low registers of the called procedure, and therefore the parameters automatically transfer from calling to called procedure.

As mentioned previously, the 10 global registers R0 through R9 are available to all procedures. Each procedure in Figure below has available a total of 32 registers while it is active. This includes 10 global registers, 10 local registers, six low overlapping registers, and six high overlapping registers. Other fixed size register window schemes are possible, and each may differ in the size of the register window and the size of the total register file. In general, the organization of register windows will have the following relationships:

- Number of global registers = G
- Number of local registers in each window = L
- Number of registers common to two windows = C
- Number of windows = W

Mr. Abdul Rehman email id: abdul@northern.edu.pk Whatsapp# 0308-7

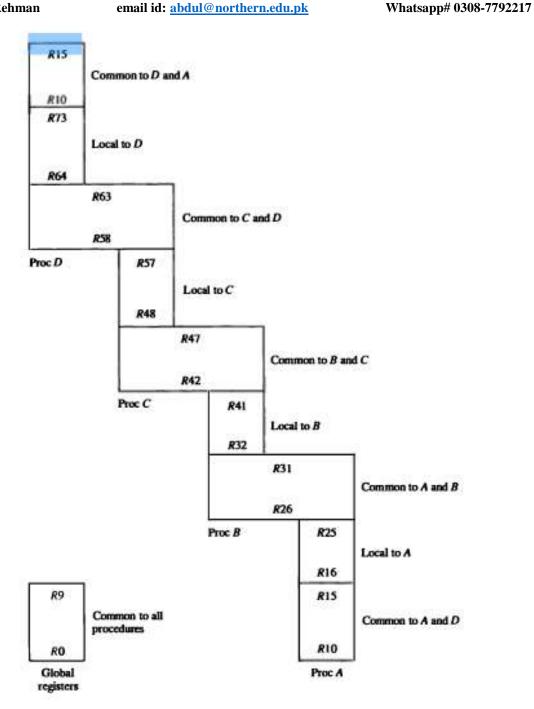


Figure: Overlapped register Window

Parallel processing:

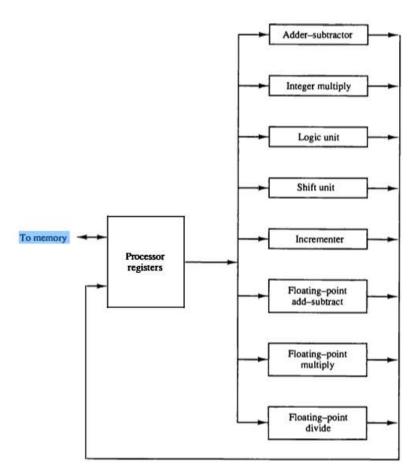
Parallel processing is a term used to denote a large class of techniques that are used to provide simultaneous data processing tasks for the purpose of increasing the computational speed of a computer system. Instead of processing each instruction sequentially as in a conventional

Mr. Abdul Rehman

email id: abdul@northern.edu.pk

Whatsapp# 0308-7792217

computer, a parallel processing system is able to perform concurrent data processing to achieve faster execution time. For example, while an instruction is being executed in the ALU, the next instruction can be read from memory. The system may have two or more ALUs and be able to execute two or more instructions at the same time. Furthermore, the system may have two or more processors operating concurrently. The purpose of parallel processing is to speed up the computer processing capability and increase its throughput, that is, the amount of processing that can be accomplished during a given interval of time. The amount of hardware increases with parallel processing and with it, the cost of the system increases. However, technological developments have reduced hardware costs to the point where parallel processing techniques a.re economically feasible.



One classification introduced by M. J. Flynn considers the organization of a computer system by the number of instructions and data items that are manipulated simultaneously. The normal

Mr. Abdul Rehman email id: abdul@northern.edu.pk

Whatsapp# 0308-7792217

operation of a computer is to fetch instructions from memory and execute them in the processor.

The sequence of instructions read from memory constitutes an instruction stream. The operations

performed on the data in the processor constitute a data stream. Parallel processing may occur in

the instruction stream, in the data stream, or in both. Flynn's classification divides computers into

four major groups as follows:

• Single instruction stream, single data stream (SISD)

• Single instruction stream, multiple data stream (SIMD)

• Multiple instruction stream, single data stream (MISD)

• Multiple instruction stream, multiple data stream (MIMD)

SISD represents the organization of a single computer containing a control unit, a processor unit,

and a memory unit. Instructions are executed sequentially and the system may or may not have

internal parallel processing capabilities. Parallel processing in this case may be achieved by

means of multiple functional units or by pipeline processing.

SIMD represents an organization that includes many processing units under the supervision of a

common control unit. All processors receive the same instruction from the control unit but

operate on different items of data. The shared memory unit must contain multiple modules so

that it can communicate with all the processors simultaneously.

MISD structure is only of theoretical interest since no practical system has been constructed

using this organization.

MIMD organization refers to a computer system capable of processing several programs at the

same time. Most multiprocessor and multicomputer systems can be classified in this category.

Pipelining

Pipelining is a technique of decomposing a sequential process into suboperations, with each

subprocess being executed in a special dedicated segment that operates concurrently with all

other segments. A pipeline can be visualized as a collection of processing segments through

which binary information flows. Each segment performs partial processing dictated by the way

the task is partitioned. The result obtained from the computation in each segment is transferred to

13

Mr. Abdul Rehman email id: abdul@northern.edu.pk Whatsapp# 0308-7792217

the next segmenting the pipeline. The final result is obtained after the data have passed through all segments. The name "pipeline" implies a flow of information analogous to an industrial assembly line. It is characteristic of pipelines that several computations can be in progress in distinct segments at the same time. The overlapping of computation is made possible by associating a register with each segment in the pipeline. The registers provide isolation between each segment so that each can operate on distinct data simultaneously.

The pipeline organization will be demonstrated by means of a simple example. Suppose that we want to perform the combined multiply and add operations with a stream of numbers.

$$A_i * B_i + C_i$$
 for $i = 1, 2, 3, ..., 7$

Each suboperation is to be implemented in a segment within a pipeline. Each segment has one or two registers and a combinational circuit as shown in Fig. 9-2. R 1 through RS are registers that receive new data with every clock pulse. The multiplier and adder are combinational circuits. The suboperations performed in each segment of the pipeline are as follows:

 $R1 \leftarrow A_i, R2 \leftarrow B_i$ Input A, and B

 $R3 \leftarrow R1*R2$, $R4 \leftarrow C_i$ Multiply and input C

 $R5 \leftarrow R3 + R4$ Add C_i to product

The five registers are loaded with new data every clock pulse. The effect of each clock is shown in Table 9-1. The first clock pulse transfers A1 and B1 into R1 and R2. The second dock pulse transfers the product of R1 and R2 into R3 and C1 into R4. The same clock pulse transfers A2 and B2 into R1 and R2. The third clock pulse operates on all three segments simultaneously. It places A, and B, into R1 and R2, transfers the product of R1 and R2 into R3, transfers C, into R4, and places the sum of R3 and R4 into RS. It takes three clock pulses to fill up the pipe and retrieve the first output from RS. From there on, each dock produces a new output and moves the data one step down the pipeline. This happens as long as new input data flow into the system. When no more input data are available, the clock must continue until the last output emerges out of the pipeline.

Mr. Abdul Rehman

email id: abdul@northern.edu.pk

Whatsapp# 0308-7792217

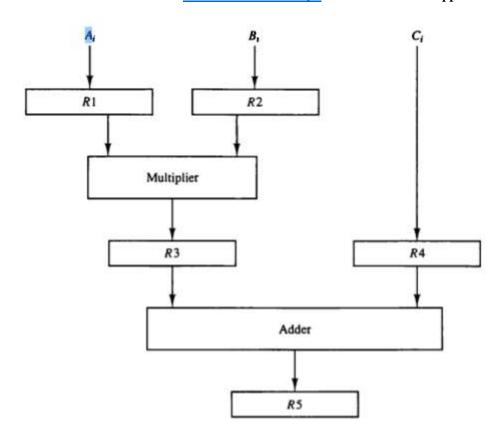


TABLE 9-1 Content of Registers in Pipeline Example

Clock Pulse Number	Segment 1		Segment 2		Segment 3
	R1	R2	R3	R4	R5
1	<i>A</i> ₁	<i>B</i> ₁	_	_	_
2	A_2	B_2	$A_1 * B_1$	C_1	_
3	A_3	B_3	$A_2 * B_2$	C_2	$A_1*B_1+C_1$
4	A_4	B_4	$A_3 * B_3$	C_3	$A_2*B_2+C_2$
5	A_5	B_5	$A_4 * B_4$	C_4	$A_3*B_3+C_3$
6	A_6	B_6	$A_5 * B_5$	C_{5}	$A_4*B_4+C_4$
7	A_7	B_7	$A_6 * B_6$	C_6	$A_5*B_5+C_5$
8	_	_	$A_7 * B_7$	C_7	$A_6*B_6+C_6$
9	_	_	_	_	$A_7 * B_7 + C_7$