

OBJECTIVES:

- To learn the basic structure and operations of a computer.
- To learn the arithmetic and logic unit and implementation of fixed-point and floating point arithmetic unit.
- To learn the basics of pipelined execution.
- To understand parallelism and multi-core processors.
- To understand the memory hierarchies, cache memories and virtual memories.
- To learn the different ways of communication with I/O devices.

UNIT I BASIC STRUCTURE OF A COMPUTER SYSTEM 9

Functional Units – Basic Operational Concepts – Performance – Instructions: Language of the Computer – Operations, Operands – Instruction representation – Logical operations – decision making – MIPS Addressing.

UNIT II ARITHMETIC FOR COMPUTERS 9

Addition and Subtraction – Multiplication – Division – Floating Point Representation – Floating Point Operations – Subword Parallelism

UNIT III PROCESSOR AND CONTROL UNIT 9

A Basic MIPS implementation – Building a Datapath – Control Implementation Scheme – Pipelining – Pipelined datapath and control – Handling Data Hazards & Control Hazards – Exceptions.

UNIT IV PARALLELISIM 9

Parallel processing challenges – Flynn's classification – SISD, MIMD, SIMD, SPMD, and Vector Architectures - Hardware multithreading – Multi-core processors and other Shared Memory Multiprocessors - Introduction to Graphics Processing Units, Clusters, Warehouse Scale Computers and other Message-Passing Multiprocessors.

UNIT V MEMORY & I/O SYSTEMS 9

Memory Hierarchy - memory technologies – cache memory – measuring and improving cache performance – virtual memory, TLB's – Accessing I/O Devices – Interrupts – Direct Memory Access – Bus structure – Bus operation – Arbitration – Interface circuits - USB.

TOTAL : 45 PERIODS

OUTCOMES:

On Completion of the course, the students should be able to:

- Understand the basics structure of computers, operations and instructions.
- Design arithmetic and logic unit.
- Understand pipelined execution and design control unit.
- Understand parallel processing architectures.
- Understand the various memory systems and I/O communication.

TEXT BOOKS:

1. David A. Patterson and John L. Hennessy, Computer Organization and Design: The Hardware/Software Interface, Fifth Edition, Morgan Kaufmann / Elsevier, 2014.
2. Carl Hamacher, Zvonko Vranesic, Safwat Zaky and Naraig Manjikian, Computer Organization and Embedded Systems, Sixth Edition, Tata McGraw Hill, 2012.

REFERENCES:

1. William Stallings, Computer Organization and Architecture – Designing for Performance, Eighth Edition, Pearson Education, 2010.
2. John P. Hayes, Computer Architecture and Organization, Third Edition, Tata McGraw Hill, 2012.
3. John L. Hennessey and David A. Patterson, Computer Architecture – A Quantitative Approach, Morgan Kaufmann / Elsevier Publishers, Fifth Edition, 2012.

Arunai Engineering College

UNIT - I

Part - A

TWO MARKS WITH ANSWERS

1. List the major components of a computer system. [APR/MAY]
What are components of a computer system? [NOV/DEC-17]

* The major components of a computer systems.

1. Computer Hardware

a) Input unit

b) output unit

c) processing unit (CPU)

i) ALU ii) control unit

iii) Registers

* Data Registers

* Address Registers

* Status Registers

* General purpose Registers

* Program Counter.

d) Memory unit.

* Main Memory

↳ RAM

↳ ROM

* Secondary Memory

2. Computer Software

a) System software

b) Application software

* General purpose Software

* Special purpose Software

2.

What are the addressing modes? [NOV/DEC-17]

* Addressing modes refers to the way in which the location of an operand is specified in an instruction is referred to as addressing modes.

* Information contained in the instruction code is the value of the operand or the address of the result/operand.

Different types of addressing modes are

1. Immediate Addressing mode
2. Register Addressing mode
3. Register Indirect Addressing mode
4. Direct Addressing mode
5. Indirect Addressing mode
6. Implied Addressing mode
7. Relative Addressing mode
8. Indexed Addressing mode
9. Base Register Addressing mode
10. Autoincrement or Autodecrement Addressing mode.

3. State the need for indirect addressing mode.

Give an Example. [APR/MAY-17]

* The effective address of the operand is the contents of a register or main memory location, location whose address appears in the instruction.

* The register or memory location that contains the address of the operand is a pointer.

* When an execution takes place in such mode, instruction may be told to go to a specific address.

* Once it's there, instead of finding an operand, it finds an address where the operand is located.

Eg: $\text{J } \$\text{SI} \quad // \text{SI} = 4008(\text{Address})$

4. What is an instruction register? [NOV/DEC-16]

* An Instruction Register (IR) is the part of a CPU control unit that holds the instruction currently being executed or decoded.

* Decoding the OP-code in the instruction register includes determining the instruction, determining where its operands are in memory, retrieving the operands from memory, allocating processor resources to execute the command etc.

5.

Give the formula for CPU execution time for a program. [NOV/DEC-2016]

* CPU execution time (CPU time) - Time the CPU spends working on a task.

$$\text{CPU Execution time for a program} = \text{CPU clock cycles for a program} \times \text{Clock cycle time}$$

6.

How to represent instruction in a computer system? [MAY/June - 2016]

* Instructions are kept in the computer as a series of high and low electronic signals and may be represented as numbers.

* Since registers are referred in instructions, there must be a convention to map register names into numbers.

* Three types of instruction formats are used in MIPS they are

i) R-Type or R-Format (Registers)

ii) I-Type or I-Format (Immediate)

iii) J-Type (Addressing in Branches and Jumps)

7. Distinguish between auto increment and auto decrement addressing mode. [MAY/JUNE-2016]

AUTO INCREMENT-Addressing mode	Auto Decrement-Addressing mode
1. The Effective address of the operand is the contents of a register in the instruction.	1. The Effective address of the operand is the content of a register in the instruction.
2. After accessing the operand the contents of this register is <u>automatically incremented</u>	2. After accessing the operand the contents of this register is <u>automatically decremented</u> .
3. Syntax: $(R_i) +$	3. Syntax, $-(R_i)$
4. Eg: $EA = (R_1) +$	4. Eg: $EA = -(R_1)$

8. What is Instruction Set Architecture? [Nov/Dec-2015]

* The Instruction Set Architecture (ISA) or simply architecture of a computer is the interface between the hardware and the lowest-level software.

* It includes anything programmers need to know to make a binary machine language program work correctly, including instructions, I/O devices, and so on.

9. How CPU execution time for a program is calculated? [NOV/DEC - 2015]

* CPU execution time or simply CPU time, which recognizes this distinction, is the time the CPU spends computing for this task and does not include time spent waiting for I/O or running other programs.

* CPU time can be further divided into

- i) User CPU time - CPU time spent in the program
- ii) System CPU time - CPU time spent in the operating system performing tasks on behalf of the program.

10. List the eight great ideas invented by Computer Architects. [APR/MAY - 2015]

- i) Design for Moore's Law
- ii) Use Abstraction to Simplify Design
- iii) Make the Common Case Fast
- iv) Performance via Parallelism
- v) Performance via Pipelining
- vi) Performance via Prediction
- vii) Hierarchy of Memories
- viii) Dependability via Redundancy.

11. Distinguish Pipelining from Parallelism. [APR/MAY-2015]

Pipelining	Parallelism
<p>i) Pipelining is a technique in which multiple instructions are overlapped in execution.</p> <p>ii) It exploits parallelism among the instructions in a sequential instruction stream.</p>	<p>i) Performance can be improved by doing operations in parallel.</p> <p>ii) Parallel execution is easier when tasks are independent, but often they need to co-operate with each other.</p>

12. State Amdahl's Law. [NOV/DEC-2014]

* Amdahl's Law is used to find the execution time of a program after making the improvement.

* It can be represented in an equation as follows:

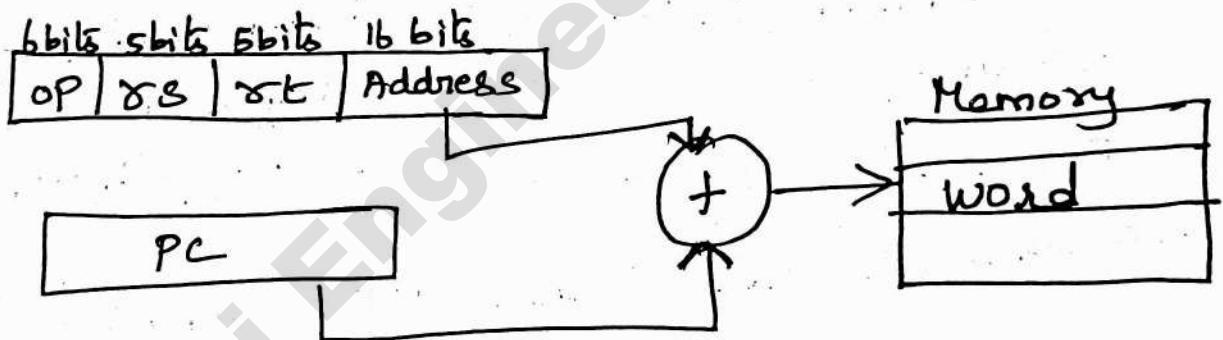
Execution time after improvement

$$\text{Speed UP} = \frac{\text{Execution time for Entire task Using Improved Machine}}{\text{Execution time for Entire task Using old original Machine}}$$

3. Brief about Relative addressing mode with an Example. [NOV/DEC-2014]

* The pc-relative addressing mode can be used to load a register with a value stored in program memory a short distance away from the current instruction.

* It can be seen as a special case of the "base + offset" addressing mode, one that selects the Program Counter (PC) as the "base register".



PC - relative addressing

Example:

`lwr $s0, $s1, Exit`

(ie) #goto Exit if $\$s0 \neq \$s1$

14

Classify the instructions based on the operation they perform and give one example to each category. [Apr/May-18] [Nov/Dec-18]

Depending on operation they perform, all instructions are divided into several groups

i) Arithmetic Instructions

eg: ADD r_1, r_2, r_3 (ie) $r_1 = r_2 + r_3$

ii) Data transfer Instructions

eg: MOV r_1, r_2 (ie) $r_1 = r_2$

iii) Logical Instructions

eg: AND r_1, r_2, r_3 (ie) $r_1 = r_2 \& r_3$

iv) Conditional Branch Instructions

eg1: CMP r_1, r_2 (ie) cond. flag = $r_1 - r_2$

v) eg2: beq $\$s1, \$s2, 25$

(ie) if ($\$s1 = \$s2$) go to $PC + \$ + 100$

vi) Unconditional Branch Instructions

eg: jr. $\$ra$ (ie) go to $\$ra$

15 Write the equation for the dynamic power required per transistor. [Apr/May-18]

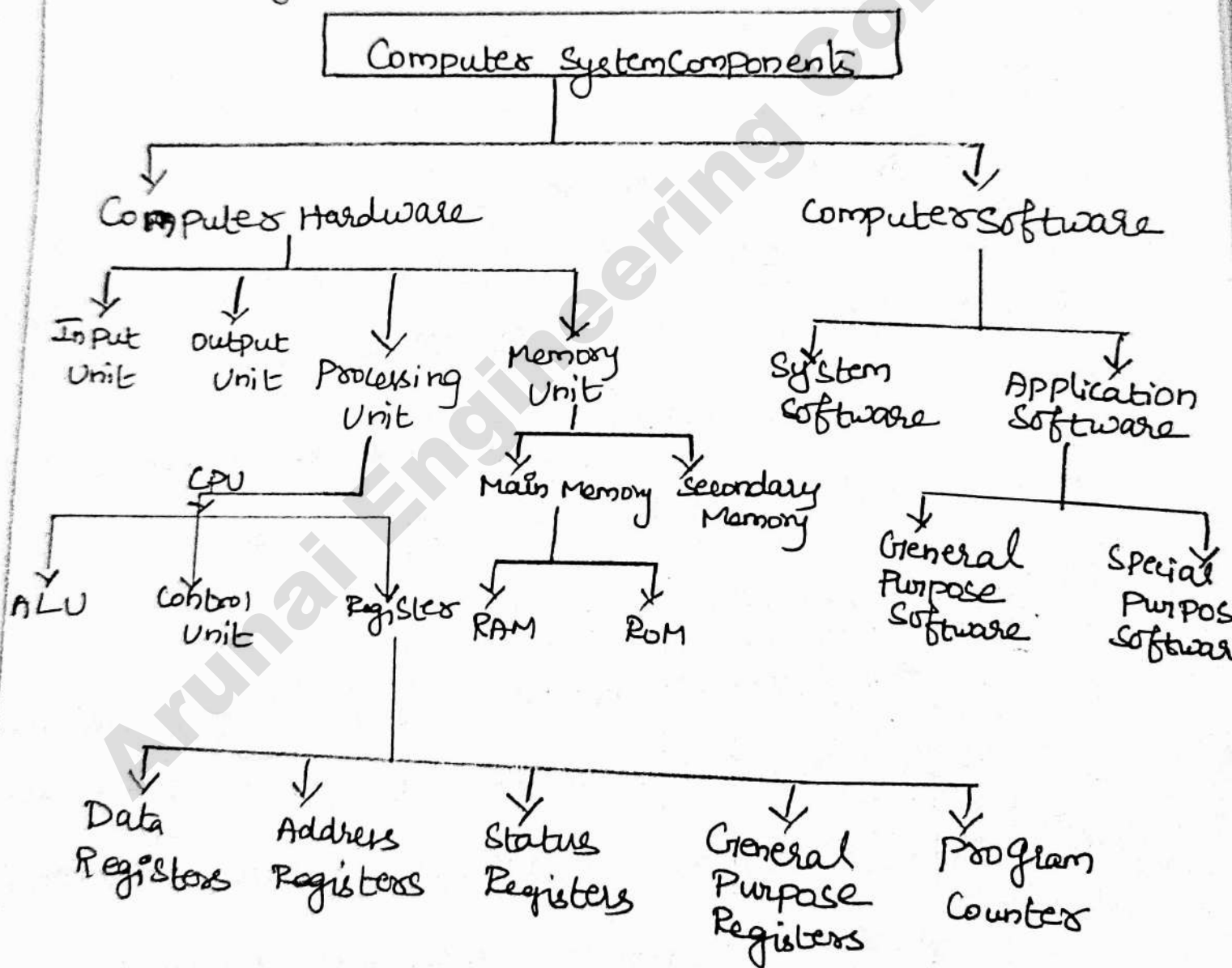
The power required per transistor is just the product of energy of a transition and the frequency of transitions:

$$\text{Power} \propto \frac{1}{2} \times \text{Capacitive load} \times \text{Voltage}^2 \times \text{Frequency switched}$$

PART-B

Explain in detail the various components of computer system with neat diagram. [Nov/Dec-2016] [Nov/Dec-2014] [Nov/Dec-2015]

* Computer System contains many kinds of components. As a computer architect must know the components and functions of it used in Computer System



Computer Hardware

* The term hardware refers to the physical components of computer such as system unit, mouse, keyboard, monitors, printer, etc.

* Computer hardware can be classified into many units such as follows

1. Input unit
2. output unit
3. processing unit
4. Memory unit

1. Input Devices

* Input devices are used to give information to the computer. Different kinds of input devices can perform different kinds of tasks such as follows.

Mouse, keyboard, Tracker ball, Scanners, Touch pads, Light Pens, Joy sticks.

2. output Devices

* output devices are used to get the results of task performed by the computer. Different kind of output devices used for different purpose.

* VDU or monitors, printers, plotters, Speakers, speech synthesizers

Processing Unit

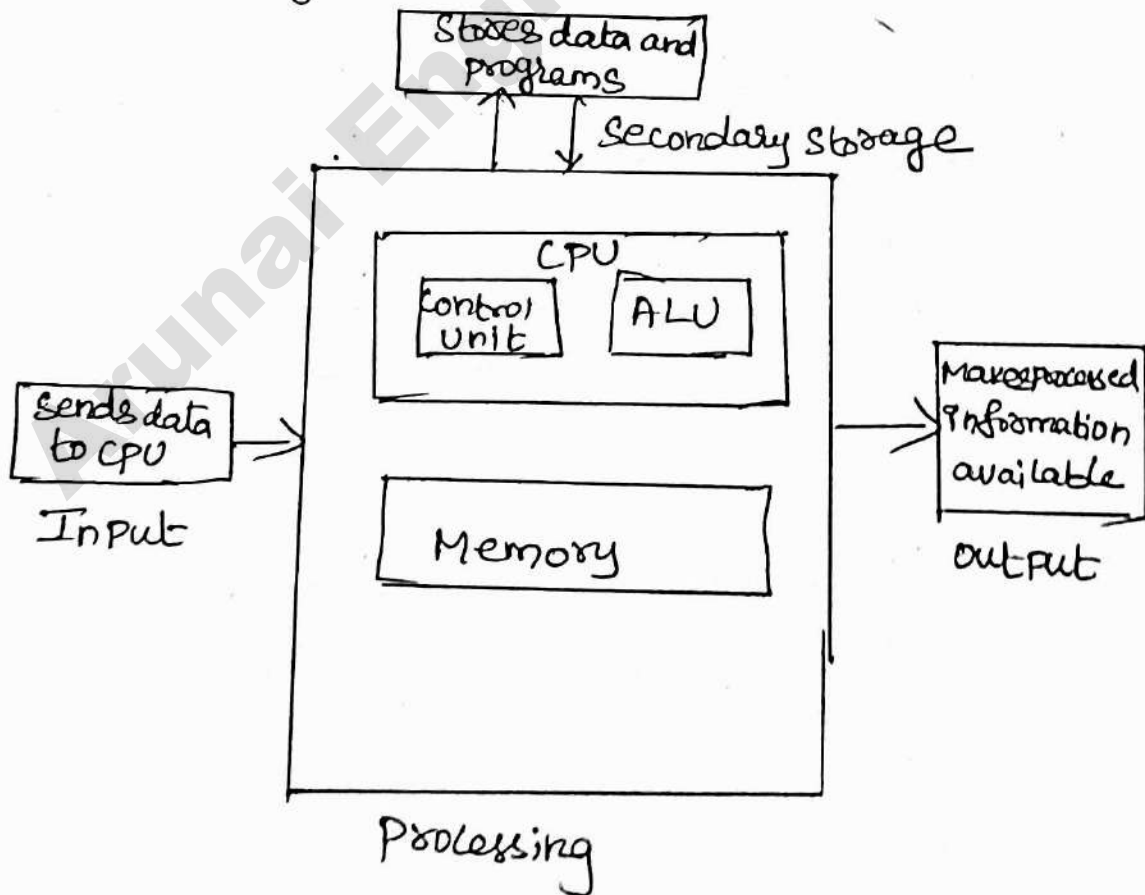
* Processing Unit is the brain of any computer system. It contains the following functional units

1. Central Processing Unit (CPU)
2. Arithmetic and Logical Unit (ALU)
3. Control Unit
4. Registers

Central Processing Unit (CPU)

* CPU is a complex set of electronic circuit which acts as a control central for the system.

* Set of electronic circuit used to execute stored program instructions.



CPU contains two major parts

1. Control Unit (CU)
2. Arithmetic and Logic Unit (ALU)

1. Control Unit (CU)

* Control Unit is a part of the hardware

* It directs the computer system to execute stored program instructions and tells the ALU what operation

* Control Unit must communicate with memory and ALU

* It sends data and instructions from secondary storage to memory as needed.

2. Arithmetic and Logical Unit (ALU)

* ALU is a digital circuit used to perform arithmetic and logic operation

* ALU performs basic arithmetic and logic operations

eg Arithmetic operations are addition, subtraction, multiplication and division.

Logic operations are comparisons of values such as NOT, AND and OR.

Registers:

* In a computer, a register is one of a small set of data holding places that are part of a computer processor.

* A register may hold a computer instruction, a storage address or any kind of data.

Data registers

* It is the register of a computer's control unit that contains the data to be stored in the computer storage.

Address registers

* It is a CPU register that either stores the memory address from which data will be fetched to the CPU or the address to which data will be sent and stored.

Status registers

* The status register is a hardware register which contains information about the state of the processor.

Program Counter

* A Program Counter is a register in a computer processor that contains the address of the instruction being executed at the current time.

General Purpose registers

* General purpose registers are not used for storing any specific type of information.

* Instead operands as well as address are stored at the time of program execution.

Memory unit

The function of the memory unit is to store programs and data. There are two classes of memory.

1. Main memory (or) primary memory
2. Secondary memory

1. Main memory

* Main memory is a fast memory that operates at electronic speeds. (Expensive)

* Programs must be stored in the memory while they are being executed.

* The memory contains a large number of semiconductor storage cells, each capable of storing one bit of information

* Main memory can be classified into two types

1. RAM - Random Access Memory (volatile (ie) Information is lost when you switch off PC)
2. ROM - Read only Memory (Non-volatile eg Network cards and video cards)

2. Secondary memory

* It is used when large amount of data and many programs have to be stored.

* It can be classifications, 1. Harddisk 2. Floppy disk
3. CD-ROM or optical disks 4. Magnetic disks 5. Tape drive 6. Flash memory

Computer Software

Software is a set of programs, which is designed to perform a well defined function. A program is a sequence of instruction written to solve a particular problem.

Types of software

1. System software
2. Application software

System software:

* System software is collection of programs designed to operate, control and extend the processing capabilities of the computer itself.

* System software is especially designed to control different operations of computer system.

Eg operating systems, utility programs, device drivers

Types of System software

1. System Control program
2. System Support program
3. System development programs

Application Software

* Application software is the software that is designed to satisfy a particular need of a particular environment.

* These software are especially designed to solve the specific problems of users and it also known as software package.

Eg student record software,

Income tax software,

Financial packages

Types of application software

1. General purpose application software
2. Special purpose application software.

What is an addressing mode? Explain the various addressing modes with suitable examples. [Nov/Dec-2015]

Explain the different types of addressing modes with suitable Example. [Nov/Dec-2016]

What is the need for addressing in a computer system? Explain different addressing modes with Examples [May-2015] [Apr/May-17]

Addressing Modes:

* Addressing mode is one of several addressing regimes delimited by their varied use of operands and/or address.

* The different ways that a processor can access data are referred to as addressing schemes or addressing modes.

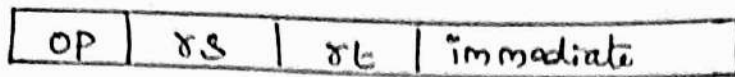
* An address computed by the processor when executing a memory access or branch instruction or when fetching the next sequential instruction is known as Effective Address (EA).

* An Effective Address can be made up of three elements, the base, index & displacement.

MIPS has the following addressing modes

- 1) Immediate addressing
- 2) Register addressing
- 3) Base or displacement addressing
- 4) PC - relative addressing
- 5) Pseudo-direct addressing

1. Immediate Addressing



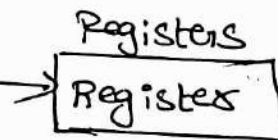
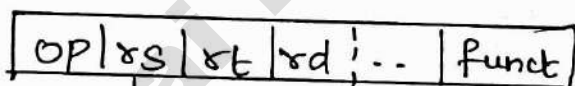
* The operand is given explicitly in the instruction.

eg MOV #20, A

This instruction copies operand 20 in the register A. The sign # in front of the value of an operand is used to indicate that this value is an immediate operand.

2. Register Addressing:

* The operand is a register. (ie) The operand is the contents of processor register. The name of register is specified in the instruction.



* Either compiler or assembler must break large constants into pieces and then reassemble them into a register.

* In this type of addressing mode the name of the register is given in the instruction where the data to be read or result is to be stored.

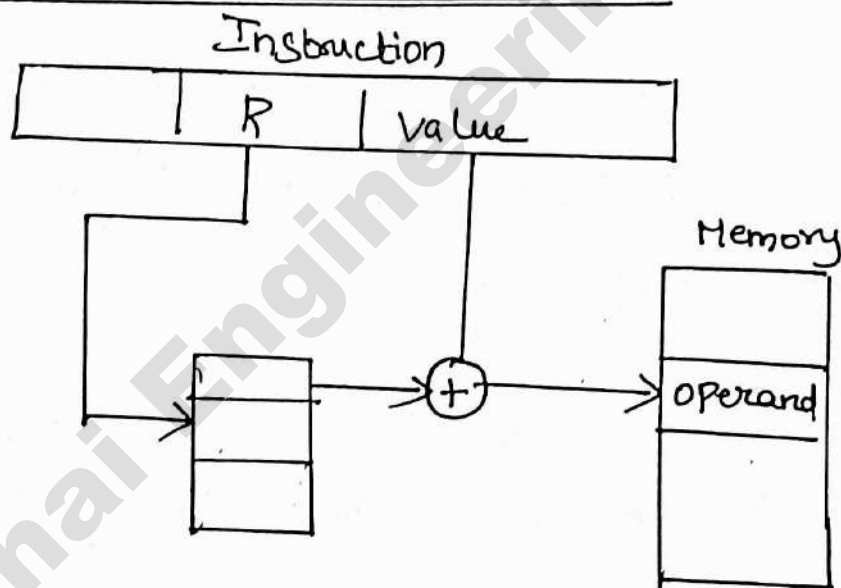
* In immediate instructions the size will be the problem for performing a load and store operations.

* To solve this problem registers are used for temporary

Eg: MOV R₂, R₁:

This instruction copies the contents of register R₂ to register R₁.

3. Base or Displacement Addressing:



* This addressing mode combines the capabilities of direct addressing and register indirect addressing.

* In this addressing mode, instruction has two address fields:

- > value
- > Referenced register

* The effective address is computed by adding contents of referenced register to value.

$$EA = \text{value} + (R)$$

* Three common variation of displacement addressing are

- i) Relative Addressing
- ii) Base register addressing
- iii) Index addressing

4.
i) Relative Addressing mode:

* The referenced register is program counter (PC) and hence this addressing mode is also known as PC-relative addressing.

* The Effective Address is determined by adding the contents of PC to the address field

$$EA = PC + \text{Address Part of Instruction}$$

* The address part is a signed number so that it is possible to have branch target location either before or after the branch instruction. This addressing mode is commonly used to specify the target address in branch instruction.

Eg JNZ BACK

This instruction causes program execution to go to the branch target location identified by the name BACK, if the branch condition is satisfied.

i) Base register addressing:

* In this addressing mode, the referenced register contains the main memory address and address field contains the displacement.

* Displacement is usually unsigned integer number

$$EA = (R) + \text{Displacement}$$

Eg: MOV[R+8], A

* This instruction copies the contents of memory whose address is determined by adding the contents of register R and displacement 8 to the register A.

iii) Index Addressing mode:

* In this addressing mode, the address field references the main memory and the referenced register contains a positive displacement from the address

$$EA = \text{Memory address} + (R)$$

$$EA = X + (R)$$

* Used for Array type data structures

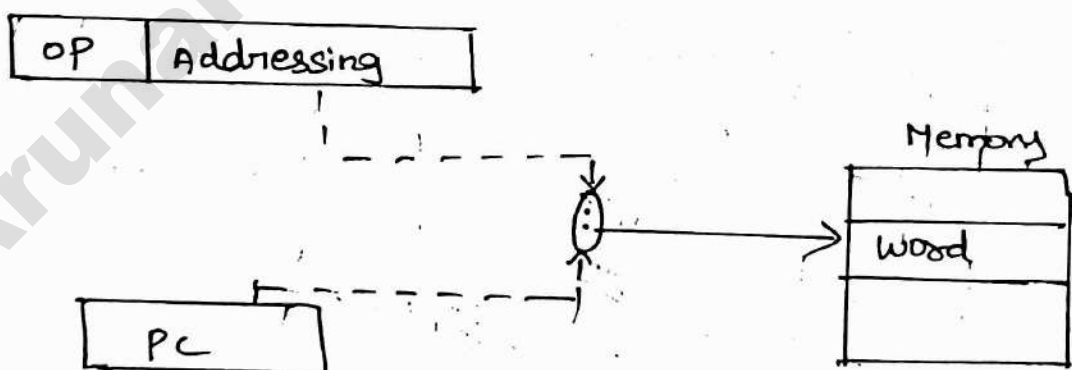
* The indexing is a technique that allows programmer to point or refer the data (operand) stored in sequential memory locations one by one. It is an efficient mechanism for performing iterative operations.

Eg $\text{MOV}[R_1 + R_I], R$

* In this instruction main memory address is given by register R_1 and the referenced register R_I gives the positive displacement. The contents of the memory address generated by the addition of main memory address and displacement is copied to register R .

5. PseudoDirect Addressing:

* The jump address is the 26 bits of the instruction concatenated with the upper bits of the PC



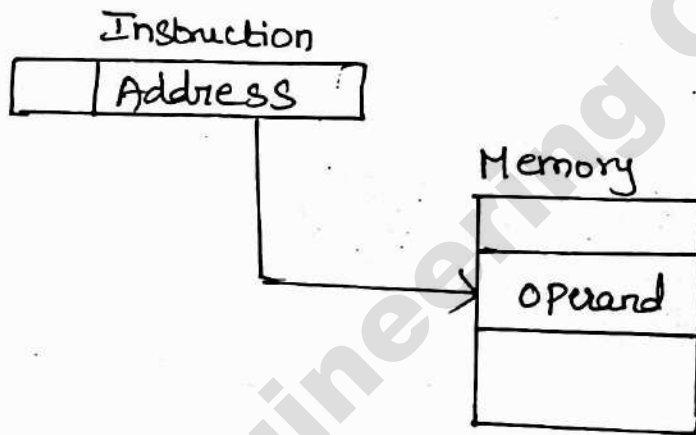
Other Addressing Modes:

6. Absolute or Direct addressing mode:

The address of the location of the operand is given explicitly as a part of the instruction.

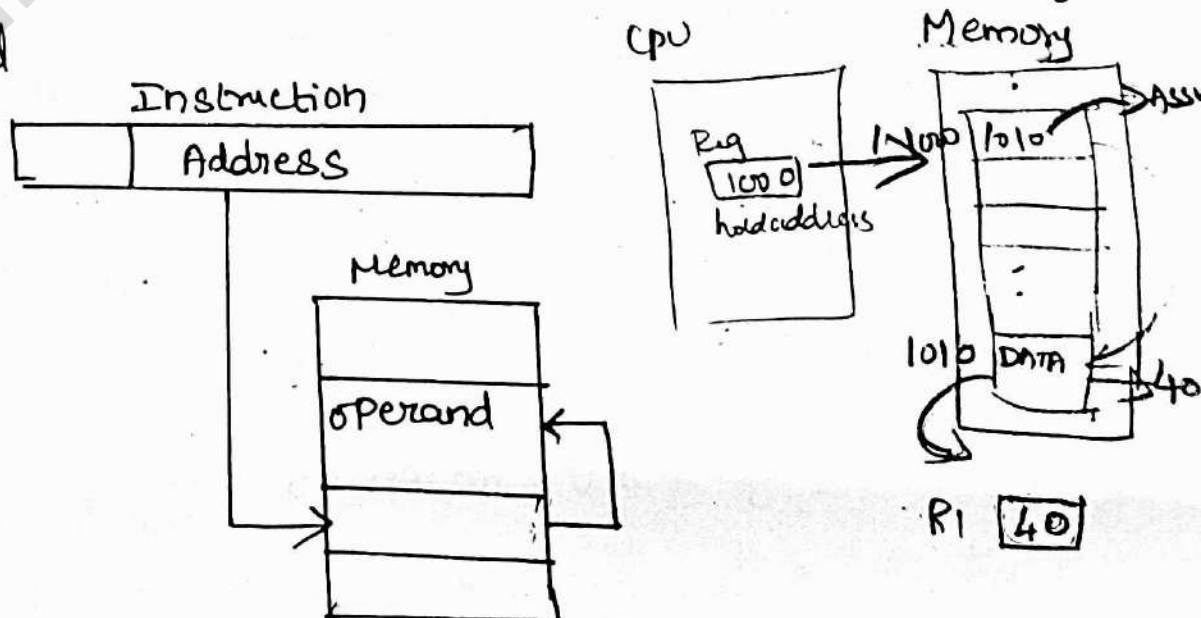
Eg MOV 2000, A

This instruction copies the contents of memory location 2000 into the A register.



7. Indirect Addressing mode:

* In this addressing mode, the instruction contains the address of memory which refers the address of the operand

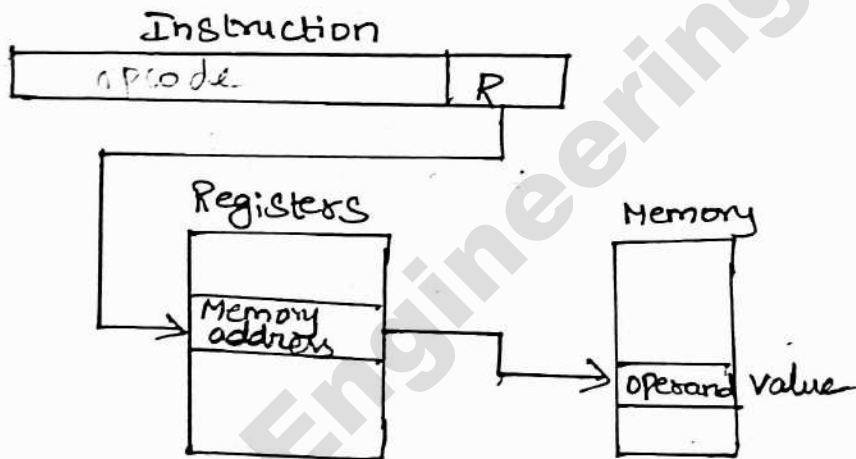


8. Register Indirect addressing mode:

* The effective address of the operand is the contents of a register or the main memory location whose address is given explicitly in the instruction.

eg MOV (R0), A

* This instruction copies the contents of memory addressed by the contents of register R0 into the register A.



9. Autoincrement addressing mode

* The effective address of the operand is the contents of a register specified in the instruction.

* After accessing the operand, the contents of this register are incremented to address the next location.

eg MOV R0, (R2)+

* The instruction copies the contents of register R0 into the memory location whose address is specified

by the contents of register R₂. After copy operation, the contents of register R₂ are automatically incremented by 1.

10. Autodecrement addressing mode:

* The contents of a register specified in the instruction are decremented and then they are used as an effective address to access a memory location.

eg MOV - (R₀), R₁

* This instruction, initially decrements the contents of register R₀ and then the decremented contents of register R₀ are used to address the memory location.

* Finally the contents from the addressed memory locations are copied into the register R₁.

11. Stack addressing mode:

* A stack is linear array of reserved memory locations. It is associated with a pointer called stack pointer (SP).

* In stack addressing mode, stack pointer always contains the address of TOP of stack (TOS) where the operand is to be stored or located.

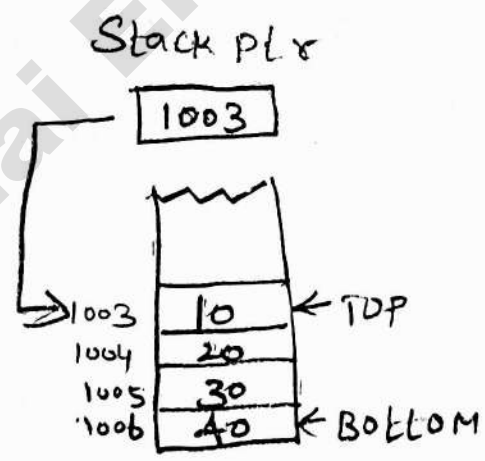
* Thus, the address of the operand (source or destination) is the contents of stack pointer.

* This addressing mode is the special case of register indirect addressing where referenced register is a stack pointer.

* Usually, stack grows in the direction of descending addresses (descending stack), starting from a high address and progressing to lower one.

* SP (stack pointer) is decremented before any items are appended (pushed) on stack

* SP is incremented after any items are popped from the stack.



State the CPU performance equation and

discuss the factors that affect performance

Explain the important measures of the performance of a computer and

Performance: derive the basic performance equation [Nov/Dec-14] [Apr/May-17]

* Performance of computer depends on many factors such as modern software systems and wide range of performance improvement in hardware side.

* Here it can also be evaluated by its speed

Response time or Execution time

* Total time required for the computer to complete the task including disk access, memory access, I/O activities, CPU Execution.

Throughput (or) Bandwidth:

* Total amount of work done in a given time

↳ To understand the relationship between throughput and Response time.

↳ To maximize performance:

* Minimise Response time or Execution time for some task

Relate Performance and Execution time for a Computer

$$\text{Performance}_x = \frac{1}{\text{Execution time}_x}$$

Then for two computers x & y

↳ performance x is greater than the performance of y

$$\text{Performance}_x > \text{Performance}_y$$

(ie)

$$\frac{1}{\text{Execution time}_x} > \frac{1}{\text{Execution time}_y}$$

(ie)

$$\text{Execution time}_y > \text{Execution time}_x$$

(ie) Execution time on y is larger than that on x

* If x is n times faster than y

$$\frac{\text{Performance}_x}{\text{Performance}_y} = n$$

Then Execution time will be

$$\frac{\text{Execution time}_y}{\text{Execution time}_x} = n$$

(ie) y is n times longer than on x

Eg: Computer A runs a program in 10s, Computer B runs same program in 15s, how much faster A is than B
We know that A is n times faster than B

$$\frac{\text{Performance A}}{\text{Performance B}} = \frac{\text{Execution time B}}{\text{Execution time A}} = n$$

Thus the performance ratio is $\frac{15}{10} = 1.5$

Then A is 1.5 times faster than B

Measuring performance:

* Time is the measure of computer's performance

It can be defined as

→ wall clock time

→ Response time

→ Elapsed time

CPU Execution time (or) CPU time:

* Actual time that the CPU spends computing for a specific task.

* It has two types

1) User CPU time:

* CPU time spent in the programs

System CPU time:

↳ The CPU time spent in the operating system performing task on behalf of the programs.

↳ Computer designers construct the computers with a clock which determines when the events to be take place

Clock cycles:

* It is also called tick, clock cycles

* It is the time for one clock period, usually of the processor clock, which runs at constant rate

Clock period:

* Length of each clock cycle called clock period

CPU performance and its factors:

* A simple formula relates the basic metrics to CPU time,

$$\text{CPU Execution time for a program} = \text{CPU Clock cycles for a program} \times \text{Clock cycle time}$$

Alternatively, clock rate and clock cycle time are inverse then,

$$\text{CPU Execution time for a program} = \frac{\text{CPU Clock cycles for a program}}{\text{Clock rate}}$$

↳ performance can be improved by

* Reducing the length of clock cycle or Number of clock cycles required for a program

Instruction Performance:

* Execution time is equals to the: number of instructions executed multiplied by average time per instructions

* The number of clock cycles required for a program can be written as

$$\text{CPU Clock cycles for a program} = \text{Instructions for a program} \times \text{Average Clock cycle Per Instruction}$$

Clock cycle Per Instruction (CPI)

* CPI is the average number of clock cycle per instruction for a program or program fragment

* CPI is the average of all instructions executed in the program.

Classic CPU Performance Equation:

* Instruction Count:

The number of instructions executed by the program

$$\text{CPU time} = \text{Instruction Count} \times \text{CPI} \times \text{Clock cycle time}$$

* Clock rate is inverse of clock cycle time

$$\text{CPU time} = \frac{\text{Instruction Count}}{\text{Clock rate}} \times \text{CPI}$$

Discuss about various techniques to represent instructions in a computer system [APR/MAY-2015]
Explain various instruction formats and illustrate the same with an example [Nov/Dec-17]
Representing Instructions in a Computer System:

* An instruction is an order given to a computer processor by a computer program.

* At the lowest level, each instruction is a sequence of 0's and 1's that describe a physical operation the computer is to perform and depending on the particular type the operation is varied.

* The specification of special storage areas called registers that may contain data to be used in carrying out the instruction or the location in computer memory of data.

* Instructions are kept in the computer as a series of high and low electronic signals and it may be represented as numbers.

* To represent the instructions we need to use the Instruction format specified by particular language.

Instruction Format

It is a form of representation of an instruction composed of fields of binary numbers

Registers

* In computer hardware registers are referred to in instructions.

* But instructions are represented as numbers so we must convert register names into numbers

* In MIPS assembly language, registers \$s0 to \$s7 map onto registers 16 to 23 and registers \$t0 to \$t7 map onto registers 8 to 15

* \$s0 means 16, \$s1 means 17, \$s2 means 18 and so on as like this \$t0 means register 8, \$t1 means register 9 and so on.

Machine Language

* Assembly language instructions use exactly 32 bits and the same size as a data word.

* All MIPS instructions are 32 bits long so we need to focus some numeric versions of instruction called machine language

* Machine Language is a binary representation used for communication within a computer system

* Instruction used in machine language are called machine code.

Eg:

Translating a MIPS assembly instruction into a machine instruction, let us consider MIPS instruction

add \$t0, \$s1, \$s2

Translate first as a combination of decimal numbers and then a binary numbers.

Answer:

The decimal representation is

0	17	18	8	0	32
---	----	----	---	---	----

* Each of these segments of an instruction is called a field

* The first and last fields (0 & 32) in combination tell the MIPS computer that this instruction perform addition

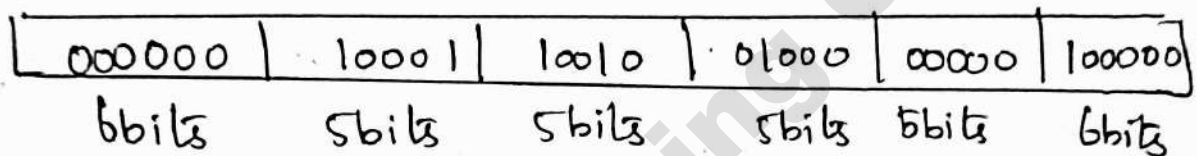
* Second field gives the number of the register that is the first source operand of the addition operation (17 \$s1)

* Third field gives the other source operand for the addition ($18 \$s_2$)

* Fourth field contains the number of register that is to receive the sum ($8 \$to$)

* Fifth field is unused in this instruction so it is set to 0

* This instruction adds register $\$s_1$ to register $\$s_2$ and places the sum in register $\$to$



Hexadecimal numbers:

* Computer can use binary numbers to read and write data. In binary number format for small value also it requires large amount of bits

* So we can use higher base that can be easily converted into binary.

* All computer data sizes are multiples of 4 in that hexadecimal numbers are popular

* Base value of hexadecimal number is 16 and it is power of 2

Hexadecimal to binary conversion

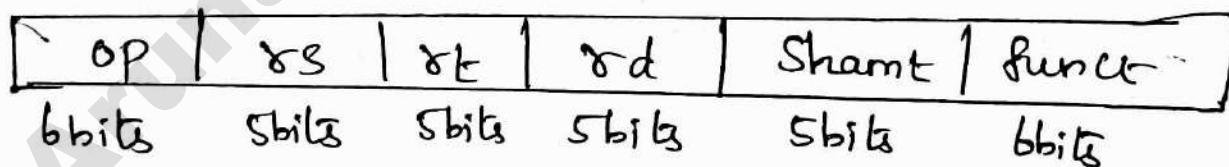
Hexadecimal	Binary	Hexadecimal	Binary
0 hex	0000 two	8 hex	1000 two
1 hex	0001 two	9 hex	1001 two
2 hex	0010 two	A hex	1010 two
3 hex	0011 two	B hex	1011 two
4 hex	0100 two	C hex	1100 two
5 hex	0101 two	D hex	1101 two
6 hex	0110 two	E hex	1110 two
7 hex	0111 two	F hex	1111 two

MIPS Field:

MIPS fields has two kind of format such as

1. R-type or R-format (for register)
2. I-type or I-format (for immediate)

1. R-format



OP - basic operation of the instruction called as opcode

↳ opcode denotes the operation and format of an instruction

rs - the first register source operand

rt - the second register, source operand

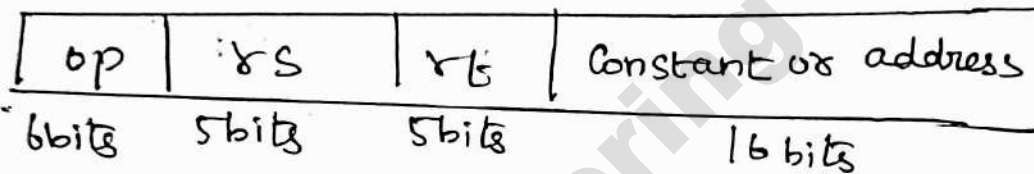
rd - the register destination operand, it gets result

shamt - Shift amount

func - This field called function or function code,

Selects the specific variant of the operation in the op field

Q. I-Format



* Used by the immediate and data transfer

Instructions

Drawbacks of different formats

* Multiple formats complicate the hardware

* It increase the complexity

5. Assume a two address format specified as source, destination. Examine the following sequence of instructions and explain the addressing modes used and the operation done in every instruction.

1) Move: (R5)+, R0

[Nov/Dec-14]

2) Add (R5)+, R0

3) Move R0, (R5)

4) Move lb(R5), R3

5) Add #40, R5

Solution:

1) Move (R5)+, R0

* This instruction copies the contents of memory location whose address is specified by register R5 to register R0.

* After copying operation, then increment R5 to the next address.

2) Add (R5)+, R0

* The instruction adds the contents of memory location whose address is specified by register R5 in the content of register R0 and the

result is stored into R0, then increment the R5 address value.

3) Move (R0), (R5)

Move the value from the address R0 to the address pointed to by R5

4) Move 16(R5), R3

* In index mode, the operand specifies both the register and an offset value x .

* The value x can be either a number, a symbol or an expression. Adding x to the value in the register gives an effective address.

* This mode can be used randomly to access elements in an array, where the register contains a base address that marks the start of the array and x is an offset from this base address.

So, "Move the value from the address pointed to by R5 at 16 bits offset of R5 register to R3."

5) Add #40, R5

* It is immediate addressing. This instruction will store value 40 in register R5.

5.11)

Consider the computer with three instruction classes and CPI measurements are given below and instruction counts for each instruction class for the same program from two different compilers are given. Assume that the computer's clock rate is 4GHz. Which code sequence will execute faster according to execution time?

[Nov/Dec-14]

Code from	CPI for this instruction class		
	A	B	C
CPI	1	2	3

Code from	Instruction Count for each class		
	A	B	C
Compiler 1	2	1	2
Compiler 2	4	1	1

Solution:

i) Instructions:

* Instruction executed by compiler 1 = $2 + 1 + 2$
 $= 5$ Instructions

* Instruction executed by compiler 2 = $4 + 1 + 1$
 $= 6$ Instructions

(ie) compiler 2 executes more number of instructions.

ii) CPU Clock cycles:

The total number of clock cycles for each sequence can be found using the following equation.

$$\text{CPU clock cycles} = \sum_{i=1}^n (\text{CPI}_i \times C_i)$$

$$\text{CPU clock cycles} = \text{Instructions executed} \times \text{corresponding CPI}$$

$$\begin{aligned} \text{CPU clock cycles for compiler 1} &= (2 \times 1) + (1 \times 2) + (2 \times 3) \\ &= 2 + 2 + 6 \\ &= 10 \text{ cycles} \end{aligned}$$

$$\begin{aligned} \text{CPU clock cycles for compiler 2} &= (4 \times 1) + (1 \times 2) + (1 \times 3) \\ &= 4 + 2 + 3 \\ &= 9 \text{ cycles} \end{aligned}$$

* It is clear that compiler 2 is faster than compiler 1 even though it executes one extra instruction.

iii) Time for executing the instruction:

$$\text{CPU execution time} = \frac{\text{CPU clock cycles}}{\text{Clock rate}}$$

Given:

$$\text{Clock rate} = 4 \text{ GHz} = 4 \times 10^9 \text{ Hz}$$

$$\left. \begin{array}{l} \text{CPU execution time for} \\ \text{Compiler 1} \end{array} \right\} = \frac{\text{CPU clock cycles for compiler 1}}{\text{Clock rate}}$$

$$= \frac{10}{4 \times 10^9}$$

$$= 2.5 \times 10^{-9} \text{ second}$$

$$\left. \begin{array}{l} \text{CPU execution time for} \\ \text{Compiler 2} \end{array} \right\} = \frac{\text{CPU clock cycles for compiler 2}}{\text{Clock rate}}$$

$$= \frac{9}{4 \times 10^9}$$

$$= 2.25 \times 10^{-9} \text{ second}$$

(re) Execution time for compiler 1 is greater than the execution time for compiler 2

So, compiler 2 will execute faster than the compiler 1

6. Explain with an example about the operations and operands of the Computer Hardware? [Nov/Dec-17]

Operations of the Computer Hardware:

* Every computer will perform different kinds of operations based on the application areas.

* Every computer must be able to perform arithmetic operations.

* The MIPS assembly language notation for performing addition operation is

```
add a, b, c
```

↳ add = addition operation
a, b, c are variables

* It instructs a computer to add the two variables b and c and put their sum in a.

Example: Suppose we need to place the sum of four variables b, c, d and e into variable a means the following sequence of instructions can be executed.

add a, b, c : the sum of b and c is placed in a
add a, a, d : the sum of b, c, and d is now in a
add a, a, e : the sum of b, c, d and e is now in a

Three instructions to sum the four variables

Category	Instruction	Example	Meaning	Comment
Arithmetic	add	add \$S1, \$S2, \$S3	$\$S1 = \$S2 + \$S3$	Three registers operands
	subtract	sub \$S1, \$S2, \$S3	$\$S1 = \$S2 - \$S3$	Three registers operands
	add immediate	addi \$S1, \$S2, 20	$\$S1 = \$S2 + 20$	Used to add constants
Data transfer	load word	lw \$S1, 20(\$S2)	$\$S1 = \text{Memory} [\$S2 + 20]$	word from memory to register
	store word	sw \$S1, 20(\$S2)	$\text{Memory} [\$S2 + 20] = \$S1$	word from register to memory
Logical	and	and \$S1, \$S2, \$S3	$\$S1 = \$S2 \& \$S3$	Three registers, operands; bit by bit AND
	or	or \$S1, \$S2, \$S3	$\$S1 = \$S2 \$S3$	Three registers operands; bit by bit OR
conditional branch	branch on equal	beq \$S1, \$S2, 25	if ($\$S1 = \$S2$) goto PC+4+100	Equal test; PC-relative branch
	branch on not equal	bne \$S1, \$S2, 25	if ($\$S1 \neq \$S2$) goto PC+4+100	Not equal test PC relative
Unconditional jump	jump	j 2500	goto 10000	Jump to target address
	jump register	jr \$ra	goto \$ra	For switch, procedure return

* To keep the hardware simple, every instruction has to have exactly three operands no more and no less

* Hardware for a variable number of operands is more complicated than hardware for a fixed number.

* Three design principles of hardware technology has

Design principles 1: Simplicity favors regularity

Design principles 2: Smaller is faster

Design principles 3: Good design demands good compromises

Design Principle 1: Simplicity favors regularity:

Compiling two C assignment statements into MIPS

* C program contains the five variables a, b, c, d and e

$$a = b + c$$

$$d = a - e$$

* The translation from C to MIPS assembly language instructions is performed by the compiler. Show the MIPS code produced by a compiler.

Solution:

* A MIPS instruction operates on two source operands and places the result in one destination operand.

* So the two simple statements above compile directly into these two MIPS assembly language instructions.

add a, b, c

sub d, a, e

Operands of the Computer Hardware:

* operand is a variable used to perform any kind of operations.

* Unlike high level languages, the operands of arithmetic instructions are restricted

* operand must be from a limited number of special locations built directly in hardware called registers.

* Registers are primitives used in hardware design that also visible to the programmer when the computer is completed.

* MIPS architecture has 32 bits registers and group of 32 bits are called word.

Design principle 2: Smaller is faster

* Registers is a very small amount of very fast memory that is built into the CPU.

* A very large number of registers may increase the clock cycle time because it takes electronic signals longer to travel it.

* Another reason for why more than 32 bit registers are not used means it takes in the instruction format.

* MIPS convention use two character names following a dollar sign to represent a register.

Example:

* Compiling a C assignment using registers. It is the compiler's job to associate program variables with registers. Take for instance, the assignment statement from our example is

$$f = (g+h) - (i+j);$$

The variables f, g, h, i and j are assigned to the registers $\$s0, \$s1, \$s2, \$s3$ and $\$s4$ respectively. What is the compiled MIPS code?

Answer:

add \$t0, \$s1, \$s2: register \$t0 contains $g+h$

add \$t1, \$s3, \$s4: register \$t1 contains $i+j$

sub \$s0, \$t0, \$t1: f gets $\$t0 - \$t1$, which is
 $(g+h) - (i+j)$

Memory operands:

* programming language have simple variables that contain single data element, but they also have more complex data structures like arrays and structures

* complex data structures contains large amount of elements than the registers in a computer.

Data transfer instructions:

* Data transfer instruction is a command that moves data between memory and registers.

Address: → A value used to define the location of a specific data element within a memory array called address.

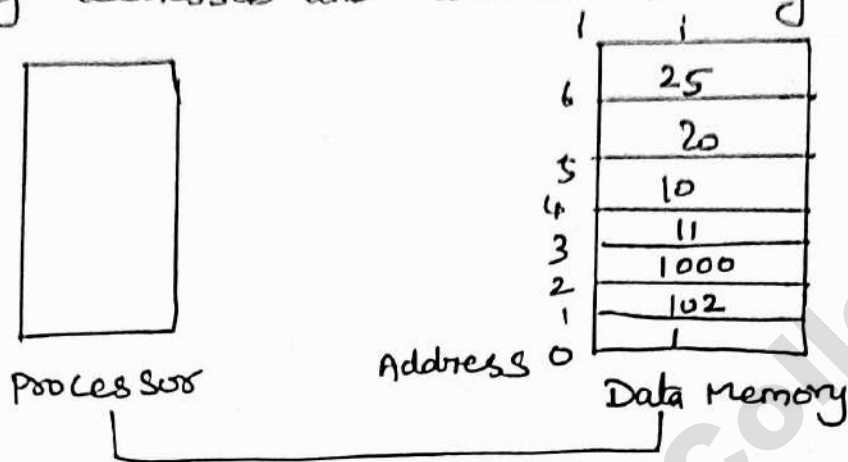
* To access a word in memory, the instruction must supply the memory address.

Memory:

* Memory is a large, single dimensional array, with the address acting as the index to that array starting at 0

Example :

Memory addresses and contents of memory at those location



* The address of the third data element is 2 and the value of memory [2] is 1000.

* The data transfer instruction copies data from memory to a register that process is called load.

Alignment Restriction:

* In MIPS, words must start at addresses that are multiples of 4. This requirement is called an alignment restriction.

Big Endian and Little Endian:

* Address of the left most byte is called "big end" and right most byte is called "little end".

Store:

* The instruction complementary to load is called store.

* Store copies data from a register to memory.

* Load copies data from memory to registers

Spilling registers:

* Many programs have more variables than the computer registers.

* Compiler tries to keep the most frequently used variables in registers and places the rest of variables in memory.

* Using loads and stores we can move variables between registers and memory.

* "The process of putting less commonly used variables into memory is called spilling registers."

Constant or Immediate Operands:

* Constant variables are used as one of the operand for many arithmetic operations in MIPS architecture.

* Constants have been placed in memory when the program was loaded.

* To avoid load instruction used in arithmetic instruction we can use one operand is a constant.

* This quick add instruction with one constant operand is called add arithmetic or addi

Eg: `addi $S3, $S3, 4; $S3 = $S3 + 4`

* Sometimes a program may need to use a constant in an operation. Those operands are called immediate operands.

7
Consider three processors P₁, P₂ and P₃ executing the same instruction set. They have clock rates of 3 GHz, 2.5 GHz and 4.0 GHz respectively and CPI of 1.5, 1.0 and 2.2 respectively. i) Which processor has the highest performance expressed in instructions per second?
[Apr/May-18] [Nov/Dec-18]

Ans:-

Performance of P₁ (Instructions/sec) = 3 GHz

$$(ie) 3 \times 10^9$$

Performance of P₂ (Instructions/sec) = 2.5 GHz

$$(ie) 2.5 \times 10^9$$

Performance of P₃ (Instructions/sec) = 2.2 GHz

$$(ie) 2.2 \times 10^9$$

CPU execution time can be calculated using the following formula.

$$CPU \text{ time} = \frac{\text{Instruction count} \times CPI}{\text{Clock rate}}$$

Let the instruction count be 'n' and assume it is same for all the three processors. Now calculate the CPU time for each processor.

$$\text{CPU time } P_1 = \frac{n \times 1.5}{3 \times 10^9}$$

$$= n \times 0.5 \times 10^{-9} \text{ sec}$$

$$\text{CPU time } P_2 = \frac{n \times 1.0}{2.5 \times 10^9} = n \times 0.4 \times 10^{-9} \text{ sec}$$

$$\text{CPU time } P_3 = \frac{n \times 2.2}{4 \times 10^9} = n \times 0.55 \times 10^{-9} \text{ sec}$$

CPU time for processor P_2 is smallest than other processors. Hence, P_2 produces highest performance than other processors.

i) If the processors each execute a program in 10 seconds, find the number of cycles and the number of instructions in each processor.

Number of clock cycles

CPU time for P_1 , P_2 and $P_3 = 10 \text{ sec}$

The formula for finding CPU time can be rewritten to find the number of clock cycles.

$$\text{CPU time} = \frac{\text{CPU clock cycles}}{\text{Clock rate}}$$

$$\text{CPU Clock cycles} = \text{CPU time} \times \text{Clock rate}$$

$$\text{CPU Clock cycles } P_1 = 10 \text{ sec} \times 3 \times 10^9 \text{ Hz} = 30 \times 10^9 \text{ cycles}$$

$$\text{CPU Clock cycles } P_2 = 10 \text{ sec} \times 2.5 \times 10^9 \text{ Hz} = 25 \times 10^9 \text{ cycles}$$

$$\text{CPU Clock cycles } P_3 = 10 \text{ sec} \times 4 \times 10^9 \text{ Hz} = 40 \times 10^9 \text{ cycles}$$

Arunai Engineering College

8 Translate the following C code to MIPS assembly code. Use a minimum number of instructions. Assume that i and k correspond to registers $\$s3$ and $\$s5$ and the base of the array `save` is in $\$s6$. What is the MIPS assembly code corresponding to this C segment?

```
while (save[i] == k)
```

[Apr/May-18]

```
  i++;
```

Ans:-

* The first step is to load `save[i]` into a temporary register. Before we can load `save[i]` into a temporary register, we need to have its address.

* Before we can add i to the base of array `save` to form the address, we must multiply the index i by 4 due to the byte addressing problem.

* We can use shift left logical, since shifting left by 2 bits multiplies by 2^2 or 4.

* We need to add the label `loop` to it so that we can branch back to that instruction at the end of the loop.

Registers allocation:

↳ i $\$S3$

↳ k $\$S5$

↳ base of $\text{save}[j]$ $\$S6$

Loop: $\text{Sll } \$t1, \$S3, 2$ # temp reg $\$t1 = i * 4$

To get the address of $\text{save}[i]$, we need to add $\$t1$ and the base of save in $\$S6$:

$\text{add } \$t1, \$t1, \$S6$ # $\$t1 = \text{address of } \text{save}[i]$

Now we can use that address to load $\text{save}[i]$ into a temporary register:

$\text{lw } \$t0, 0(\$t1)$ # Temp reg $\$t0 = \text{save}[i]$

The next instruction performs the loop test, exiting if $\text{save}[i] < k$:

$\text{bne } \$t0, \$S5, \text{Exit}$ # go to Exit if $\text{save}[i] < k$

The next instruction adds 1 to i :

$\text{addi } \$S3, \$S3, 1$ # $i = i + 1$

The end of the loop branches back to the while test at the top of the loop.

We just add the Exit label after it

Loop # go to Loop

Exit:

9 Consider two different implementations of the same instruction set architecture. The instructions can be divided into four classes according to their CPI (Class A, B, C and D). P₁ with a clock rate of 2.5 GHz and CPIs of 1, 2, 3 and 3 respectively and P₂ with a clock rate of 3 GHz and CPIs of 2, 2, 2 and 2 respectively. Given a program with a dynamic instruction count of 1.0×10^6 instructions divided into classes as follows: 10% Class A, 20% Class B, 50% Class C and 20% Class D, which implementation is faster? What is the global CPI for each implementation? Find the clock cycles required in both cases.

Answer:

Class	A	B	C	D	
P ₁	1	2	3	3	2.5 GHz
P ₂	2	2	2	2	3 GHz
IC	10%	20%	50%	20%	

P1:

CPU execution time can be calculated using the following formula.

$$\text{CPU time} = \frac{\text{Instruction Count} \times \text{CPI}}{\text{Clock rate}}$$

Class A: 1×10^5 instruction

Class B: 2×10^5 instruction

Class C: 5×10^5 instruction

Class D: 2×10^5 instruction.

$$\text{Total time}_{P1} = \frac{(1 \times 10^5 \times 1 + 2 \times 10^5 \times 2 + 5 \times 10^5 \times 3 + 2 \times 10^5 \times 3)}{2.5 \times 10^9}$$

$$C_{PA} = \frac{(1.0 \times 10^5) \times 1}{2.5 \times 10^9} = 4 \times 10^{-5} \text{ s}$$

$$C_{PB} = \frac{(2 \times 10^5) \times 2}{2.5 \times 10^9} = 1.6 \times 10^{-4} \text{ s}$$

$$C_{PC} = \frac{(5 \times 10^5) \times 3}{2.5 \times 10^9} = 6.0 \times 10^{-4} \text{ s}$$

$$C_{PD} = \frac{(2 \times 10^5) \times 3}{2.5 \times 10^9} = 2.4 \times 10^{-4} \text{ s}$$

$$\text{Total time} = 1.04 \times 10^{-3} \text{ s}$$

Clock cycles required in both cases:

$$\text{CPU Clock cycles} = \text{Instruction count} \times \text{CPI}$$

$$\text{Clock cycles}_A = (1.0 \times 10^5)(1) + (2 \times 10^5)(2) + (5 \times 10^5)(3) + (2 \times 10^5)(3)$$

$$\text{Clock cycles}_A = 2.6 \times 10^6 \text{ cycles}$$

$$\text{Clock cycles}_B = (1.0 \times 10^5)(2) + (2 \times 10^5)(2) + (5.0 \times 10^5)(2) + (2 \times 10^5)(2)$$

$$\text{Clock cycles}_B = 2.0 \times 10^6 \text{ cycles}$$

Arunai Engineering College

$$\text{Total time}_{P_2} = \frac{(1 \times 10^5 \times 2 + 2 \times 10^5 \times 2 + 5 \times 10^5 \times 2 + 2 \times 10^5 \times 2)}{(3 \times 10^9)}$$

(10)

$$CPU_A = \frac{(1.0 \times 10^5) \times 2}{3 \times 10^9} = 6.67 \times 10^{-5} \text{ s}$$

$$CPU_B = \frac{(2 \times 10^5) \times 2}{3 \times 10^9} = 1.33 \times 10^{-4} \text{ s}$$

$$CPU_C = \frac{(5 \times 10^5) (2)}{3 \times 10^9} = 3.33 \times 10^{-4} \text{ s}$$

$$CPU_D = \frac{(2 \times 10^5) (2)}{3 \times 10^9} = 1.33 \times 10^{-4} \text{ s}$$

$$\text{Total time} = 6.657 \times 10^{-4} \text{ s}$$

P_2 is the best implementation

The Global CPI for Each Implementation

$$CPI_A = \frac{(2.5 \times 10^9) (1.04 \times 10^{-3})}{10^6} = 2.6$$

$$CPI_B = \frac{(3.0 \times 10^9) (6.657 \times 10^{-4})}{10^6} = 1.9971$$

PART - A

Define Little Endian arrangement [NOV/DEC-2014]
when the lower byte addresses are used for
the less significant bytes [the rightmost byte] of
the word, addressing is called little-endian

2. What is the overflow/underflow conditions for
addition and subtraction

Addition: If two operands having same sign it
will cause the overflow in addition operation.

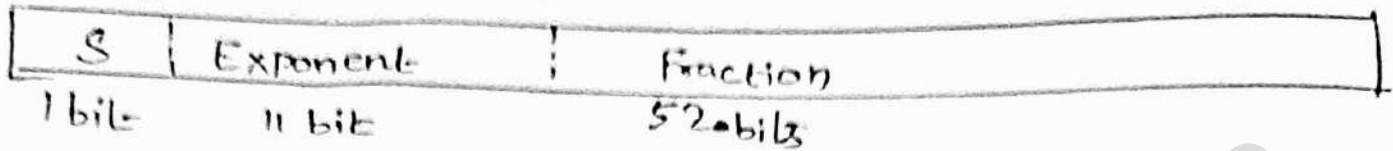
Eg first operand 32 bit word and second operand is
a 32-bit word, so the result also must be 32-bit word.
If the result exceeds more than 32 bit word it causes
overflow

Subtraction: overflow occurs in subtraction based on
two cases. [May-2015]

Case 1: If we subtract a negative number from a
positive number and the result is negative, the overflow
will occur.

Case 2: If we subtract a positive number from a
negative number and the result is positive, then overflow

3 state the representation of double precision floating point number. [Nov/Dec-2015]



* The representation of double precision floating point number take two MIPS words.

where

S is the sign of the number

4. what do you mean by subword parallelism [May-2015]

* Suppose the size of the data item within the vector is 128-bits - with this data by partitioning the carry chains, a processor can use parallelism to perform simultaneous operations on short vectors of

1. Sixteen 8 bit operands
2. Eight 16 bit operands
3. Four 32 bit operands
4. Two 64 bit operands

* The parallelism which occur within a wide word, in case it extends means that is called as subword parallelism. It is also known as data level parallelism or vector or SIMD parallelism.

5. Subtract $(11011)_2 - (10011)_2$ using 2's complement. [Nov/Dec-2017]

$$\begin{array}{r} 11011 \\ (-) 10011 \end{array}$$

(i) Take 2's complement for 10011

First take 1's complement

$$10011 = 01100$$

Now, take 2's complement 01100

$$\begin{array}{r} 01100 \\ + 00001 \\ \hline 01101 \end{array}$$

Now perform Binary Subtraction

$$\begin{array}{r} 11011 \\ 01101 \\ \hline 10100 \end{array}$$

Indicate Negative Sign

Answer = $\boxed{10100}$

Q Divide $(1001010)_2 \div (1000)_2$ [NOV/DEC-2017]

$$\begin{array}{r} 1001 \leftarrow \text{Quotient} \\ \overline{1000 \overline{) 1001010} \leftarrow \text{Dividend}} \\ \underline{1000} \\ 1010 \\ \underline{1000} \\ 10 \leftarrow \text{Remainder} \end{array}$$

1000
↓
Divisor

Arunai Engineering College

7) Subtract $(11010)_2 - (10000)_2$ using 1's complement and 2's complement method. [APR/MAY-2017]

1's Complement method

i) Take 1's complement of subtrahend: 01111

ii) Add with minuend:

$$\begin{array}{r} 11010 \\ + 01111 \\ \hline \text{Carry} \rightarrow 1 \quad 01001 \end{array}$$

iii) If there is a carry add with result.

$$\begin{array}{r} \text{(ie)} \quad 01001 \\ \quad \quad \quad 1 \\ \hline 01010 \end{array}$$

iv) The Result is 01010

2's Complement method

i) Take 2's complement of subtrahend: 10000

ii) Add with minuend:

$$\begin{array}{r} 11010 \\ + 10000 \\ \hline \text{Carry} \rightarrow 1 \quad 01010 \end{array}$$

iii) If there is a carry discard it: 01010

iv) The result is 01010

8) write the rules to perform addition on floating point numbers. [APR/MAY - 2017]

Step 1: Align the decimal numbers of the smallest Exponent number that matches the larger Exponent.

(e) Shift the smaller number right until Exponents match.

Step 2: Add/ Subtract the mantissas, depending on sign

Step 3: Normalize the result if necessary.

Step 4: Check for overflow.

Step 5: Rounding off to appropriate number of bits.

9. What is a guard bit and what are the ways to truncate the guard bits? [NOV/DEC-2016]

Guard bits :

The first of two extra bits kept on the right during intermediate calculations of a floating point format

There are several ways to truncation

1. To remove the guard bits and make no changes in the retained bits (Chopping)

2. Non Neumann rounding - If the bits to be removed are all 0s, they are simply dropped, with no change to the retained bits

3. Rounding Procedure - Achieves the closest approximation the number being truncated and is an unbiased technique.

10 What is arithmetic overflow. [NOV/DEC-2016]

* Arithmetic overflow, it is the condition that occurs when a calculation produces a result that is greater than a given register or storage location can store or represent.

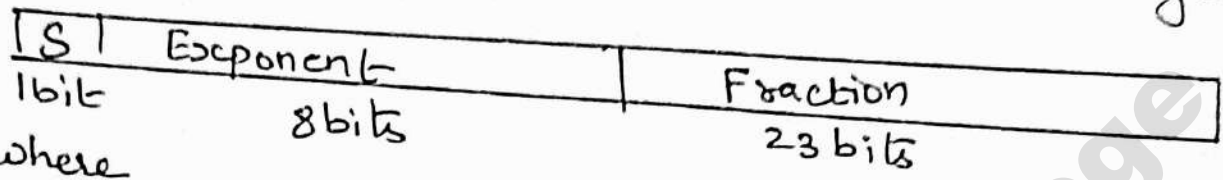
* It occurs in which positive exponent becomes too large to fit in the exponent field.

11 Define ALU [MAY-JUNE-2016]

* ALU - Arithmetic Logic Unit is a digital circuit used to perform arithmetic operations like, addition, multiplication and logic operations like AND, OR, NOT etc.

* It represents the fundamental building block of the central processing unit (CPU) of a computer.

12 Show the IEEE 754 binary representation of the number $(-0.75)_{10}$ in single precision
 [Ans: (Apr/May-18)]



where

S - sign of the floating-point number
 $0 \rightarrow$ positive, $1 \rightarrow$ negative

Exponent - The value of the 8-bit exponent field

Fraction - 23-bit number

i) Convert decimal number $(-0.75)_{10}$ into binary

$$\begin{array}{r} 0.75 \times 2 \\ \hline 1.50 \times 2 \\ \hline 1.0 \end{array}$$

The binary equivalent of

$$(-0.75)_{10} = (-0.11)_2 \quad \text{(ie) } -0.11_2 \times 2^0$$

ii) Normalize the number

$$\text{(ie) } (-0.11)_2 = -1.1 \times 2^{-1}$$

iii) Single precision representation

For the given number

$$(-1)^S \times (1 + \text{fraction}) \times 2^{(\text{Exponent} - 127)}$$

$$\text{(ie) } (-1)^0 \times (1 + 0.1_2) \times 2^{(-1 + 127)}$$

$$= (-1)^0 \times (1 + 0.1_2) \times 2^{126}$$

The given number can be written in IEEE single precision format as

Sign bit	Exponent (8bit)	Fraction (23bits)
0	0111 1110	1000000 00000000 0000 0000

13 Define a datapath in a CPU [Apr/May-18]

The 'data path' is a representation of the flow of information (data, instructions)

Data path elements are

- i) Memory unit - It is used to store the instructions of a program and supply instructions given an address
 - ii) program Counter (PC) - It is a register that holds the address of the current instruction
 - iii) Adder - An adder is a device which is used to increment the PC to the address of the next instruction
- * This adder can be built from the ALU.

14 Perform $X - Y$ using 2's complement arithmetic for the given two 16-bit binary numbers. [Nov/Dec-18]

$$X = 0000\ 1011\ 1110\ 1111 \text{ and}$$

$$Y = 1111\ 0010\ 1001\ 1101$$

$$X - Y \text{ (ie) } \begin{array}{cccc} 0000 & 1011 & 1110 & 1111 \\ & 1111 & 0010 & 1001 & 1101 \end{array}$$

Take 2's complement for 1111 0010 1001 1101

First, take 1's complement

$$\{1111\ 0010\ 1001\ 1101\} \rightarrow 0000\ 1101\ 0110\ 0010$$

Now, take 2's complement

$$\begin{array}{cccc} 0000 & 1101 & 0110 & 0010 \\ & & & + 1 \\ \hline Y = 0000 & 1101 & 0110 & 0011 \end{array}$$

Now, Perform Binary subtraction,

$$X = 0000\ 1011\ 1110\ 1111$$

$$Y = 0000\ 1101\ 0110\ 0011$$

$$\begin{array}{cccc} 0000 & 1001 & 0101 & 0010 \\ \hline \end{array}$$

Answer = 0001 1001 0101 0010

Divide $(12)_{10}$ by $(3)_{10}$ Using the Restoring and non restoring division algorithm with step by step intermediate results and Explain [Nov/Dec-2014]

Discuss in detail about division algorithm in detail with diagram and examples [Nov/Dec-2015] [Nov/Dec-2016]

Illustrate the division algorithm with an Example [Apr/May-17]
[Nov/Dec-17]

Division:

* Division is a reciprocal operation of multiplication.

* Division operation is less frequent and is faster. It will perform invalid operation when we perform dividing by 0.

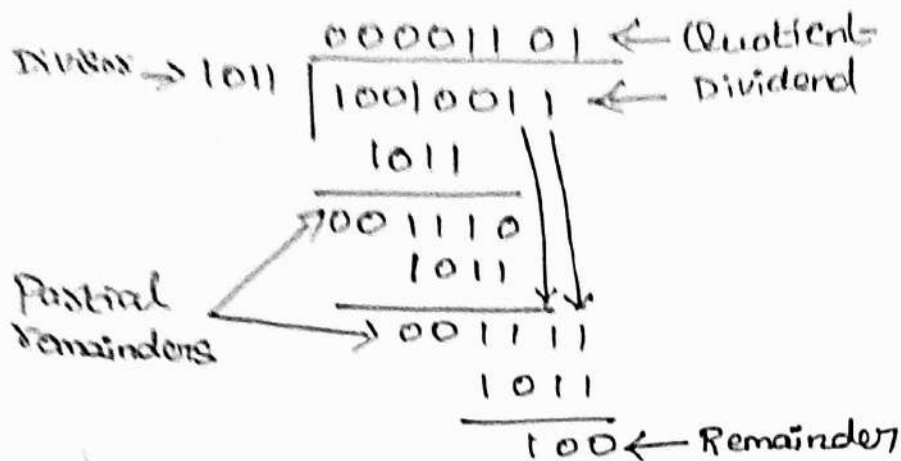
* Dividers use two operands called the dividend and divisor and produce the results called quotient and remainder.

* Hence the another way to express the relationship between the components.

$$\text{Dividend} = \text{Quotient} * \text{Divisor} + \text{Remainder}.$$

Let us assume that both the dividend and divisor are positive and hence the quotient and the remainder are non-negative.

Example



Algorithm for Restoring Division:

Step 1: Shift Remainder register left logically (S₁₁) one binary position

Step 2: Subtract divisor from left (upper half) of the remainder register and the result in the left half of the Remainder register

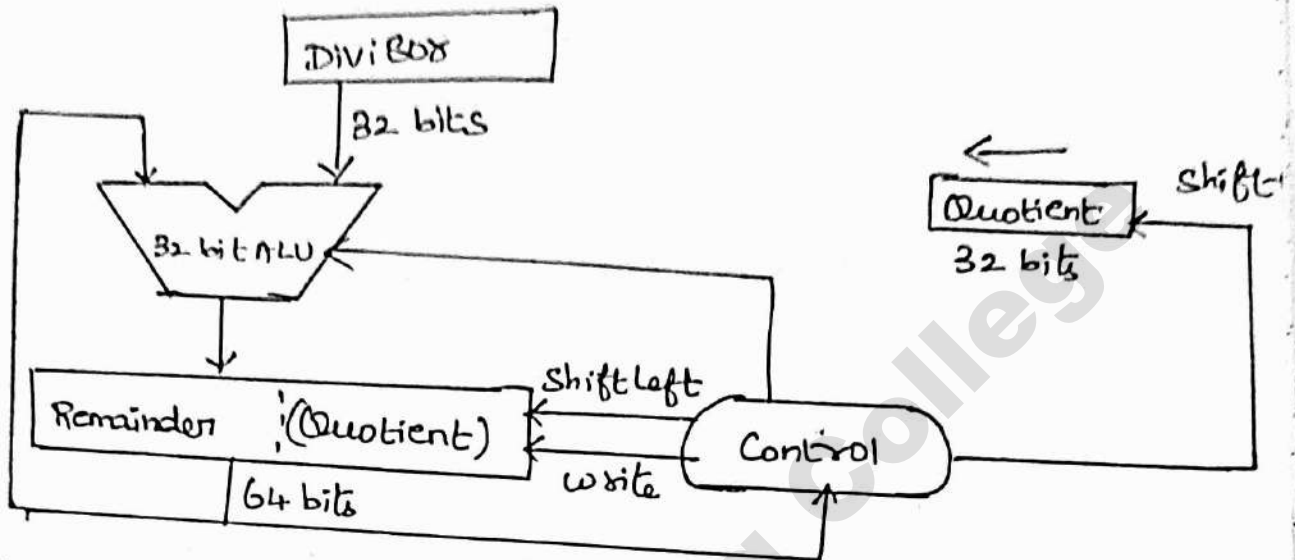
Step 3: If the sign bit of Remainder register is 1, then shift Quotient register (Q) left by one binary position and set Q₀ to 0 and add divisor back to remainder register (i.e. restore remainder register); otherwise set Q₀ to 1.

Step 4: Repeat step 1, 2 and 3 for < 33 times to obtain the remainder and quotient result.

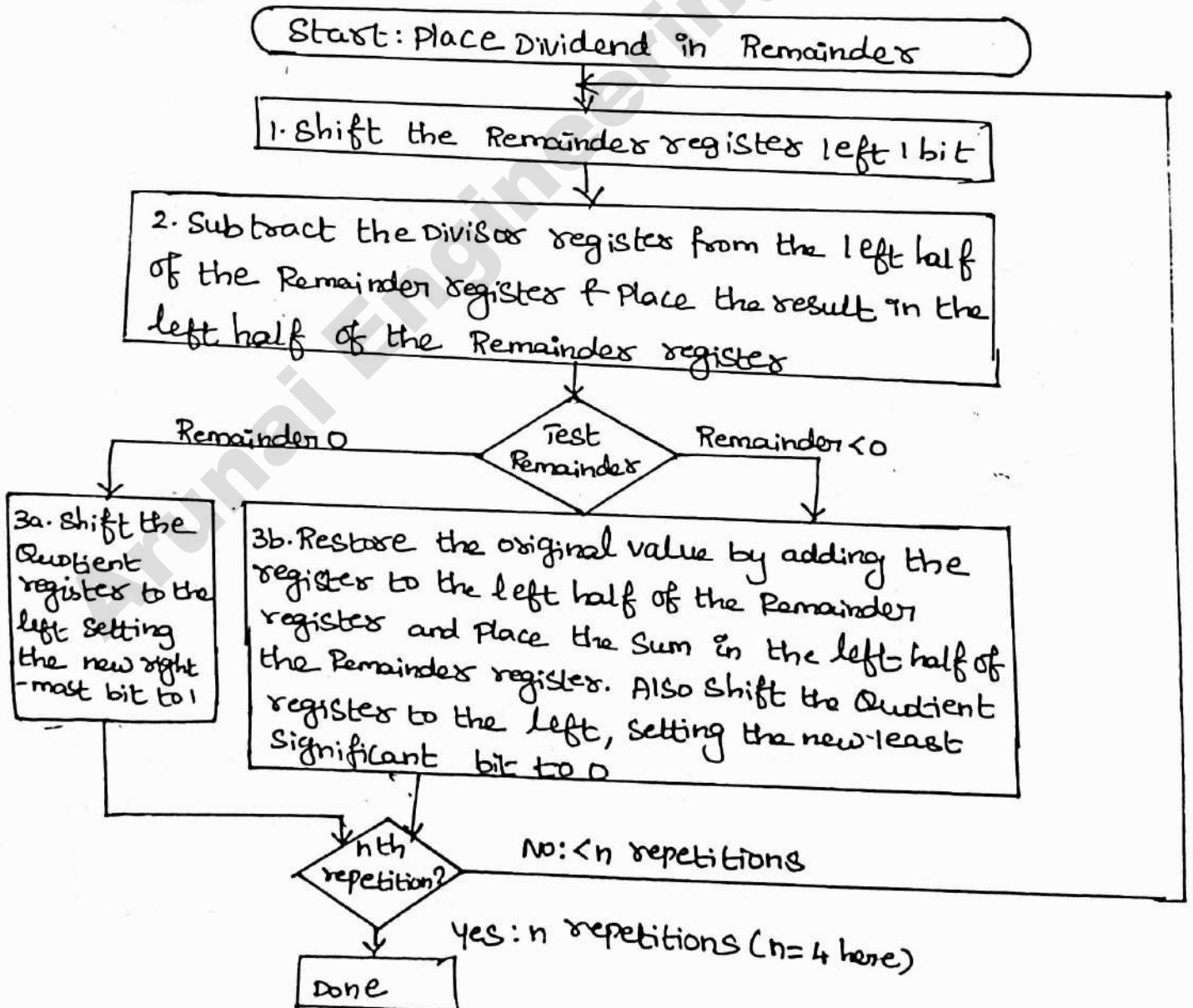
For simplicity, use 4 bit version to perform division, same concept can be expanded to 32 bit number.

Restoring Division (or) Second Version of Division Hardware of 32-bit Boolean numbers representation

Schematic Diagram of ALU circuitry



Flowchart



Algorithm for Non-Restoring Division:

Step 1: a) If the sign of Remainder Register (Rem Reg) is 0, Shift Rem Reg left logically by one binary position and subtract the divisor from the left half of the Rem Reg and place the result in left half of Rem Reg. Otherwise, shift Rem Register left and add divisor to the left of the rem reg.

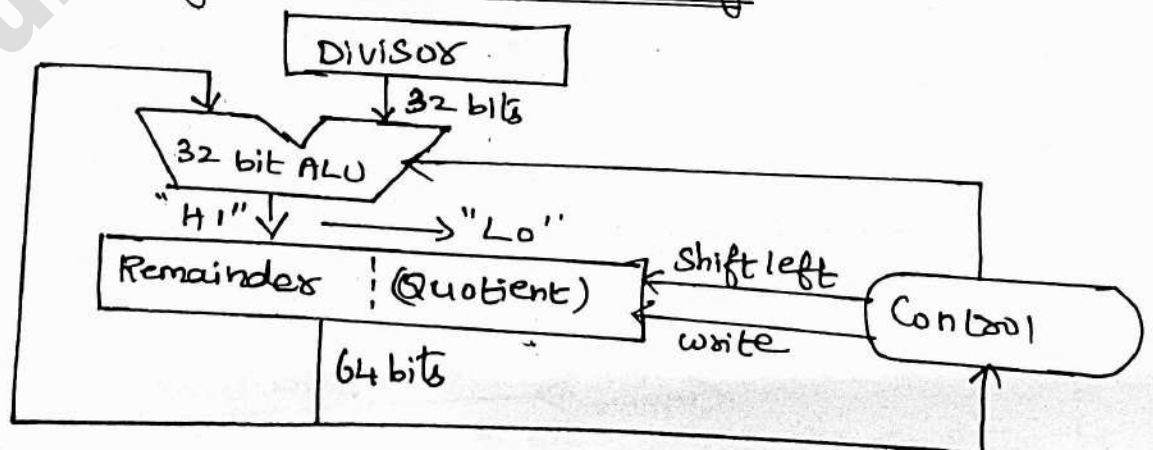
b) If the sign of Rem Reg is < 0 , set Q_0 to 0 otherwise set Q_0 to 1.

Step 2: Repeat steps 1 and 2 for < 33 times

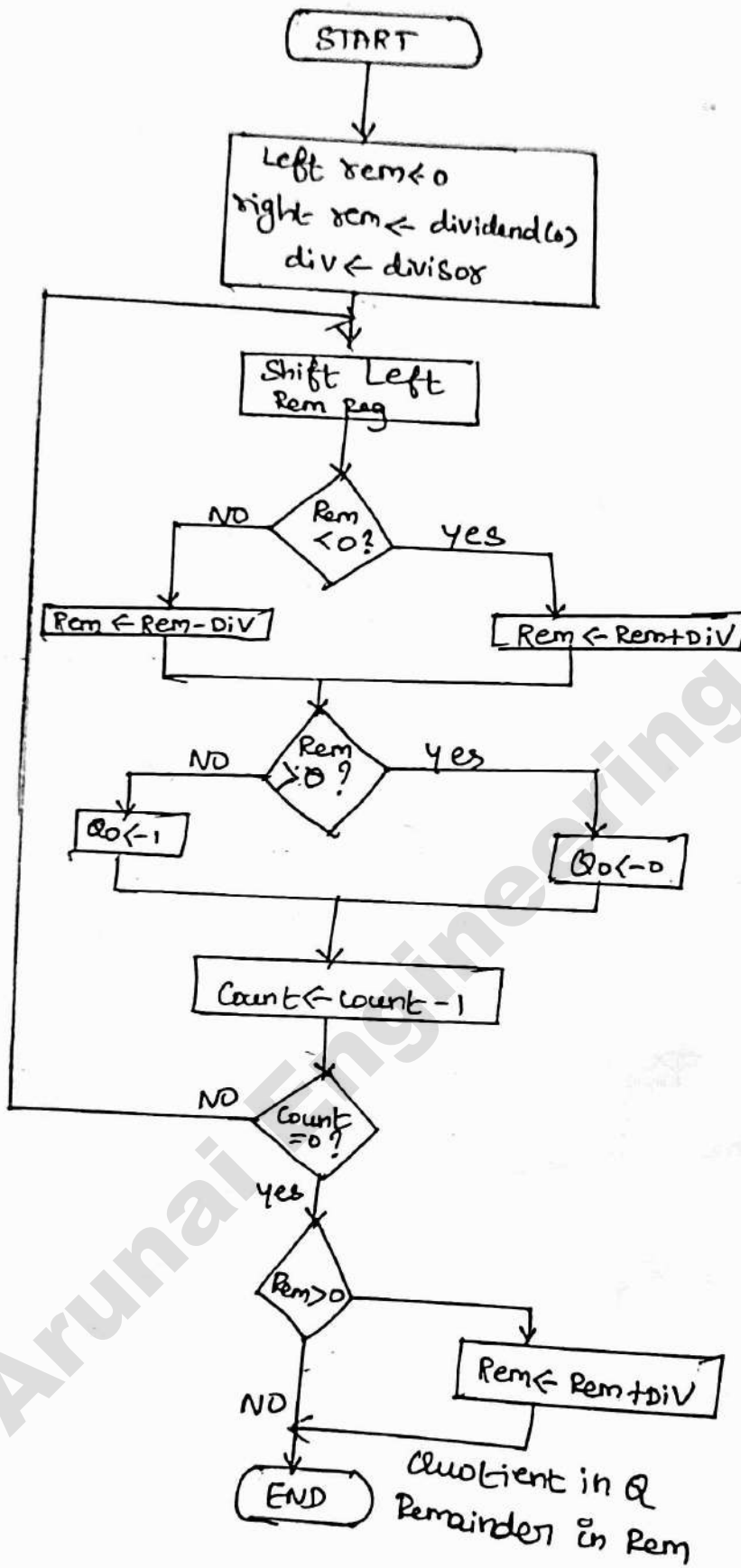
Step 3: At the end of n cycle, If the sign of Rem Reg is 1, add divisor to Rem Reg. This step is required to leave the proper positive remainder in Rem Reg at the end of n cycles.

Non-Restoring Division or Third version of Division Hardware for 32-bit Boolean number representations

Schematic Diagram of ALU circuitry



Flowchart:



Iteration	Steps	Q (Quotient)	Remainder (Rem)	
			Left of the rem reg Set to 0	Dividend
0	Initial value	0000	0000	1100
1	Shift Rem Reg	0000	0001	1000
	Left half Rem reg ← left Rem - Divisor		0011	
			1110	
	Rem < 0, restore Rem reg (+divisor), Shift Q, Q ₀ = 0	0000	0001	1000
2	Shift Rem Reg	0000	0011	0000
	Left half Rem reg ← left Rem - Divisor		0011	
			0000	
	Rem > 0, Shift Q, Q ₀ = 1	0001	0000	0000
3	Shift Rem Reg		0000	0000
	Left half Rem reg ← left Rem - Divisor		0011	
			1101	
	Rem < 0, restore Rem reg (+divisor), Shift Q, Q ₀ = 0	0010	0000	0000
4	Shift Rem Reg	0010	0000	0000
	Left half Rem reg ← left Rem - Divisor		0011	
			1101	
	Rem < 0, restore Rem reg (+divisor), Shift Q, Q ₀ = 0	0100 Quotient (4) ₁₀	0000 Remainder (0) ₁₀	0000

Non-Restoring Algorithm:

Step 1: a) If the sign of Remainder Register (Rem Reg) is 0, shift Rem Reg left logically by one binary position and subtract the divisor from the left half of the Rem Reg and place the result in left half of Rem Reg. Otherwise, shift Rem Register left and add divisor to the left of the Rem Reg.

b) If the sign of Rem Reg is 0, set Q_0 to 0, otherwise set Q_0 to 1

Step 2: Repeat steps 1 and 2 for n times

Step 3: At the end of n cycle, If the sign of Rem Reg is 1, add divisor to Rem Reg. This step is required to leave the proper positive remainder in Rem Reg at the end of n cycles

$$\text{Dividend: } (12)_{10} = 1100_2$$

$$\text{Divisor: } (3)_{10} = 0011_2$$

Iteration	Steps	Remainders neg. set's (Rem Reg)	
		Left Of Rem Reg Set to 0	Right Half Of Rem Reg Dividend/Q (Quotient)
0	Initial values	0000	1100
1	Shift Rem Reg	0001	100Q ₀
	Left half Rem reg ← left Rem - Divisor	0011	
		1110	
	Rem < 0, Q ₀ = 0	1110	1000
2	Shift Rem Reg	1101	000Q ₀
	Left half Rem reg ← left Rem + Divisor	0011	
		0000	
	Rem > 0, Q ₀ = 1	0000	0001
3	Shift Rem Reg	0000	001Q ₀
	Left half Rem reg ← left Rem - Divisor	0011	
		1101	
	Rem < 0, Q ₀ = 0	1101	0010
4	Shift Rem Reg	1010	010Q ₀
	Left half Rem reg ← left Rem + Divisor	0011	
		1101	
	Rem < 0, Q ₀ = 0	1101	0100
	Sign of Rem Reg = result in negative. Restore: Remainder add with divisor 1101 + 0011 = 0000 Remainder (0) ₁₀		Quotient (4) ₁₀

2 Explain Booth's Algorithm for the multiplication of signed two's complement numbers. [Nov/Dec-16]

Define Booth Multiplication algorithm with suitable example [May/June-2016]

Multiply the following pair of signed numbers using Booth's bit-pair recoding of the multiplier. $A = +13$ (multiplicand) and $B = -6$ (multiplier) [Nov/Dec-2014]

Booth's Algorithm:

* Booth's Algorithm is a powerful algorithm for signed-number multiplication, which treats both positive and negative numbers uniformly.

* In this algorithm, multiplier bits are encoded from right to left, before the bits are used for getting partial products.

* To encode the multiplier bits, two adjacent bits b_i (current bit) and b_{i-1} (previous bit) is examined as follows.

1. If i^{th} bit b_i is 0 and $(i-1)^{\text{th}}$ bit b_{i-1} is 1, then take b_i as +1 (i.e. for (0,1) pair, multiplicand is added to the partial product

2. If i^{th} bit b_i is 1 and $(i-1)^{\text{th}}$ bit b_{i-1} is 0, then take b_i as -1 (i.e., for (1,0) pair, multiplicand is subtracted to the partial product

3. If i^{th} bit b_i is 0 and $(i-1)^{\text{th}}$ bit b_{i-1} is 0, then take b_i as 0

4. If i^{th} bit b_i is 1 and $(i-1)^{\text{th}}$ bit b_{i-1} is 1, then take b_i as 0 (i.e.) For (0, 0) and (1, 1) pair, then neither addition nor subtraction is performed.

5. When least significant bit $b_0 = 1(0)$, assume that it had b_{i-1} as 0 (implied zero) thus take $b_0 = -1(0)$

Multiplier		version of multiplicand selected by bit i
Bit i	Bit $i-1$	
0	0	$0 \times M$
0	1	$+1 \times M$
1	0	$-1 \times M$
1	1	$0 \times M$

Booth Multiplier Recoding Table

* Booth's algorithm effectively skips over sequence of 1's and 0's in multiplier. As a result, the total number of addition/subtraction steps required to multiply two numbers decrease

* The process of inspecting the multiplier bits required by Booth's algorithm can be viewed as encoding the multiplier using three digits 0, +1 and -1 where 0 means shift the partial product to the right (ie no addition or subtraction is performed)

* While +1 means add multiplicand before shifting and -1 means subtract multiplicand from the partial product before shifting.

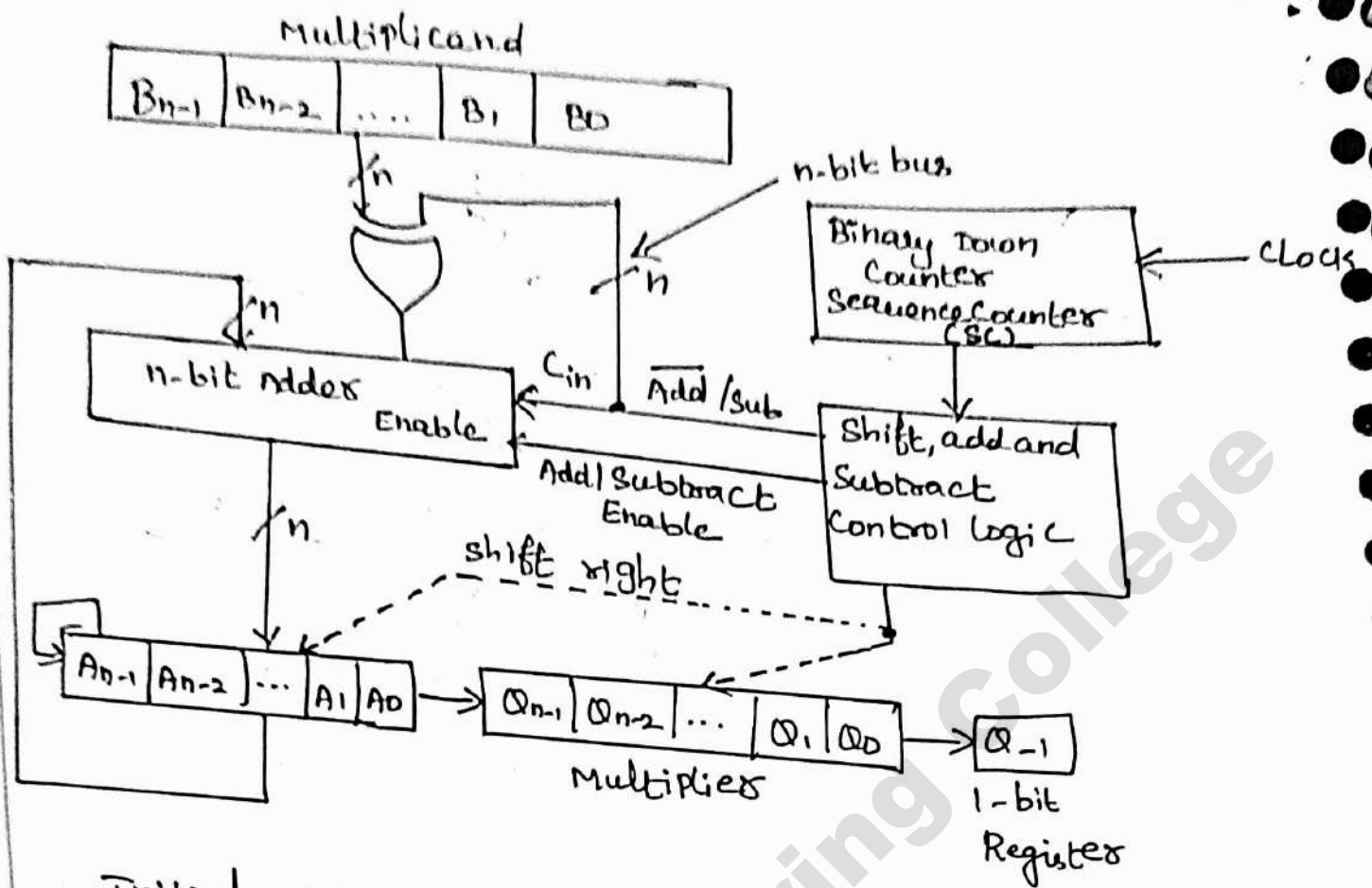
* The number thus produced is called as signed digit number and this process of converting a multiplier into a signed digit form is called as multiplier recoding

Hardware Implementation of Booth's Algorithm:

* It consists of n-bit adder, shift and subtract control logic and four registers A, B, Q and Q-1

* Multiplier and multiplicand are loaded into register Q and register B respectively and register A and Q-1 are initially set to 0.

* The Sequence Counter, SC is set to a number n equal to the number of bits in the multiplier.



Initial settings: $A \leftarrow 0$ and $Q_{-1} = 0$

The n -bit adder performs addition of two inputs.

- i) A register
- ii) Multiplicand.

↳ In case of addition, $\overline{\text{Add/Sub}}$ line is 0, therefore $C_{in} = 0$ and multiplicand is directly applied as a second input to the n -bit adder.

↳ In case of subtraction, $\overline{\text{Add/Sub}}$ line is 1, therefore $C_{in} = 1$ and multiplicand is complemented and then applied to the n -bit adder. As a result, the 2's complement of multiplicand is added in the A register.

↳ The shift, add and subtract control logic scans bits Q_0 and Q_{-1} one at a time and generates the control signals.

Flowchart of Booth's Multiplication Algorithm

Algorithm:

Step 1: Load $A = 0$, $Q_{-1} = 0$

$B =$ Multiplicand

$Q =$ Multiplier

$SC = n$

Step 2: check the status of $Q_0 Q_{-1}$

if $Q_0 Q_{-1} = 10$ perform $A \leftarrow A - B$

if $Q_0 Q_{-1} = 01$ perform $A \leftarrow A + B$

Step 3: Arithmetic Shift right: A, Q, Q_{-1}

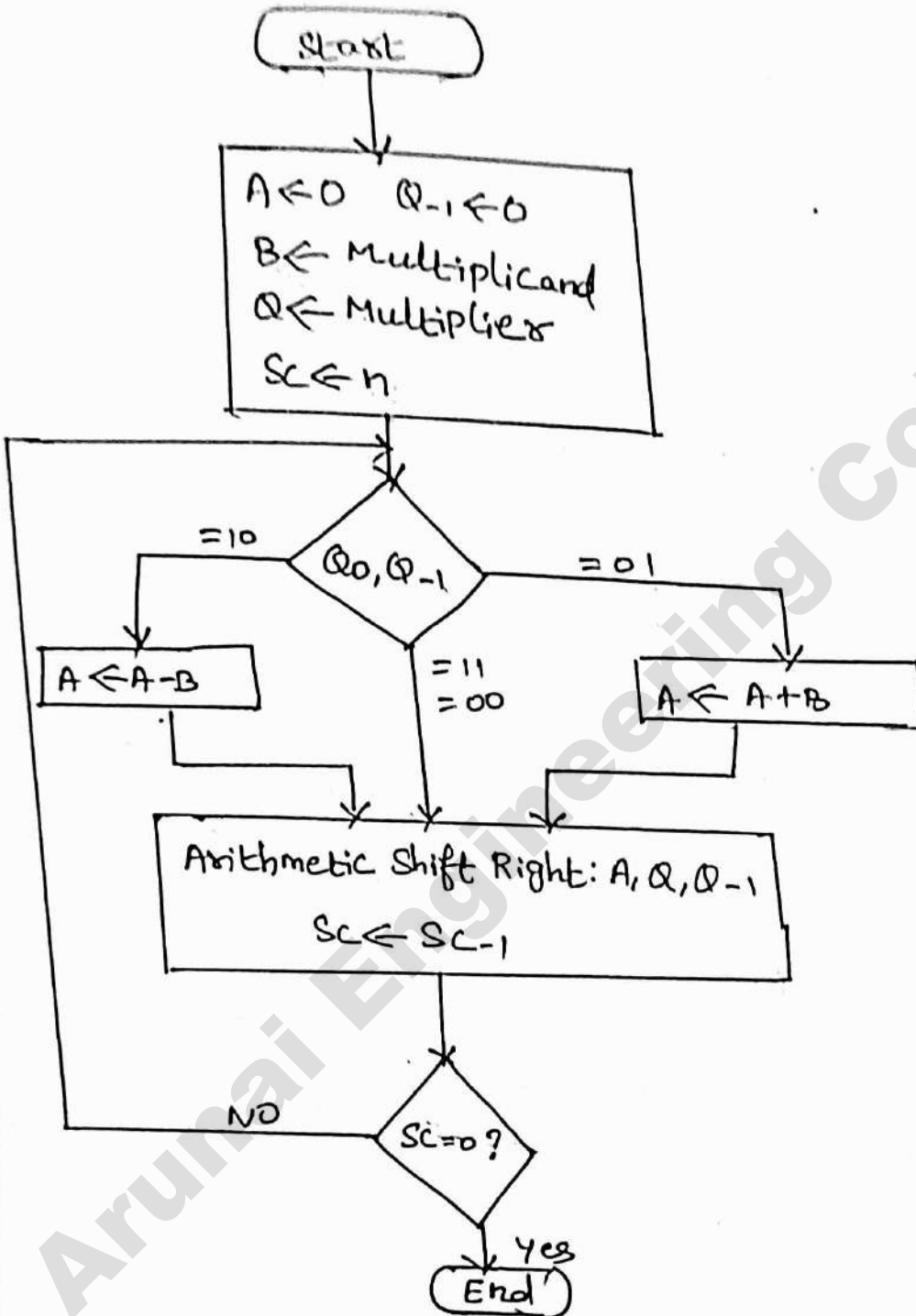
Step 4: Decrement Sequence Counter if not zero, repeat

Step 2 through 4

Step 5: Stop

Booth's Algorithm for signed Multiplication

Flowchart:



Example:

Multiply the following pair of signed numbers using Booth's bit-pair recoding of the multiplier A = +13 (Multiplicand) and B = -6 (Multiplier)

Multiplicand:

0	1	1	0	1
---	---	---	---	---

 +13

Multiplier:

0	0	1	1	0
1	1	0	1	0

 6
-6 (2's complement form)

1	1	0	1	0
0	-1	+1	-1	0

0 implied zero multiplier
Recoded multiplier

				0	1	1	0	1	Multiplicand
				0	-1	+1	-1	0	Recoded multiplier
Sign Extension →	0	0	0	0	0	0	0	0	
+	1	1	1	1	0	0	1	1	2's complement of multiplicand
+	0	0	0	1	1	0	1		
+	1	1	0	0	1	1			2's complement of multiplier
+	0	0	0	0					
	1	1	0	1	1	0	0	1	0

(-78) in 2's complement form

3: Explain how floating point addition is carried out in a computer system. Give an example for a binary

floating point addition [APR/MAY-2015] [APR/MAY-17]

Design an arithmetic element to perform the basic floating point operation [APR/MAY-17]

Floating Point Addition:

Floating Point:

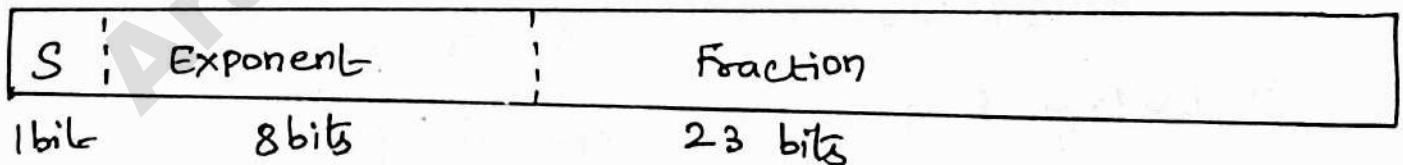
Computer arithmetic represents numbers in which the binary point is not fixed called floating point.

Floating point representation

* A designer of a floating point representation must find a compromise between the size of the fraction and the size of the exponent.

* Fraction is the value generally between 0 and 1 it is called the mantissa

* Exponent is the numerical representation system of floating point arithmetic



MIPS floating point number representation

where S = sign of the floating point number
1 - negative, 0 - positive

General floating point numbers are of the form

$$(1)^S F 2^E$$

F - involves the value in the fraction field

E - involves the value in the exponent field

Overflow (Floating point)

* overflow is a situation in which a positive exponent becomes too large to fit in the exponent field.

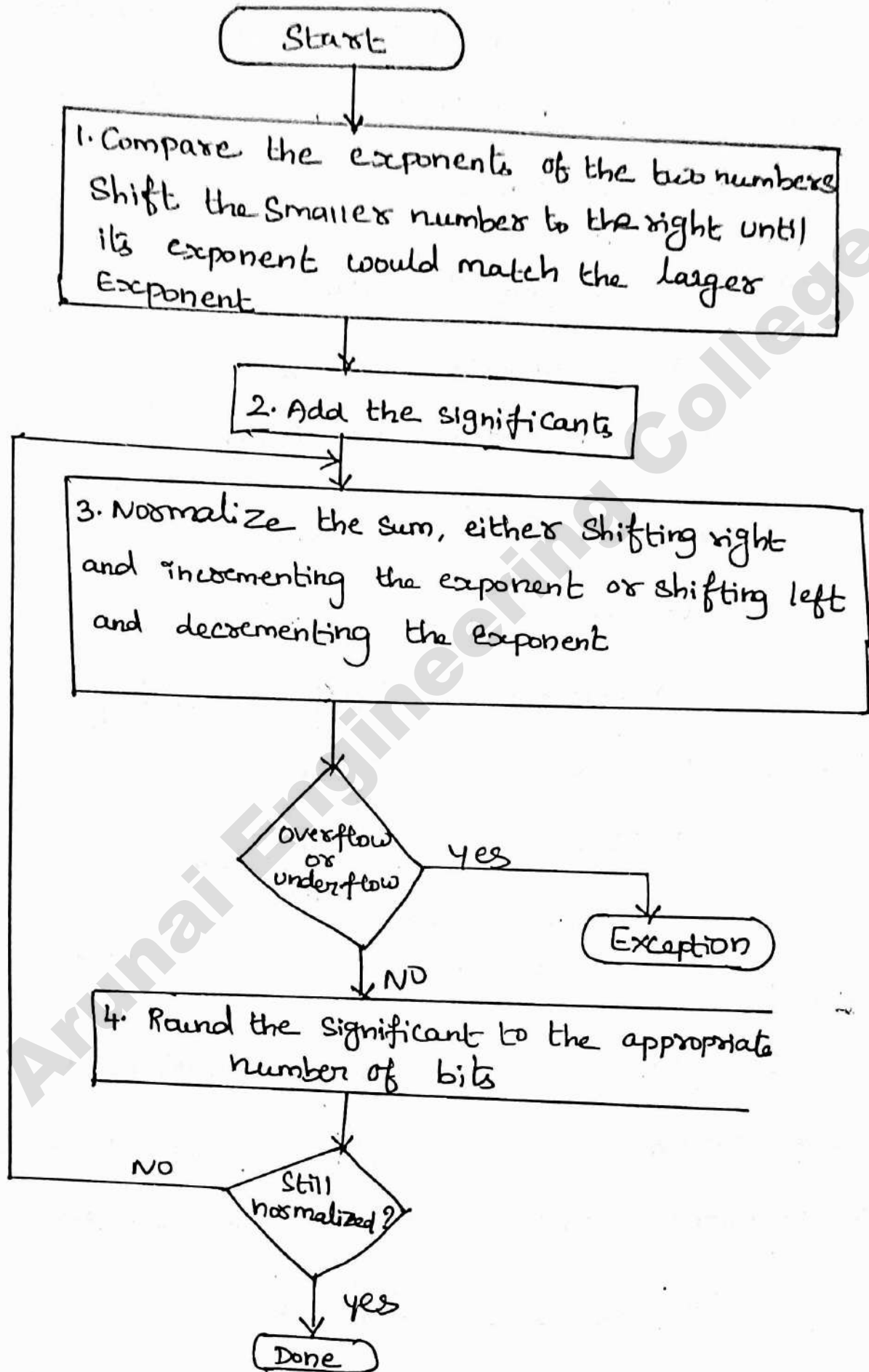
* overflow means the exponent is too large to be represented in the exponent field.

Underflow (Floating Point)

* Underflow is a situation in which a negative exponent becomes too large to fit in the exponent field.

* Underflow occurs when the negative exponent is too large to fit in the exponent field.

Flowchart for Floating Point Addition:



Binary floating point addition:

Perform floating point addition using the numbers 0.5_{10} and 0.4375_{10} use the floating point addition algorithm.

Assume 4 bits of precision

$$\text{First number} = 0.5_{10}$$

$$\text{Second number} = 0.4375_{10}$$

Let 0.5_{10} can also be written as

$$0.5_{10} = \frac{1}{2}_{10} = \frac{1}{2^1}_{10}$$

$$0.1_{10} = 0.1_{10} 2^0 = 1.000_{10} 2^1$$

0.4375_{10} can also be written as

$$0.4375_{10} = \frac{7}{16}_{10} = \frac{7}{2^4}_{10}$$

$$0.011_{10} = 0.011_{10} 2^0 = 1.110_{10} 2^2$$

Step 1:

The significant of the numbers with the lesser exponent ($1.11_{10} 2^2$) is shifted right until its exponent matches the larger number

$$1.110_{10} 2^2 = 0.111_{10} 2^1$$

Step 2:

Add the significant

$$1.000 \times 2^1 + (0.111 \times 2^{-1}) \times 0.001 \times 2^1$$

Step 3:

↳ normalize the sum, checking the overflow or underflow

$$0.001 \times 2^1 = 0.010 \times 2^2 = 0.100 \times 2^3 = 1.000 \times 2^4$$

Since $127 \leq 126$, there is no overflow or underflow

Step 4:

Round the sum

$$1.000 \times 2^4$$

The sum already fits exactly in 4 bits, so there is no change to the bits due to rounding

$$1.000 \times 2^4 = \frac{0.0001000}{2^4} \times 2^4 = 0.0001 \times 2^0$$

$$= \frac{1}{2^4} \text{ ten}$$

$$= \frac{1}{16} \text{ ten}$$

$$\boxed{= 0.0625 \text{ ten}}$$

$$1.000 \times 2^4$$
$$0.0625 \times 2^0$$

$$0.1 \times 2^1$$
$$0.01 \times 2^2$$
$$0.001 \times 2^3$$
$$0.0001 \times 2^4$$

Explain the sequential version of multiplication algorithm and its hardware. [APR/MAY-2015]

Explain in detail about the multiplication algorithm with suitable example and diagram [NOV/DEC-2015].
[APR/MAY-2017]

Multiplication:

- * The first operand is called Multiplicand
- * The second operand is called Multiplier
- * Final Result is called Product

Basic rules for Binary digits multiplication

$$0 \times 0 = 0$$

$$0 \times 1 = 0$$

$$1 \times 0 = 0$$

$$1 \times 1 = 1$$

Eg multiplying 1010 by 1000

Multiplicand	1010
Multiplier	× 1000
	0000
	0000
	0000
	1010
Product	1010000

* Take the digits of the multiplier one at the time from right to left, multiplying the

multiplicand by the single digit of the multiplier.

* Shifting the intermediate product one digit to the left of the earlier intermediate products.

* The number of digits in the product is larger than the number of digits in either the multiplicand or the multiplier.

* Multiplication is usually implemented by some form of repeated addition.

* To compute $x \times y$ is to add the multiplicand y to itself x times, where x is the multiplier.

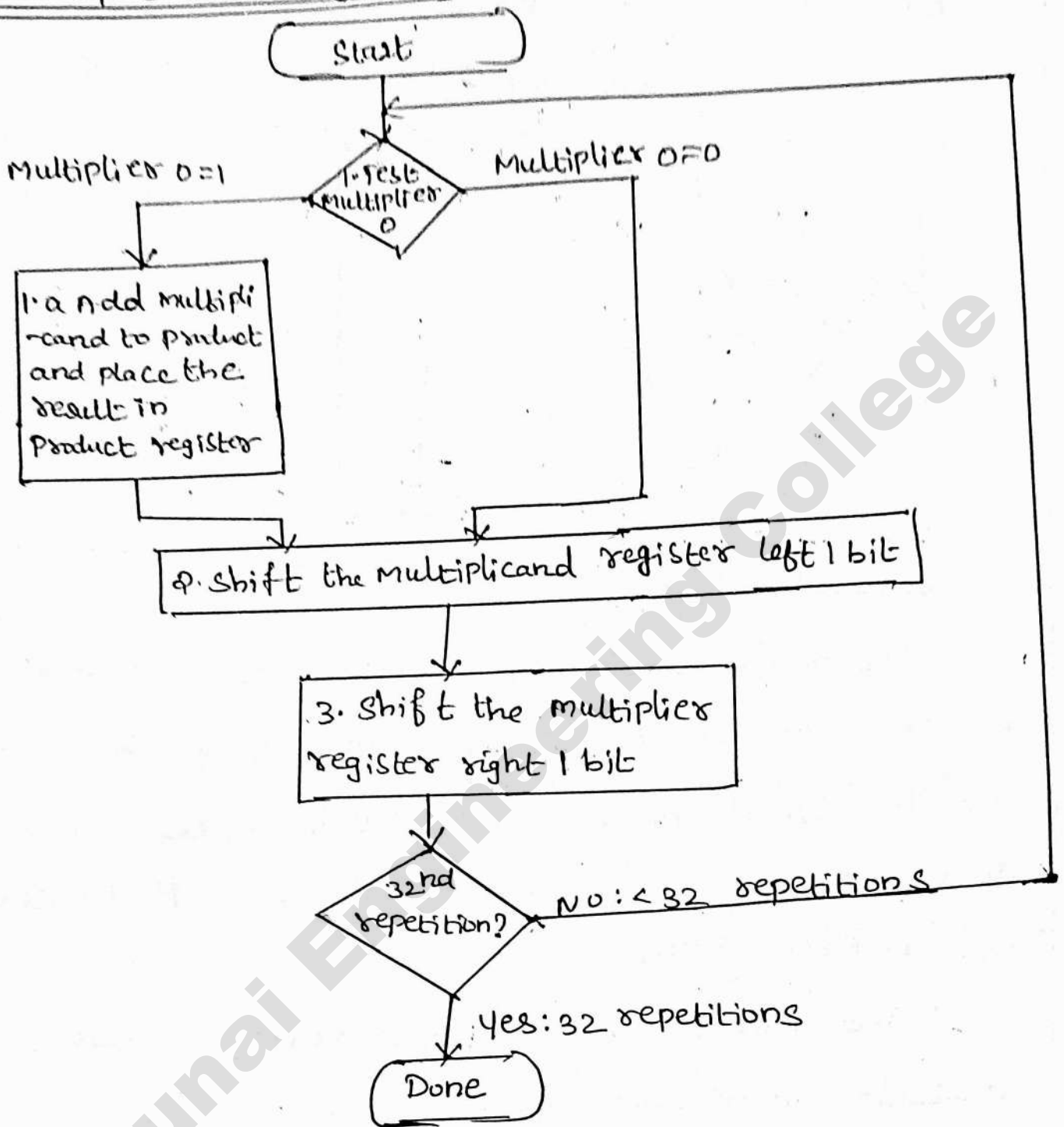
Example:

$$4 \times 8 = 32$$

Here $x=4$ and $y=8$ now we add y itself x number of times, that will produce the product

$$\begin{array}{r} 8 \\ 8 \\ 8 \\ 8 \\ \hline 32 \\ \hline \end{array}$$

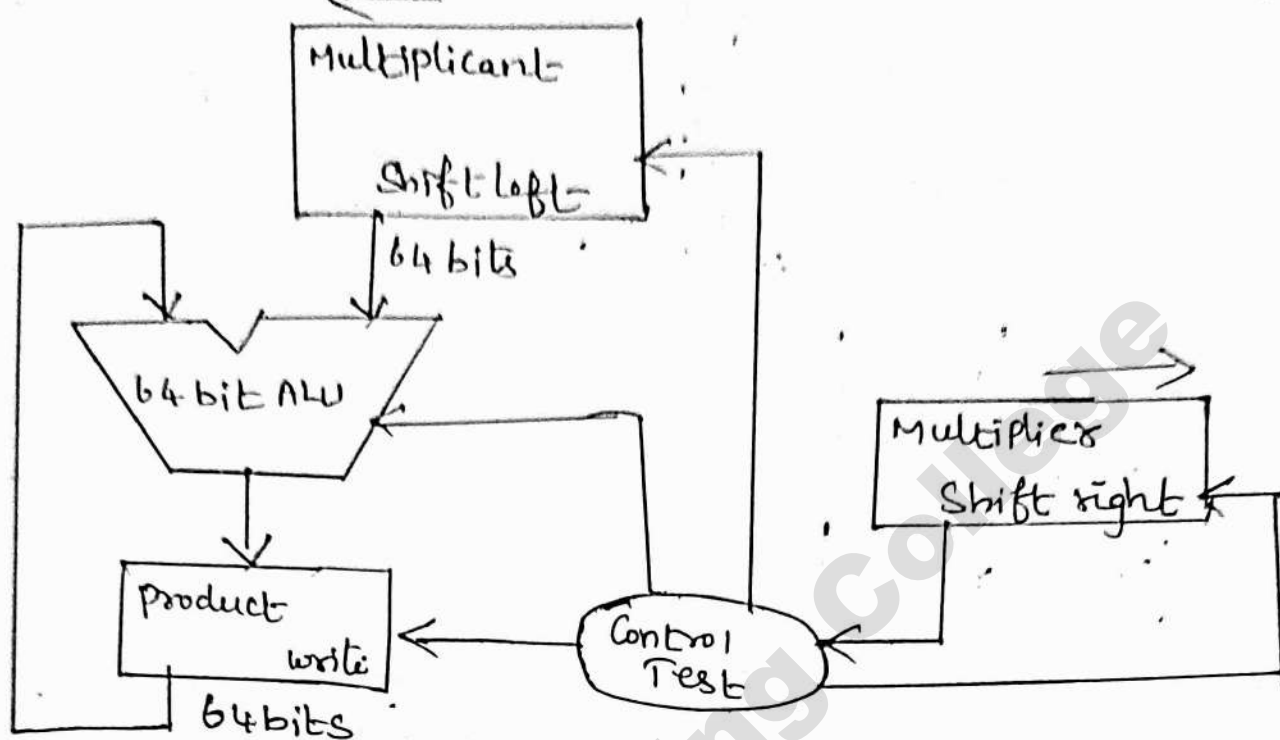
Multiplication Algorithm:



The first Multiplication Algorithm:

Step 1: The least significant bit of the multiplier (multiplier 0) determines whether the multiplicand is added to the product register

Multiplication Hardware:



* The multiplicand registers, ALU and product registers are all 64 bit and multiplier is 32 bits

* The 32 bit multiplicand starts in the right half of the multiplicand registers and is shifted left 1 bit on each step.

* The multiplier algorithm starts with the product initialized to 0.

* The multiplier is shifted in the opposite direction at each step.

* Control decides when to shift the multiplicand and multiplier registers and when to write new values into the product registers.

Step 2: The left shift, has the effect of moving the intermediate operands to the left.

Step 3: The shift right gives the next bit of the multiplier to examine in the following iteration.

Step 4: These three steps are repeated 32 times to obtain the product.

Clock cycles for multiplication algorithm

* This algorithm requires almost 100 clock cycles to multiply two 32 bit numbers.

* multiply take multiple clock cycle without affecting performance

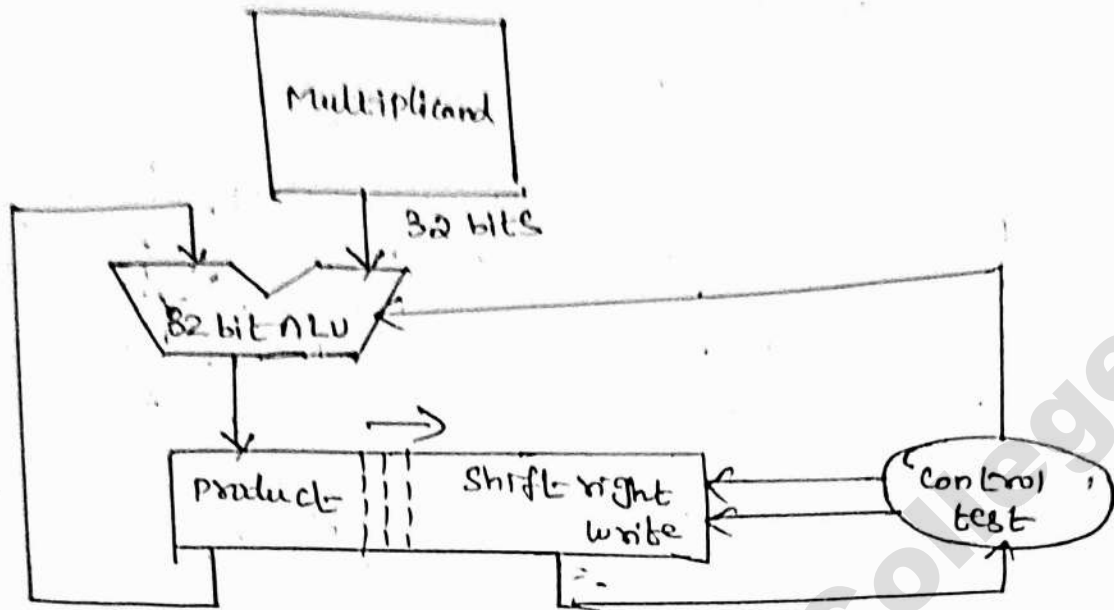
* To increase the speed of process by performing the operations in parallel manner.

* The multiplier and multiplicand are shifted while the multiplicand is added to the product if the multiplier bit is 1.

Refined version of multiplication Hardware:

* The hardware is used to ensure that it tests and right bit of the multiplier and gets the pre-shifted version of the multiplicand

* Hardware will take 1 clock cycle per step



* The multiplicand register, ALU and multiplier registers are all 32-bits, the product register only has 64 bits. The product is shifted right.

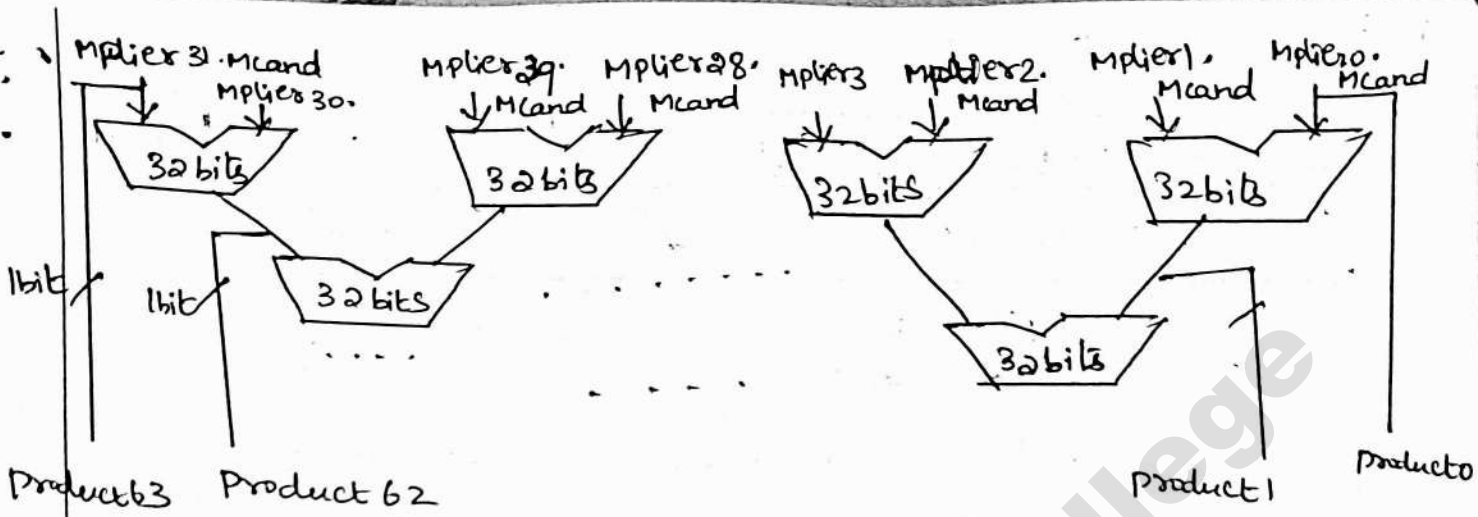
* The separate multiplier register also disappeared and the multiplier is placed instead in the right half of the product register

Signed Multiplication:

* For signed multiplication first convert the multiplier and multiplicand into positive numbers and then remember the original signs

Faster Multiplication:

* Moore's Law has provides much more resources is there the hardware designers can build much faster multiplication hardware.



* Faster multiplications are possible by providing one 32 bit address for each bit of the multiplier

* one input is the multiplicand ANDed with a multiplier bit and the other is the output of a prior adder.

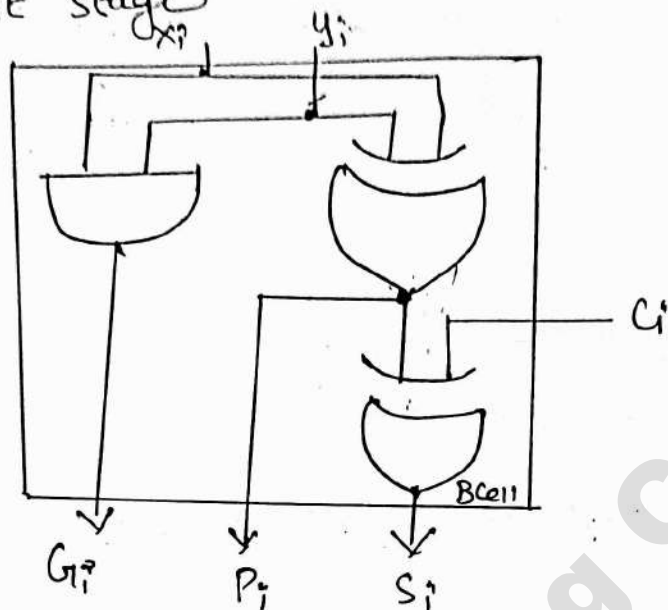
Example:

Using 4 bit numbers to save space, multiply 2ten 3ten
 or 0010two 0011two

Solution :

Iteration	Step	Multiplier	Multiplicand	Product
0	Initial values	001 1	0000 0010	0000 0000
1	1: 1 \Rightarrow prod = prod + M _{cand} 2: shift left multiplicand 3: shift right multiplier	0011 0011 ↓ ↓ ↓ 0 0 0 1	0000 0010 ↓ ↓ ↓ ↓ 0000 0100 0000 0100	0000 0010 0000 0010 0000 0010
2	1: 1 \Rightarrow prod = prod + m _{cand} 2: shift left multiplicand 3: shift right multiplier	0001 0001 ↓ ↓ ↓ 0 0 0 1	0000 0100 // // // // 0000 1000 0000 1000	0000 0110 0000 0110 0000 0110
3	1: 0 \Rightarrow no operation 2: shift left multiplicand 3: shift right multiplier	0000 0000 ↓ ↓ ↓ 0 0 0 0	0000 1000 0001 0000 0001 0000	0000 0110 0000 0110 0000 0110
4	1: 0 \Rightarrow no operation 2: shift left multiplicand 3: shift right multiplier	0000 0000 0000	0001 0000 ↓ ↓ ↓ ↓ ↓ ↓ 0010 0000 0010 0000	0000 0110 0000 0110 0000 0110

The following basic cell can be used in each bit stage



* Expanding C_i in terms of $i-1$ and substitute into C_{i+1} , then following expression will be

$$C_{i+1} = G_i + P_i G_{i-1} + P_i P_{i-1} C_{i-1}$$

* The final expression for carry variable is

$$C_{i+1} = G_i + P_i G_{i-1} + P_i P_{i-1} G_{i-2} + \dots + P_i P_{i-1} \dots P_1 G_0 + P_i P_{i-1} \dots P_0 C_0$$

* Thus all carries can be obtained three gate delays after the input signals x, y, C_0 are applied

5

Briefly Explain Carry lookahead adder [Nov/Dec-2014]

* Carry Look ahead Adder (CLA) is a fast adder which speed up the generation of carry signals.

* It uses two functions

i) Carry generate, G_i

ii) Carry propagate, P_i

The Logic Expression for a sum S_i & carry out C_{i+1} of stage i are

$$S_i = x_i \oplus y_i \oplus C_i$$

$$C_{i+1} = x_i y_i + x_i C_i + y_i C_i$$

$$C_{i+1} = x_i y_i + (x_i + y_i) C_i$$

$$C_{i+1} = G_i + P_i C_i$$

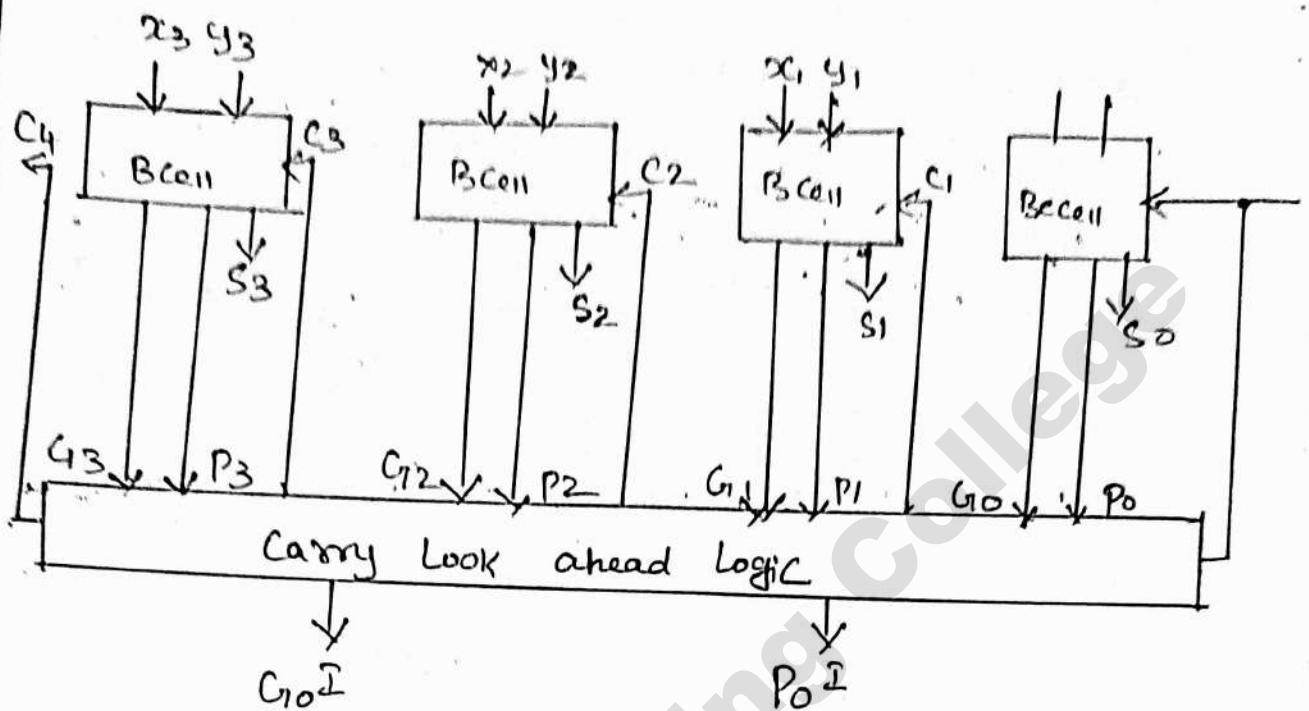
where $G_i = x_i y_i$

$P_i = x_i + y_i$

* If the generate function for stage i is equal to 1, then $C_{i+1} = 1$

* The Propagate function means that an input carry will produce an output carry when either x_i or y_i is 1

Design of 4-bit adder



The carries can be implemented as

$$C_1 = G_0 + P_0 C_0$$

$$C_2 = G_1 + P_1 G_0 + P_1 P_0 C_0$$

$$C_3 = G_2 + P_2 C_2$$

$$C_3 = G_2 + P_2 G_1 + P_2 P_1 G_0 + P_2 P_1 P_0 C_0$$

$$C_4 = G_3 + P_3 C_3$$

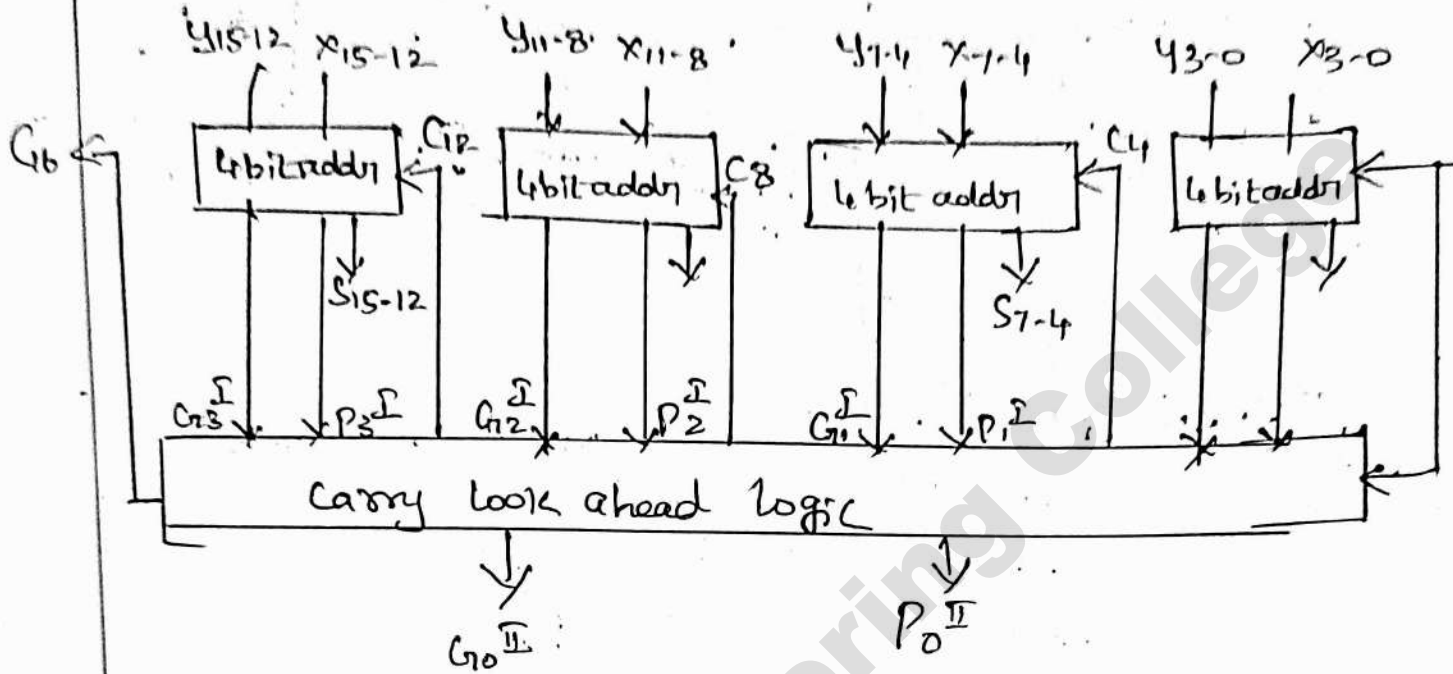
$$C_4 = G_3 + P_3 G_2 + P_3 P_2 G_1 + P_3 P_2 P_1 G_0 + P_3 P_2 P_1 P_0 C_0$$

* The carries are implemented in carry look-ahead logic. An adder of this form called carry look-ahead adder.

* Delays through the adder is 3 gate delays for all carry bits.

* 4 gate delays for an sum bits which is less than ripple carry adder

Construction of 16 bit adder built from 4 bit adder



* By using look ahead logic, higher level carry look ahead circuit is developed.

* The propagate and generate functions will be

$$P_0^I = P_3 P_2 P_1 P_0$$

$$G_0^I = G_3 + P_3 G_2 + P_3 P_2 G_1 + P_3 P_2 P_1 G_0$$

The carry C_{16} will be

$$C_{16} = G_3^I + P_3^I G_2^I + P_3^I P_2^I G_1^I + P_3^I P_2^I P_1^I G_0^I + P_3^I P_2^I P_1^I P_0^I C_0^I$$

6. What is meant by subword parallelism? Explain

[Apr/May-17]

* Parallelism will improve the performance of computers.

* A subword is a lower precision unit of data contained within a word.

* In subword parallelism, pack multiple subwords into a word and then process whole word with the appropriate subword boundaries.

* This technique results in parallel processing of subwords.

* Subword parallelism also known as data level parallelism or vector or SIMD (single instruction multiple data) processing.

* Since the same instruction applies to all the subwords within the word, this is a form of small scale SIMD processing.

* It is possible to apply subword parallelism to non contiguous subwords of different sizes within a word.

* In practice, however, implementations are much simpler if we allow only a few subword size and if a single instruction operates on contiguous

Subwords that are all of the same size.

* Data parallel programs that benefit from subword parallelism tend to process data that are of the same size.

Eg If the word size is 64 bits, some useful subword sizes are 8, 16 and 32 bits

* An instruction operates on 8-bit subwords, four 16-bit subwords, two 32-bit subwords or one 64-bit subword in parallel.

* Data parallelism refers to an algorithm execution of the same program module on different sets of data.

* Subword parallelism is an efficient and flexible solution for media processing, because the algorithms exhibit a great deal of data parallelism on lower precision data.

Advantage:

1. It allows general-purpose processors to exploit wider word sizes even when not processing high-precision data.

2. The processor can achieve more sub-word parallelism on lower precision data rather than wasting much of the word-oriented data paths and registers.

7

i) Add the numbers $(0.5)_{10}$ and $(0.4375)_{10}$ using the floating point addition. [Nov/Dec-17]

Solution:

Assume 4 bits of precision

First number = 0.5_{10}

Second number = 0.4375_{10}

Convert the decimal numbers into binary

$$0.5_{10} = 0.1_2 \times 2^0 = 1.000_2 \times 2^{-1}$$

$$0.5 \times 2 = 1.000 \times 2^1$$

$$0.4375 = 0.0111_2 \times 2^0 = 1.110_2 \times 2^{-2}$$

Step 1:

Equalize Exponents

$$0.4375 \times 2 = 0.875$$

$$0.875 \times 2 = 1.75$$

$$0.75 \times 2 = 1.50$$

$$0.50 \times 2 = 1.00$$

$$\begin{array}{r} 0.0111_2 \times 2^0 \\ \underline{0.111 \times 2^1} \\ 1.11 \end{array}$$

* The significant of the numbers with the lesser exponent ($1.110_2 \times 2^{-2}$) is shifted right until its exponent matches the larger number

$$1.110_2 \times 2^{-2} = 0.111_2 \times 2^{-1}$$

Step 2:

Add the significant

7 i) Add the numbers $(0.5)_{10}$ and $(0.4375)_{10}$ using the floating point addition. [Nov/Dec-17]

Solution:

Assume 4 bits of precision

First number = 0.5_{10}

Second number = 0.4375_{10}

Convert the decimal numbers into binary

$$0.5_{10} = 0.1_2 \times 2^0 = 1.000_2 \times 2^{-1}$$

$$\begin{array}{l} 0.5 \times 2 \\ \hline 1.000 \times 2^{-1} \end{array}$$

$$0.4375 = 0.0111_2 \times 2^0 = 1.110_2 \times 2^{-2}$$

Step 1:

Equalize Exponents

$$\begin{array}{l} 0.4375 \times 2 = 0.875 \\ 0.875 \times 2 = 1.75 \\ 0.75 \times 2 = 1.50 \\ 0.50 \times 2 = 1.00 \end{array}$$

$$\begin{array}{l} 0.0111_2 \times 2^0 \\ \hline 0.111 \times 2^{-1} \\ \hline 1.11 \end{array}$$

* The significant of the number with the lesser exponent ($1.110_2 \times 2^{-2}$) is shifted right until its exponent matches the larger number

$$1.110_2 \times 2^{-2} = 0.111_2 \times 2^{-1}$$

Step 2:

Add the significant

$$1.000_2 \times 2^{-1} + 0.111_2 \times 2^{-1} = 0.001_2 \times 2^{-1}$$

Step 3:

Normalization :-

$$\begin{array}{r} 1.000 \times 2^{-1} \\ 0.111 \times 2^{-1} \\ \hline 1.111 \times 2^{-1} \end{array}$$

* Normalize the sum, checking the overflow

or underflow.

$$0.1111 \times 2^0$$

$$0.001_2 \times 2^{-1} = 0.010_2 \times 2^{-2} = 0.100_2 \times 2^{-3} = 1.000_2 \times 2^{-4}$$

* Since $127 \leq 4 \leq 126$, there is no overflow or underflow

Step 4:

Rounding:

↳ Round the sum

$$\boxed{1.000_2 \times 2^{-4}}$$

$$0.1111$$

$$\frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{16}$$

$$\begin{array}{r} 8+4+2+1 \\ \hline 16 \\ \hline 15 \\ \hline 16 \end{array}$$

↳ The sum already fits exactly in 4 bits, so there is no change to the bits due to rounding.

$$1.000_2 \times 2^{-4} = \overset{2^4}{\underbrace{0.0001000}_2} = 0.0001_2$$

$$0.1111 \times 2^0$$

$$= 1/2^4$$

$$= 1/16_{10}$$

$$= 0.0625_{10}$$

$$\begin{array}{r} \underbrace{1111}_{15} \\ \frac{1}{2} \quad \frac{1}{4} \quad \frac{1}{8} \quad \frac{1}{16} \\ \hline 0.0625 \end{array}$$

$$= \frac{1}{16}$$

This sum is what we expect from adding 0.5_{10} and 0.4375_{10}

7
ii) Multiply the numbers $(0.5)_{10}$ and $(0.4375)_{10}$ using the floating point multiplication. [Nov/Dec -17]

Solution:

* In binary, the task is multiplying

$$0.5_{10} = 0.1_2 \times 2^0 = 1.000_2 \times 2^{-1}$$

$$0.4375_{10} = 0.0111_2 \times 2^0 = 1.110_2 \times 2^{-2}$$

Step 1: Calculate the exponent

* Adding the exponents without bias

$$1 + (-2) = 3 \quad -1 - 2 = -3$$

Step 2: Multiply the significands

$$1.000_2 = 0.5_{10} \rightarrow M_1$$

$$1.110_2 = 0.4375_{10} \rightarrow M_2$$

$$\begin{array}{r} 1.000_2 \\ \times 1.110_2 \\ \hline 0000 \\ 1000 \\ 1000 \\ 1000 \\ \hline 1110000 \end{array}$$

* The product is $1.110000_2 \times 2^3$, but we need to keep it to 4 bits only, so it is $1.110_2 \times 2^3$

Step 3: Normalization

* Now check the product is normalized format and also check the exponent for overflow or underflow.

* The product is already normalized and there is no overflow or underflow. $1272 - 32 = 126$

Step 4: Rounding

* Rounding the product makes no change.

$$1.110_2 \times 2^3$$

Step 5: Sign

* Signs of the original operands differ, make the sign of the product negative. Hence the product is

$$-1.110_2 \times 2^3$$

* Converting to decimal form to check our results

$$\begin{aligned} 1.110_2 \times 2^3 &= 0.\overset{2^3}{001110}_2 = 0.00111_2 \\ &\quad \downarrow \\ &1.11_2 = 7/2_{10} = 7/32_{10} \\ &= 0.21875_{10} \end{aligned}$$

* The product of 0.5_{10} and 0.4375_{10} is $= 0.21875_{10}$

$$\begin{array}{r} 0.5 \times 0.4375 \\ \hline 25 \\ 15 \\ \hline 0.21875 \end{array}$$

8i) Perform $x+y$ and $y-x$ using 2's complements for given the two binary numbers

$$x = 0000\ 1011\ 1110\ 1111 \quad \text{and} \quad y = 1111\ 0010\ 1001\ 1101$$

Ans:

[APR/MAY -18]

Perform $x+y$

$$\begin{array}{r} x = 0000\ 1011\ 1110\ 1111 \\ y = 1111\ 0010\ 1001\ 1101 \\ \hline 1111\ 1110\ 1000\ 1100 \end{array}$$

Perform $y-x$

Step 1: Find 2's complement of x ;

Step 2: Add the result with y ;

i) Find 2's complement of x :

$$x = 0000\ 1011\ 1110\ 1111$$

$$\begin{array}{l} \text{1's complement of } x = 1111\ 0100\ 0001\ 0000 \\ \text{2's complement} \end{array}$$

$$\begin{array}{r} 1111\ 0100\ 0001\ 0000 \\ \hline x = 1111\ 0100\ 0001\ 0001 \end{array}$$

ii) Add with y

$$y = 1111\ 0010\ 1001\ 1101$$

$$x = 1111\ 0100\ 0001\ 0001$$

$$\begin{array}{r} 1111\ 0100\ 0001\ 0001 \\ \hline 1110\ 0110\ 1010\ 1110 \end{array}$$

Discard

8ii) Multiply the following signed 2's complement numbers using the Booth Algorithm. $A = 001110$ and $B = 111001$ where A is multiplicand and B is multiplier [Apr/May-18]

Ans :

$A = 001110 \leftarrow$ Multiplicand

$B = 111001 \leftarrow$ Multiplier

Recode the multiplier using Booth Algorithm

1 1 1 0 0 1 0 \leftarrow Implied zero

↖ ↗ ↖ ↗ ↖ ↗ ↖ ↗
1 1 1 0 0 1 0

0 0 -1 0 +1 -1 Recoded multiplier

0 0 1 1 1 0 X

0 0 -1 0 +1 -1

	1	1	1	1	1	1	1	1	0	0	1	0
	0	0	0	0	0	0	0	1	1	1	0	
	0	0	0	0	0	0	0	0	0	0		
	1	1	1	1	1	0	0	1	0			
	0	0	0	0	0	0	0	0				
	0	0	0	0	0	0	0					
Result \rightarrow	1	1	1	1	1	0	0	1	1	1	1	0

9 Multiply the following signed numbers using Booth algorithm. $A = (-34)_{10} = (1011110)_2$ and $B = (22)_{10} = (0010110)_2$ where B is multiplicand and A is multiplier. [Nov/Dec-18]

Ans:-

$A = (1011110)_2$ Multiplier

$B = (0010110)_2$ Multiplicand

- * Recode the multiplier using Booth's algorithm
- * Before recoding the multiplier add the implied zero to the immediate right of LSB.

(ie) A is multiplier

1 0 1 1 1 1 0 0 ← Implied zero

Using the recoding table we can recode the multiplier with implied zero as follows

↖ ↗ ↖ ↗ ↖ ↗ ↖ ↗ ↖ ↗
1 0 1 1 1 1 0 0

-1 +1 0 0 0 -1 0 ← Recoded multiplier

0 0 1 0 1 1 0 ← Multiplicand

x -1 +1 0 0 0 -1 0 ← Recoded Multiplier

	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	1	1	1	1	1	1	1	1	0	1	0	1	0	0
	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	0	0	0	0	1	0	1	1	0					
	1	1	1	0	1	0	1	0						
1	1	1	1	1	0	1	0	0	0	1	0	1	0	0

↑ Discard ↑ Sign(-)

← 2's complement of multiplicand

← 2's complements of multiplicand

← Binary Answer

$= (0100010100)_2$

↓
 $= 748_{10}$

Arunai Engineering College

10

Explain floating point addition algorithm with a neat block diagram. [Nov/Dec-18]

Floating point addition:

Step 1: Exponent of the two floating-point numbers must be made equal. To do this, compare the exponents of the two numbers.

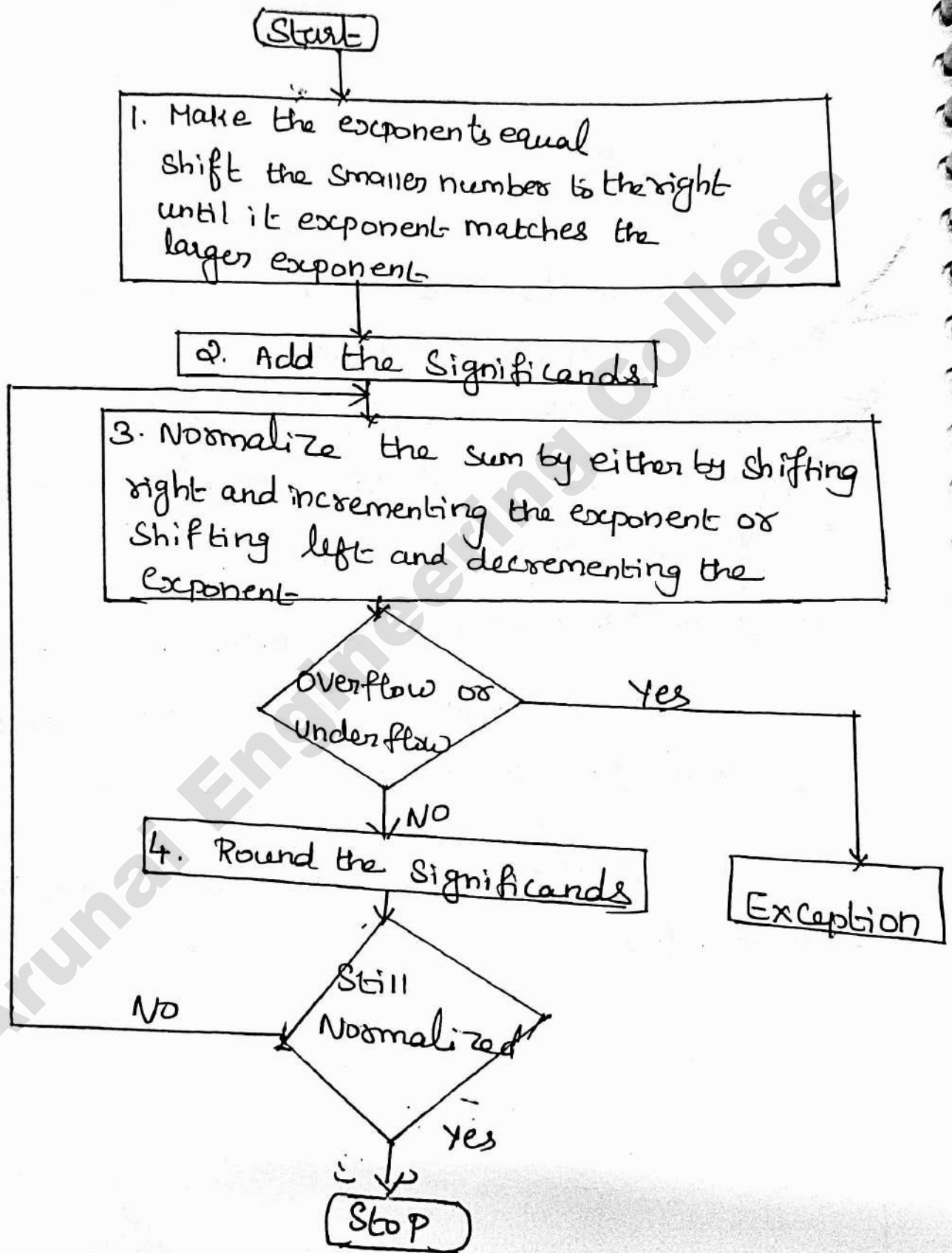
* Shift the small number to the right until its exponent is matched with the larger exponent.

Step 2: Add the two significands.

Step 3: Normalize the sum, either by shifting right and incrementing the exponent or shifting left and decrementing the exponent.

Step 4: Round the significand to the appropriate number of bits.

Floating-Point Addition Algorithm



Perform Floating-point addition on the following two decimal numbers

$$9.999_{10} \times 10^1 + 1.610_{10} \times 10^{-1}$$

Assume that we can store only four decimal digits of the significand and two decimal digits of the exponent.

Answer:

Step 1: Equalize Exponents

* Adjust the smaller number, the second number is the smaller. Thus, the first step shifts the significand of the smaller number to the right until its corrected exponent matches that of the larger number.

* But we can represent only four decimal digits so, after shifting, the number is

$$1.610 \times 10^{-1} = 0.161 \times 10^0 = 0.016 \times 10^1$$

Step 2: Add significands

The next step is addition of the significands

$$\begin{array}{r} 9.999_{10} \\ (+) 0.016_{10} \\ \hline 10.015_{10} \end{array}$$

The sum is $10.015_{10} \times 10^1$

Step 3: Normalization:

This sum is not in normalized scientific notation, so we need to adjust it.

$$10.015_{10} \times 10^1 = 1.0015_{10} \times 10^2$$

Thus, after the addition we may have to shift the sum to put it into normalized form, adjusting the exponent appropriately.

Step 4: Rounding

The number, $1.0015_{10} \times 10^2$ is rounded to four digits in the significant to

$$1.002 \times 10^2$$

UNIT-III
PART-A

1. Mention the various types of pipelining [NOV/DEC-2017]

* The various types of pipelining are

1. Arithmetic pipeline - The arithmetic logic units of a computer can be segmented for pipeline operations in various data formats.

2. Instruction pipeline - An instruction pipeline increases the performance of a processor by overlapping the processing several different instructions. This is done by dividing the the instruction execution process into several stages like fetch, decode, execute, memory access and write back stages

2. Mention the various phase in executing an instruction. [NOV/DEC-2017]

The various phase in executing an instruction, they are

1. IF: Instruction fetch
2. ID: Instruction Decode and register file read
3. EX: Execution or address calculation
4. MEM: Data Memory access
5. WB: Write Back.

3. Name the control signals required to perform arithmetic operations. [APR/MAY-2017]

* Two ALUop signals

* Seven other signals

- > RegDst - which field for write register
- > RegWrite - write to Register file
- > ALUSrc - Source for second ALU input
- > PC Src - Source for PC ($PC + 4$ or target address)
- > MemRead - Read input address from memory
- > MemWrite - write input address/data to memory
- > MemToReg - source of write register post data input.

* Branch control signal - set when instruction is branch

4. Define hazard. Given an example for data hazard. [APRIL/MAY - 2017]

Hazard:

There are situations in pipelining when the next instruction cannot execute in the following clock cycle. These events are called Hazards.

Types of Hazards

1. Data Hazards
2. Structural Hazards
3. Control Hazards

Data Hazards:

* When a planned instruction cannot execute in the proper clock cycle because data (r_i) needed to execute the instruction is not yet available.

* This type of hazard is called Data Hazard

Eg: $A \leftarrow 3 + A$ Assume $A = 5$
 $B \leftarrow 4 \times A$

5. what is meant by pipeline bubble? [NOV/DEC-2016]

* A bubble or pipeline stall is a delay in execution of an instruction in an instruction pipeline in order to resolve a hazard.

* The values are preserved until the bubble has passed through the execution stage.

6. what is a data path? [NOV/DEC-2016]

* A datapath is a collection of functional units such as ALU, multipliers that perform data processing operations, registers and buses.

1. What is a hazard? what are its types [Nov/Dec-2015]

* In pipelining there are situations when the next instruction cannot execute in the following clock cycle. These events are called hazard

* There are three hazards present in the pipelining such as

1. Structural hazard
2. Data hazard
3. Control hazard

2. What is meant by branch prediction?

* Branch prediction is a method of resolving a branch hazard. It assumes a given outcome for the branch and proceeds from that assumption rather than waiting to ascertain the actual outcome

* A simple form of branch prediction is we have to assume branch is not taken.

* This assumption is only possible in five stage pipeline. For deeper pipelines this assumption not suitable, it increases branch penalty.

9. What is Exception? [Nov/Dec - 2014] [May/June - 2016]

* Exceptions are internally generated unscheduled events that disrupt program execution and they are used to detect overflow.

10. What is a branch prediction buffer? [May/June - 2015]

* Dynamic branch prediction is a prediction of branches using runtime information. To implement dynamic branch prediction method we have to use one buffer is called branch prediction buffer. It is also called branch history table.

* This buffer is small memory that is indexed by the lower portion of the address of the branch instruction.

* It also contains one or more bits indicating whether the branch was recently taken or not.

11. What are R-type Instructions? (May - 2015)

* The arithmetic logic instructions read operands from two registers, perform an ALU operation on the contents of the registers and write the result to a register we call these instructions are R-type instructions.

* This R-type instruction class includes add, sub, AND, OR and SLT

Eg OR \$t1, \$t2, \$t3 reads \$t2 and \$t3, performs logical OR operation and saves the result in \$t1

2. what are the advantages of pipelining? [May/June-2016]

1. The advantage of pipelining is, it saves time in the execution of instruction, because it uses overlapping of instruction

2. To increase the speed and performance of the processor.

3. what is meant by pipelining

* Pipelining is an implementation technique in which multiple instructions are overlapped in execution. Pipelining is most important technique used to increase the speed and performance of the processor.

14. How many stages available in pipelining datapath?

* There are five stages available in pipelining data path.

1. IF - Instruction fetch
2. ID - Instruction Decode and Register file Read
3. Ex - Execution or address calculation
4. MEM - Data Memory Access
5. WB - Write Back.

15. Define Data hazard

* Data hazard occur when a planned instruction cannot execute in the proper clock cycle because data that is needed to execute the instruction is not yet available.

* Data hazard mainly occur in executing arithmetic instructions

16. How to resolve a hazard?

* Pipeline stall also called bubble. A stall initiated in order to resolve a hazard

17 What is an exception? Give one example for MIPS exception. [Nov/Dec-18] [Apr/May-18]

* Exceptions are also called as interrupt. It is an unscheduled event that disrupts program execution. It is also used to detect an overflow condition.

Example:

- i) Arithmetic overflow
- ii) Undefined Instruction
- iii) Invoke the OS from user program.

18 Write the two steps that are common to implement any type of instruction. [Nov/Dec-18]

For every instruction, the first two steps are identical

- 1) Send the Program Counter (PC) to the instruction memory and fetch the instruction from that memory.
- 2) Read one or two registers, using fields of the instruction.

19

What is the ideal CPI of a pipelined processor?
[Apr/May-18]

The ideal CPI of a pipelined processor
almost 1.

Arunai Engineering College

PART-B

Explain how the instruction pipeline works? what are the various situations where an instruction pipeline can stall? Illustrate with an Example [NW1DEC-15] [NW1DEC-16]

Pipelining:

* Pipelining is an implementation technique in which multiple instructions are overlapped in execution.

* Pipelining is most important technique used to increase the speed and performance of the processor.

Example:

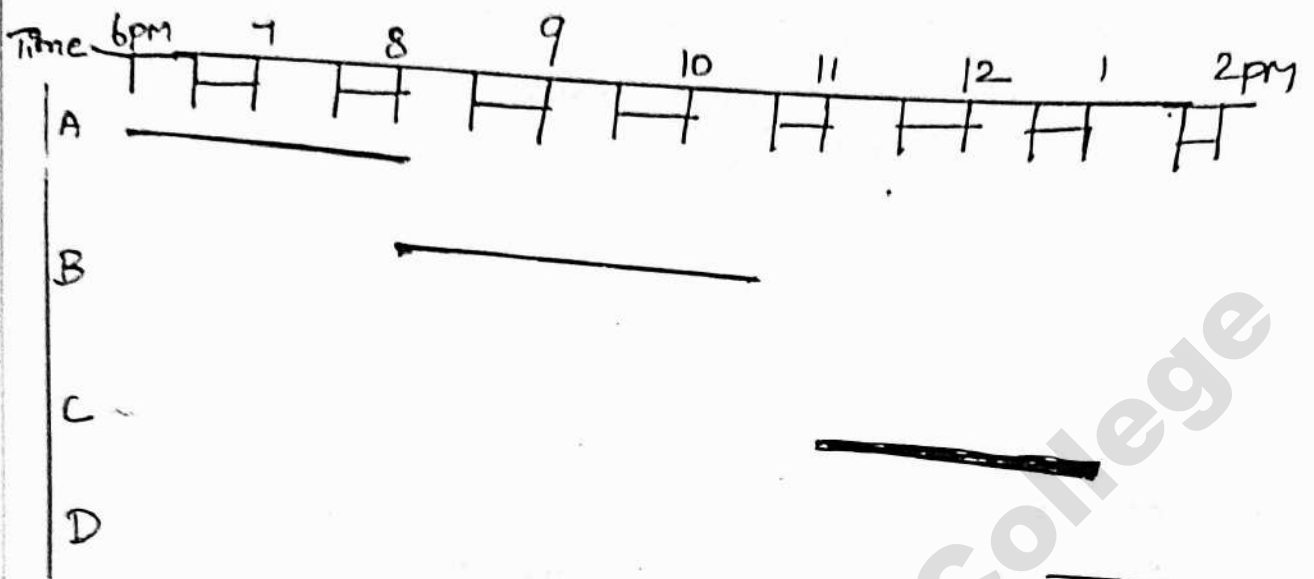
Consider the laundry process through we get the basic idea about technique.

Case 1: Non-pipelined approach to laundry

1. place ^{one} dirty load of clothes in the washer
2. when the washer is finished, place the wet load in the dryer
3. when the dryer is finished, place the dry load on a table and fold
4. when folding is finished, ask your roommate to put the clothes away.

When your roommate is done, start over with

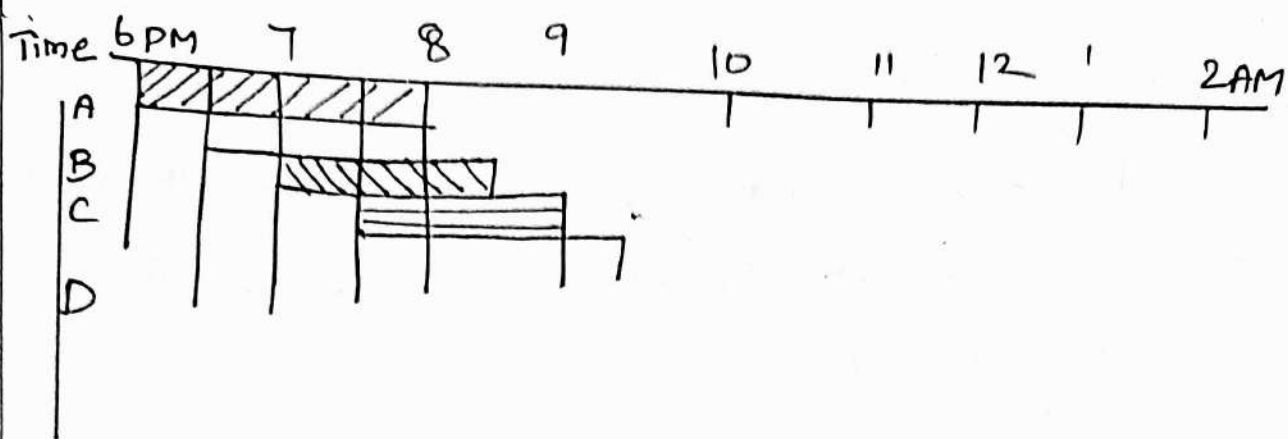
the next dirty load



Task
Order

Case 2: Pipelined approach to laundry:

1. Place dirty load of clothes in the washer.
2. As soon as the washer is finished, place in the dryer.
3. Load the washer with the second dirty load
4. First load is dry then move the wet load to the dryer and put the next dirty load into the washer.
5. Roommate put the first load away, you start folding the second load, the dryer has the third load and you put the fourth load into the washer



Task
Order

*These two processes is explained below

↳ A, B, C and D indicates the name of the Person.

↳ Each has dirty clothes to be washed, dried, folded and put away.

↳ Laundry task has four stages washing, drying, folding and putting away.

↳ Each task takes 30 minutes to complete it.

So non-pipelining laundry process takes 8 hours for ⁴ loads of wash while pipelined laundry takes ¹ just 3.5 hours

↳ Pipelined laundry is potentially four times faster than non-pipelined laundry process.

Pipelining in Processor:

The same principles apply to processors where we pipeline instruction execution. MIPS instructions classically take five steps

1. Fetch instruction from memory
2. Read registers while decoding the instruction.
The regular format of MIPS instructions allows reading and decoding to occur simultaneously
3. Execute the operation or calculate an address
4. Access an operand in data memory.
5. Write the result into a register.

Single cycle VS pipelined Performance:

Let us assume operation times for major functional units.

↳ Memory access 200 ps

↳ ALU operation 200 ps

↳ Register file read or write 100 ps

In single clock cycle model - Each every instruction takes exactly one clock cycles, so it will produce the slow speed to execute the instruction.

↳ The time between the first and fourth Instruction in the non-pipelined design is 3×800 or 2400 PS

↳ The time between the first and fourth Instruction in the pipelined design is 3×200 or 600 PS

$$\text{Time between instruction}_{\text{pipelined}} = \frac{\text{Time between instruction}_{\text{(nonpipelined)}}}{\text{Number of Pipe stages}}$$

* Pipelining improves performance by increasing instruction throughput and decreasing the execution time of an individual instruction.

Designing Instruction Sets for pipelining:

Design of the MIPS instruction set, which was designed for pipelined execution.

Some Important factors such as

1. Length of the Instruction
2. Instruction format
3. Memory operands
4. Operand alignment

1. All MIPS instructions are the same length. This restriction makes it much easier to fetch instructions in the first pipeline stage and to decode them in the second stage.

2. MIPS has only a few instruction formats, with the source register fields being located in the same place in each instruction.

Due to this symmetry the second stage can begin reading the register file at the same time that the hardware is determining what type of instruction was fetched.

3. Memory operands only appear in loads or stores in MIPS. Due to this restriction we can use the execute stage to calculate the memory address and then access memory in the following stage.

4. Since operands are aligned in memory, data can be transferred between processors and memory in a single pipeline stage.

2. Explain the different types of pipeline hazards with suitable Example. [May-2015]

Explain the pipeline hazard in detail [Nov/Dec-2017]

Pipeline Hazards:

* These are situations in pipelining when the next instruction cannot execute in the following clock cycle. These events are called hazards

Types of Hazards:

1. Structural Hazards
2. Data Hazards
3. Control Hazards

1. Structural Hazards

* The first hazard is called a structural hazard.

These hazards are because of conflicts due to insufficient resources when even with all possible combination, it may not be possible to overlap the operation.

* The performance of pipelined processor depends on whether the functional units are pipelined and whether they are multiple execution units to allow all possible combination of instructions in the pipeline.

* If for some combination, pipeline has to be stalled to avoid the resource conflicts then there is a Structural hazard.

* In other words, we can say that when two instructions require the use of a given hardware resource at the same time, the structural hazard occurs.

* The most common case in which this hazard may arise is in access to memory.

* One instruction may need to access memory for storage of the result while another instruction or operand needed is being fetched.

* If instructions and data reside in the same cache unit, only one instruction can proceed and the other instruction is delayed.

* To avoid such type of structural hazards many processors use separate caches for instruction and data.

In our Laundry Example structural hazards will occur in two cases

1. If we use a washer-dryer combination instead of a separate washer and dryer.

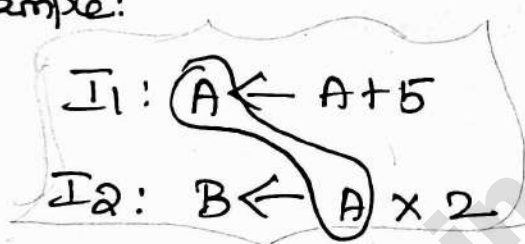
2. If our roommates doing something else without put clothes away.

Q. Data or Data dependant Hazards:

* Data hazards occurs when the pipeline must be stalled because one step must wait for another to complete.

(ie) when a planned instruction cannot execute in the proper clock cycle because data that is needed to execute the instruction is not yet available.

Example:



Suppose $A = 10$ means

* When these two operations are performed in the given order, the result is 30

* If they are performed concurrently the value of A used in computing B would be the original value 10, leading to an incorrect result.

* Data used in the I₂ depends on the result of I₁, This is called data hazard or data dependent hazard.

These are some solutions to rectify data hazards. They are

1. Forwarding or Bypassing
2. Hazard detection unit
3. Reordering code to avoid pipeline stalls
4. Hazard detection by software.

Generally in MIPS instruction set has five stages for pipelining process.

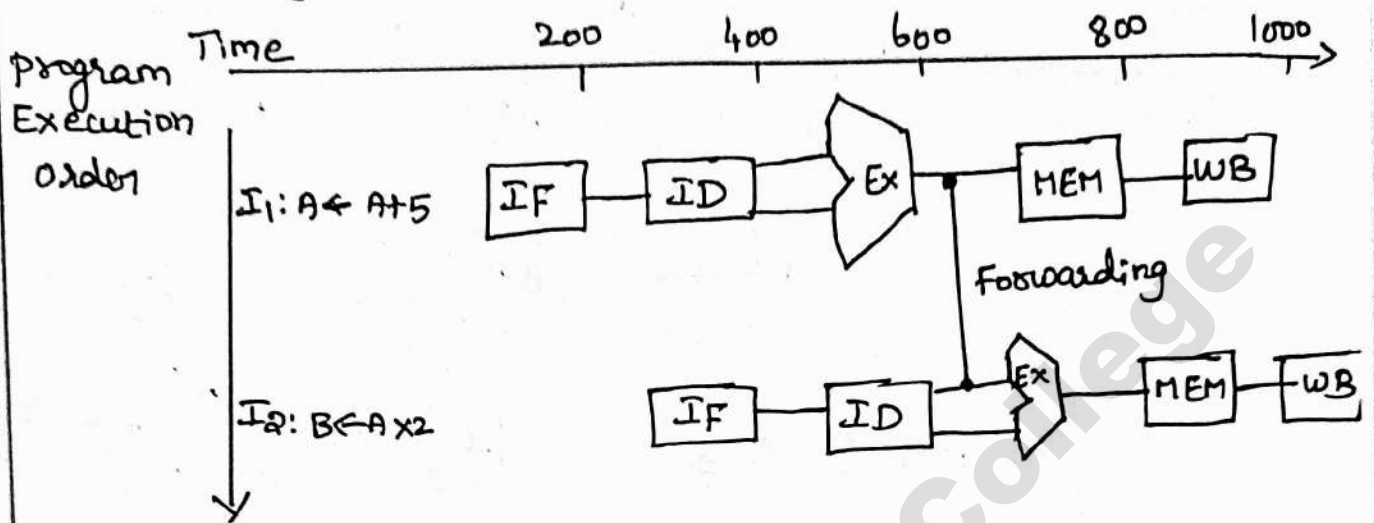
1. IF - Instruction field.
2. ID - Instruction Decode / Register file read stage
3. EXE - Execution stage
4. MEM - Memory access stage
5. WB - write Back stage.

1. Forwarding or Bypassing:

* The primary solution is based on the observation that, we don't need to wait for the instruction to complete before trying to resolve the data hazard.

* Adding extra hardware to retrieve the missing item early from the internal resources is called forwarding or Bypassing.

The graphical representation for forwarding.



* Forwarding paths are valid only if the destination stage is later in time than source stage.

* Forwarding path cannot prevent all pipeline stalls, so to focus new data hazard is called Load use data hazard.

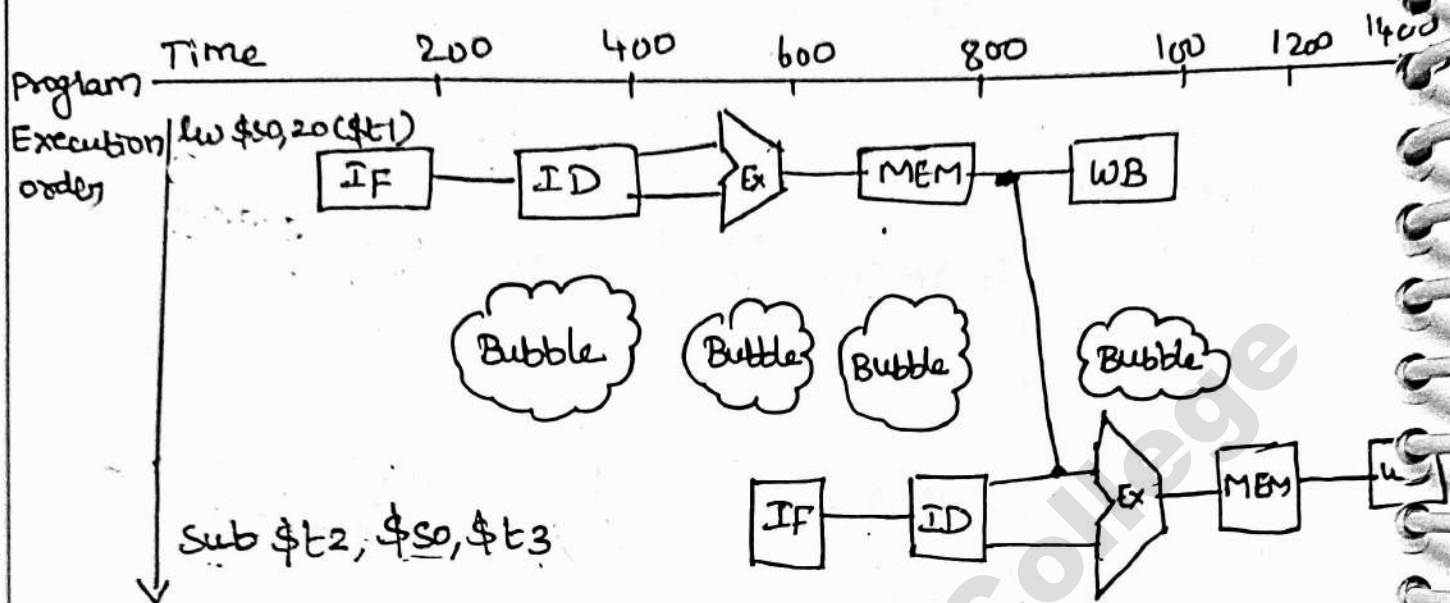
* It is a specific form of data hazards in which the data being loaded by a load instruction has not yet become available when it is needed by another instruction.

Pipeline Stall:

* It is also called bubble.

* Pipeline stall is a stall initiated in order to resolve a hazard.

Stall even with forwarding



Hazard Detection Unit:

* The dependence between load and the following instructions causes data hazard, it cannot be solved by forwarding.

* To solve the problem, in addition to a forwarding unit, need hazard detection unit. It operates during the ID stage so that it can insert the stall between the load and its use.

Checking for load instructions, the control for the hazard detection unit is this single condition

if (ID|EX.Mem Read and
 ((ID|EX.Register Rt = IF|ID.Register Rs) or
 (ID|EX.Register Rt = IF|ID.Register Rt)))

3. Reordering Code to Avoid Pipeline Stall:

Consider the following code segment in C

a = b + e;

c = b + f;

* Here is the generated MIPS code for this segment assuming all variables are in memory and are addressable as offsets from \$t0:

lw \$t1, 0(\$t0)

lw \$t2, 4(\$t0)

add \$t3, \$t1, \$t2

sw \$t3, 12(\$t0)

lw \$t4, 8(\$t0)

add \$t5, \$t1, \$t4

sw \$t5, 16(\$t0)

Solution:

* Both add instructions have a hazard because of their respective dependence on the immediately preceding lw instruction. Moving up the third lw instruction to become the third instruction eliminates hazards.

lw \$t1, 0(\$t0)

lw \$t4, 8(\$t0)

lw \$t2, 4(\$t0)

add \$t3, \$t1, \$t2

sw \$t3, 12(\$t0)

add \$t5, \$t1, \$t4

sw \$t5, 16(\$t0)

* on a pipelined processor with forwarding, the reordered sequence will complete in two fewer cycles, than the original version.

4. Hazard Detection By software:

* An alternative method for handling data hazards is to leave the task of detecting data dependencies to the software.

* The compiler can introduce the necessary delay needed between the two instructions by inserting NOP (NO operation) instructions.

* NOP is a case where pipeline stage executing an instruction but that does not make any changes in the state.

* NOP acts like a bubble in the pipeline stage.

Disadvantage:

* Insertion of NOP leads to larger code size.

3. Control Hazards

* The purpose of the instruction fetch unit is to supply the execution units with a steady stream of instructions.

* This stream is interrupted when pipeline stall occurs either due to cache miss or due to branch instruction. Such situation is known as instruction hazard or control hazard or branch hazard.

Branch Penalty:

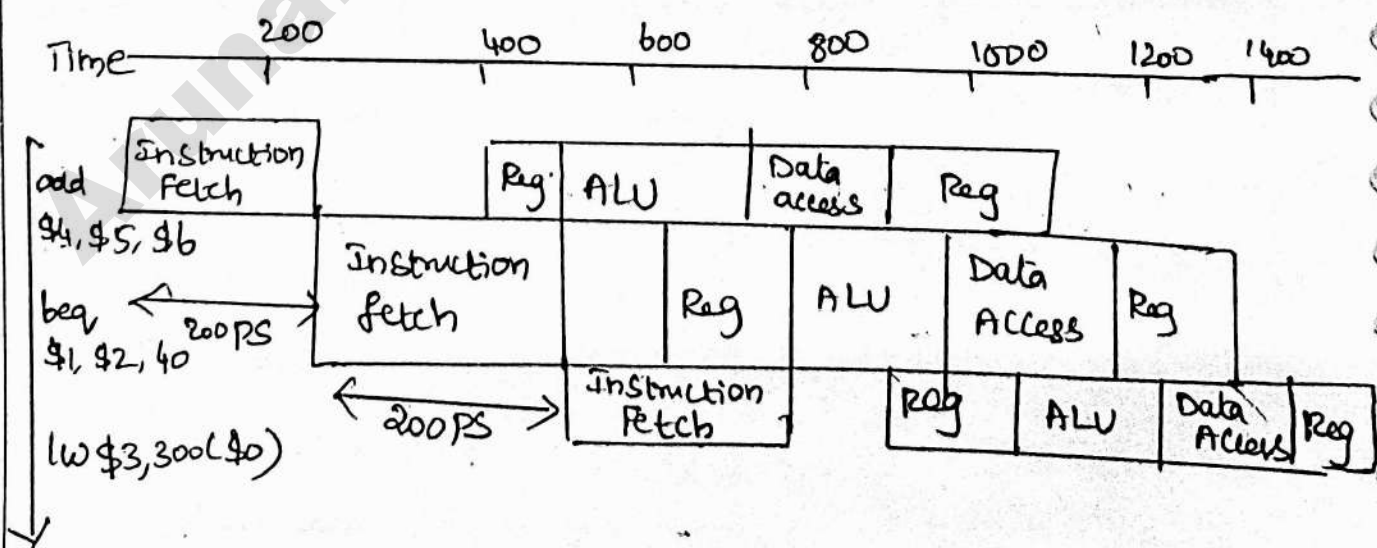
The time lost as a result of a branch instruction is often referred as the branch penalty.

* The pipeline cannot possibly know what the next instruction should be, because it only received the branch instruction from memory.

* To avoid stall, we fetch a branch after that waiting until the pipeline determines the outcome of the branch and knows what instruction address to fetch from.

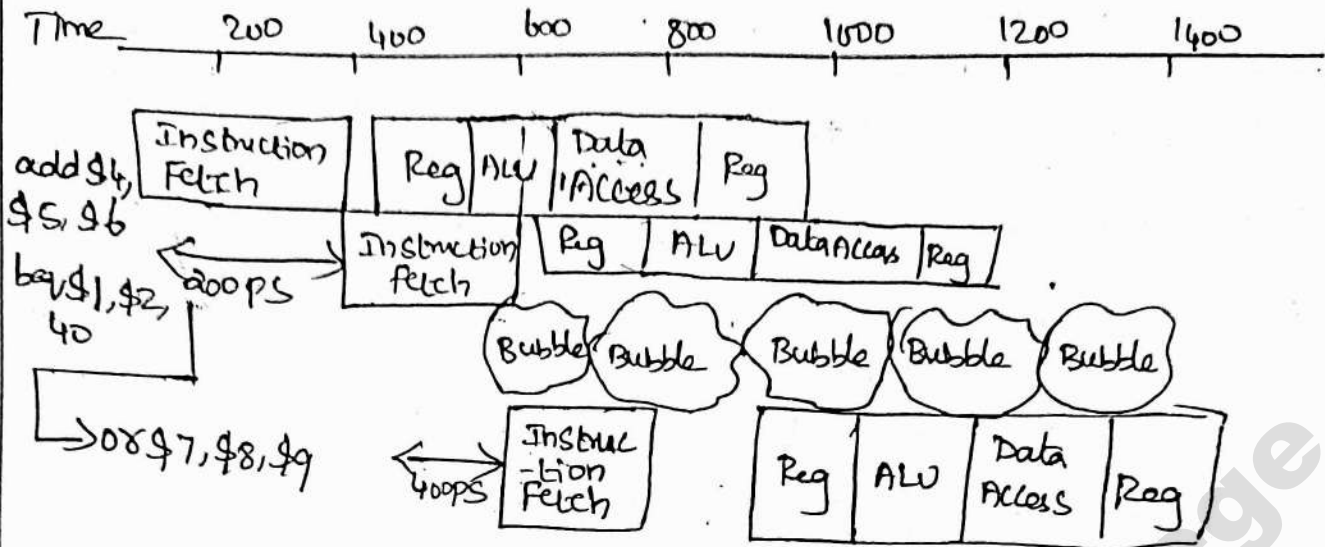
Branch Prediction:

* Branch prediction is a method of resolving a branch hazard that assume a given outcome for the branch and proceeds from that assumption rather than waiting to ascertain the actual outcome
 Eg pipeline when the branch is not taken



Pipeline when the branch is taken

3-11



Dynamic Hardware Predictors:

* Dynamic Hardware predictors make their guesses depending on the behaviours of each branch and may change predictions for a branch over the life of a program.

* one popular approach to dynamic prediction of branches is keeping a history for each branch as taken or not taken

* Using the recent past behavior to predict the future behavior.

Advantages of pipeline:

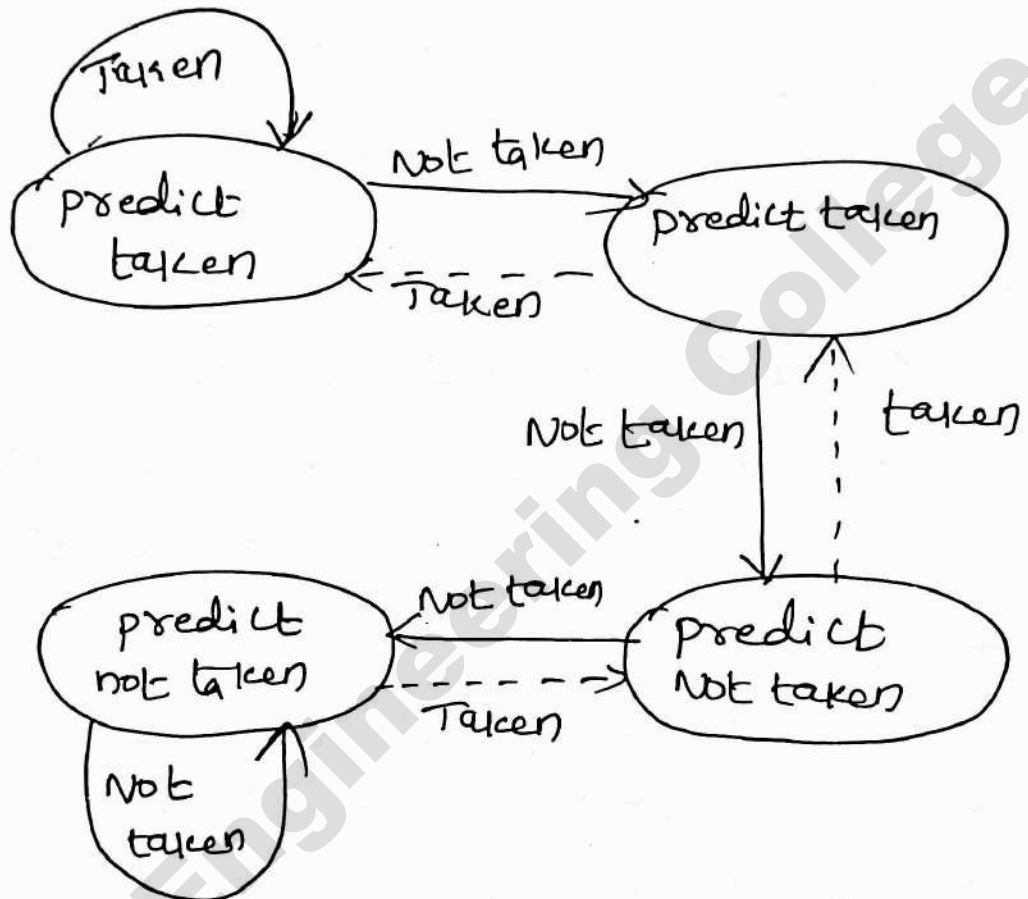
1. It increases the number of simultaneously executing instructions

* It increases the rate at which instructions are started and completed.

* It improves instruction throughput rather than individual instruction execution or latency.

Two-bit prediction scheme:-

To overcome the drawback of branch prediction buffer method, new method called two bit prediction scheme is used.



Finite state machine for 2 bit prediction scheme

Explain in detail how exceptions are handled in MIPS architecture? [APR/MAY-15]

* Control is the most challenging aspect of processor design because of two reasons:

1. It is the hardest part to get right

2. It is the hardest part to make fast

* One of the hardest parts of control is implementing exceptions and interrupts.

* Handling exceptions and interrupts are more complex task than handling branches or jumps

* Exception and interrupt are initially created to handle unexpected events from within the processor.

Exception:

* Exception is an unscheduled event that disrupts program execution and used to detect overflow. It is also called Interrupt. Exception refers to any unexpected change in control flow without distinguishing whether the cause is internal or external

Interrupt: Interrupt is an exception that comes from outside of the processor.

* Interrupt refer to any unexpected change occurred only when the event is externally caused.

Type of Event	From where?	MIPS Terminology
I/O device request	External	Interrupt
Arithmetic overflow	Internal	Exception
Invoke the operating system from user program	Internal	Exception
Using an undefined instruction	Internal	Exception
Hardware malfunctions	Either	Exception or Interrupt

To detect two kinds of exceptions we need to implement control.

Two kinds of Exceptions arise

1. From the positions of the instruction set
2. Implementation

* First we have to detect the condition for exception.

* Once the exception is detected next we need to take appropriate action for it.

* Detecting exceptional conditions and taking the appropriate action is on the critical timing path of a processor.

* Processor must determine the clock cycle time and thus performance. To design a control unit we must have proper attention to exceptions.

* Because exception can reduce performance and it leads complex task to getting the correct design.

Exceptions in MIPS Architecture.

In MIPS architecture two kinds of exceptions occur at

1. Execution of an undefined instruction
2. An arithmetic overflow

* Once exception occur then the processor must save the address of the offending instruction in the exception program counter (EPC) and then transfers control to the operating system at some specified address.

* After transferring control to the operating system it must take appropriate action to provide some service to the user program.

* Operating system will provide the following services to user program

1. Taking some predefined action in response to an overflow.

2. Stopping the execution of the program and reporting an error.

* After performing necessary action is required because of exception, the operating system can terminate the program or may continue its instruction.

* Operating system use the exception program Counter (EPC) to determine where to restart the execution of the program.

* To handle the exceptions by operating system it must know the reason for the exception.

* Operating system must know the reason for exception because the only it can handle properly.

* There are two main methods used to communicate the reason for an exception.

↳ MIPS architecture has the following two methods to find the reason for exception.

1. Status register
2. Vectored Interrupts.

Status register:

* It is also called the cause register. It holds a field that indicates the reason for the Exception.

Vectored Interrupts:

In a vectored interrupt, the address to which control is transferred is determined by the cause of the exception.

Exception type	Exception vector address (in hex)
Undefined Instruction	8000 0000 hex
Arithmetic overflow	8000 0180 hex

* operating system can know the reason for exception by the address at which it is initiated.

* The addresses are separated by 32 bytes or eight instructions.

* operating system must record the reason for the exception and may perform some limited processing in this sequence.

* when the exception is not vectored then a single entry point for all exception can be used.

* If single entry point is used for all exceptions then the operating system must decode the status register to find the reason for exception.

* To handle exceptions we can add a few extra registers and control signals to their basic implementation.

Implementation of exception in MIPS architecture:

* To implement exception system in MIPS architecture assume it has single entry point for all exception and has address value is 8000 0180

* This address value indicates it is an arithmetic overflow exception.

* To handle this exception we need to add two additional registers to our current MIPS implementation.

* Two addition registers are

i) EPC register:

A 32 bit register used to hold the address of the affected instruction. This register is needed even when exceptions are vectored

ii) Cause register:

* This register is used to record the cause of the exception. In MIPS architecture this register is 32 bits long and some bits are currently unused.

* 10 bits are used to represent an undecoded instruction and 12 bits used to represent arithmetic overflow

5. Explain the basic MIPS Implementation with necessary multiplexers and control lines. [NOV/DEC-2015]

Basic MIPS Implementation:

* In examining the implementation we can get more information to increase the processor performance. Through the implementation the following things can be identified by the user.

1. How the instruction set architecture determines?
2. How various implementation strategies affect the clock rate?
3. How implementation affects the CPI cycles per instruction of the computer?

* MIPS has three kinds of core instructions such as

1. The memory-reference instructions load word (lw) and store word (sw)
2. The arithmetic-logic instructions add, sub, AND, OR and slt
3. Branch instructions \rightarrow jump (j) and branch equal (beq)

To implement the above three types we have same method but independent of the exact class of instruction.

* For every instruction, the first two steps are identical.

1. Send the program counter (PC) to memory that contains the code and fetch the instruction from that memory.

2. Read one or two registers, using fields of the instruction to select the registers to read. For the load word instruction, we need to read only one register, but most other instructions require that we read two registers.

* These two steps are common for all the instruction set.

* After these two steps, the action required to complete the instruction depend on the instruction class.

↳ MIPS instruction set has simplicity and regularity and it will simplify the implementation by making the execution of many of the instruction classes similar.

Example:

* All instruction classes, except jump, use the arithmetic logical unit (ALU) after reading the registers.

* The memory reference instruction use the

↳ ALU for an address calculation

↳ Arithmetic logical instructions for the operation execution

↳ Branches for comparison.

* After using the ALU, the actions required to complete the task differs from various instruction classes.

* A memory reference instruction will need to access the memory either to read data for a load or write data for a store.

* An arithmetic logical or load instruction must write the data from the ALU or memory back into a register.

* Branch instruction need to change the next instruction address based on the comparison, otherwise the PC should be incremented by 4 to get the address of the next instruction.

* All instructions start by using the program counter to supply the instruction address to the instruction memory.

* After the instruction is fetched, the register operands used by an instruction are specified by fields of the instruction.

* Once the register operands have been fetched, they can be operated to do the following tasks:

1. To compute a memory address
2. To compute an arithmetic result
3. To compare for a branch.

* If the instruction is arithmetic logical instruction then the result from the ALU must be written to a register.

* If the operation is a load or store, the ALU result is used as an address to either store a value from the registers or load a value from memory into the registers.

* The result from the ALU or memory is written back into the register file.

* Branch instruction require the use of the ALU output to determine the next instruction address, which comes either from the ALU or from an adder that increments the current PC by 4.

* It Shows most of the flow of data through the processor but it omits two important aspects of instruction execution.

1. It Shows data going to a particular unit as coming from two different sources.

2. Several of the units must be controlled depending on the type of instruction

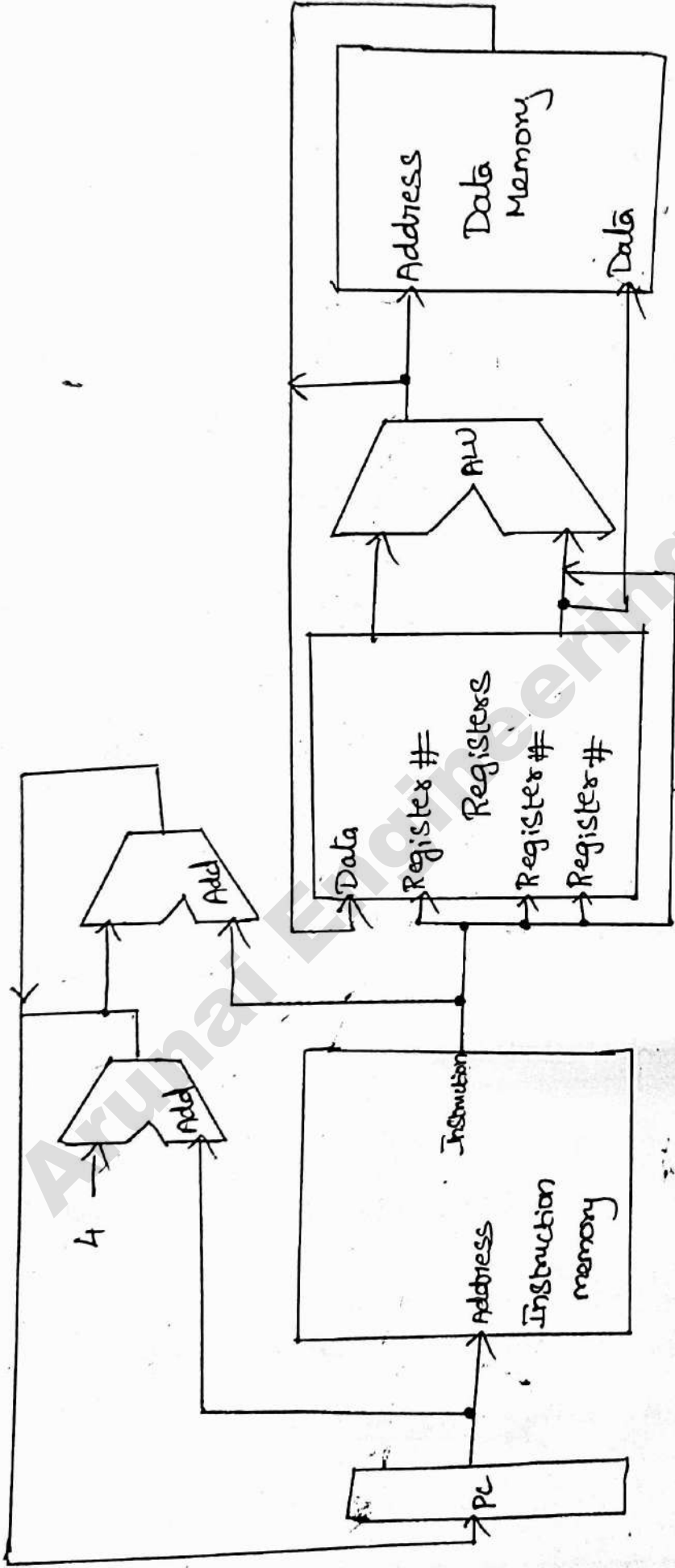
Problems in first aspect:

* Data going to a particular unit as coming from two different sources.

Example:

The value written into the PC can come from one of two address.

The data written into the register file can come from either the ALU or the data memory.



An abstract view of the Implementation of the MIPS Subset showing the major functional units and the major connection between them

Thicklines - Interconnecting the functional units represent buses, which consist of multiple signals
 Arrows → To guide the reader in knowing how information flows

* The second input to the ALU can come from a register or the immediate field of the instruction.

* These data lines cannot simply be wired together for practical purpose. so we must add a logic element that chooses from among the multiple sources and move one of those sources to its destination.

* This selection is commonly done with a device called a Multiplexer and it is also called data Selector.

* Multiplexer is a combinational circuit which selects from among several inputs based on the setting of its control lines.

* The control lines are set based on information taken from the instruction being executed.~

Problem in second aspect:

* In that figure several of the units must be controlled depending on the type of instruction.

* The data memory must read on a load and write on a store.

* The register file must be written only on a load or an arithmetic logical instruction.

* ALU must perform several operations. To overcome these problems we can use another circuit with some modification made in previous circuit

* The basic implementation of MIPS required three multiplexers and one major unit as control lines

Control unit:

* Control unit has the instruction as an input and it is used to determine how to set the control lines for the functional units and two of the multiplexers

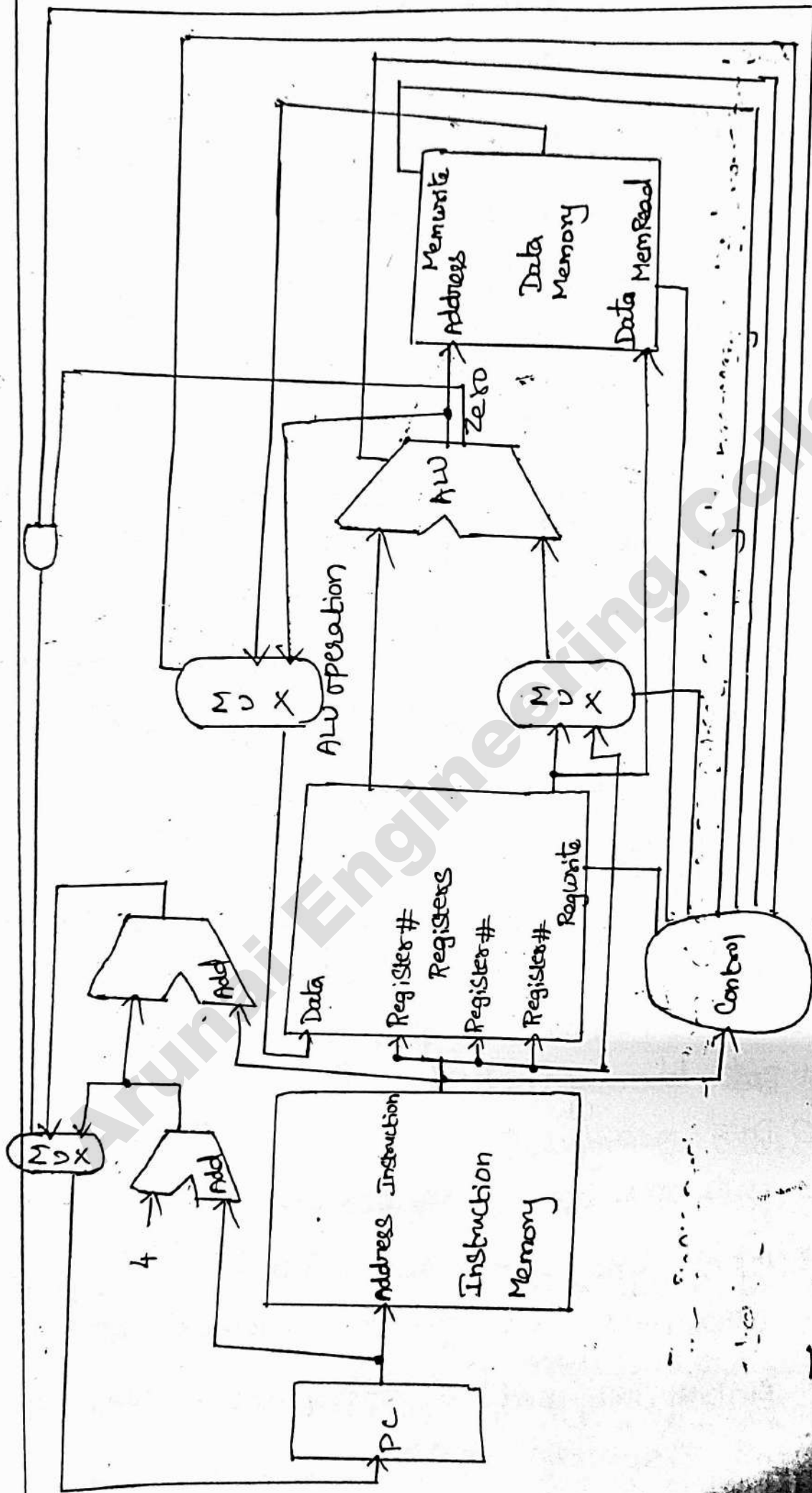
Function of third Multiplexer:

* Third multiplexer determines whether $PC + 4$ or branch destination address is written into the PC.

* It is set based on the zero output of the ALU and it is used to perform the comparison of a beq instruction

* The regularity and simplicity of the MIPS instruction set a simple decoding process.

* Simple decoding process used to determine how to set the control line.



The Basic Implementation of the MIPS Subset, Including the necessary Multiplexers and Control Lines

Functions of the Multiplexers:

* The top multiplexer controls what the value replaces the PC [$PC+4$ or branch destination address].

* The multiplexer is controlled by the gate that "AND" together with the zero output of the ALU and control signal.

* control signal indicates the branch instruction. Middle multiplexer is used to steer the output of the ALU or the output of the data memory for writing into a register file.

* Bottom multiplexer is used to determine the second ALU input is from the registers or from the offset field of the instruction.

Function of control line:

* Control lines are straight forward and determine the operation performed at the ALU

ALU can perform the following operations

1. Data memory read
2. Data memory write
3. write operation on registers

* control line determines whether the ALU perform which operations among three mentioned operations

* Control unit used to control actions taken for different instruction classes.

6. Explain in detail the operation of the datapath.
[Nov/Dec-17]

operation of the datapath:

* To consider three kinds of instruction classes, so far such as

1. R-type instructions
2. Load and store instructions
3. Branch instruction.

* It is important to know the flow of control through the datapath for these three different instruction classes

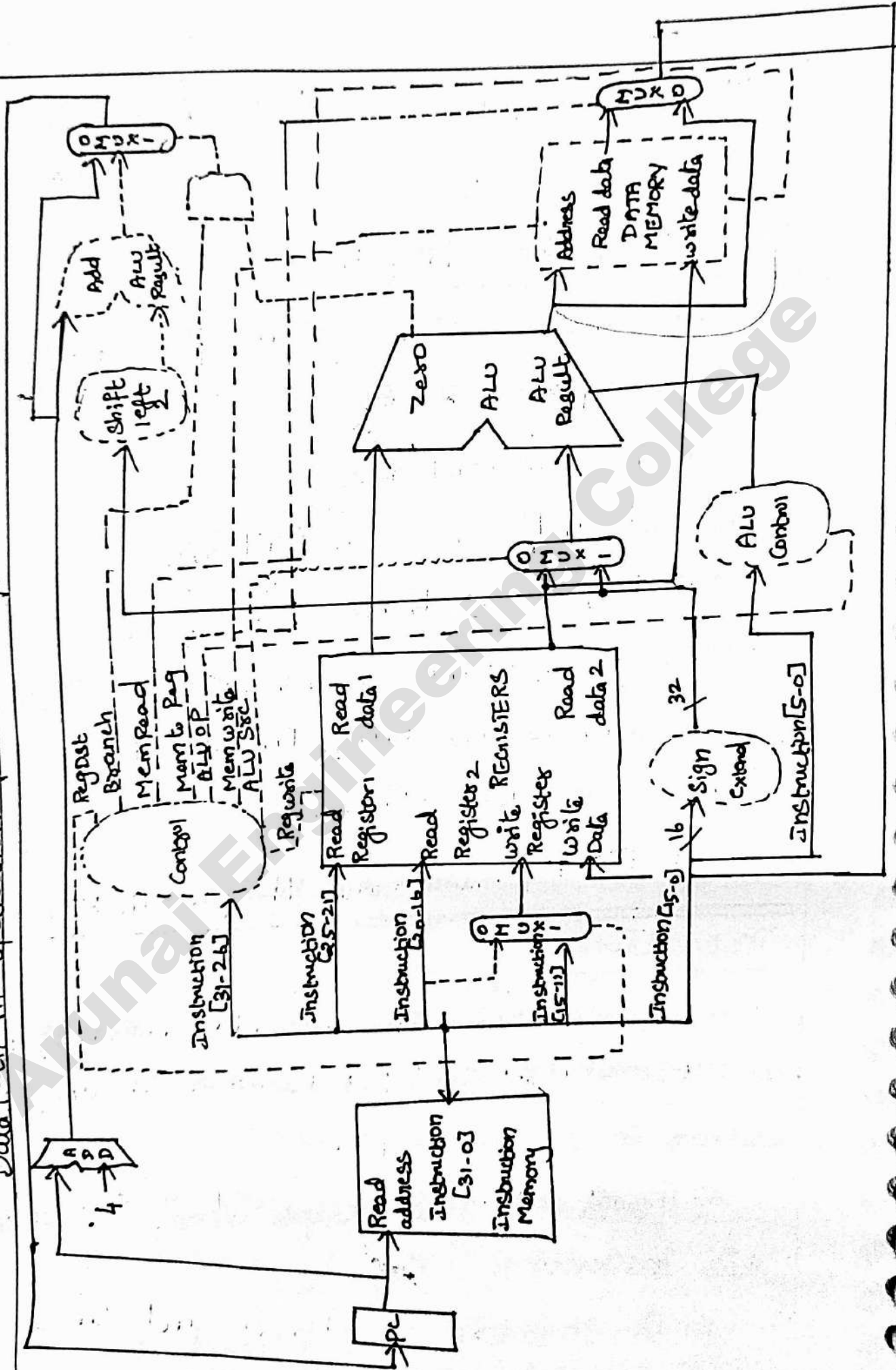
1. operation of the datapath for an R-type instruction:

* Instructions such as `add $t1, $t2, $t3` and remaining four operations (`sub`, `AND`, `OR`, `sll`) occurs in one clock cycle.

* These are four steps used to execute this instruction. They are:

1. The instruction is fetched and the PC is incremented.

Data Path in operation for an R-type Instruction



2. Two registers, $\$t2$ and $\$t3$, are read from the register file.

3. The ALU operates on the data read from the register file, using the function code to generate the ALU function.

4. The result from the ALU is written into the register file in the destination register ($\$t1$).

2. operation of the datapath for a load word instruction.

* Instructions such as $lw \$t1, offset(\$t2)$ comes under this category.

* Five steps involved in execution of load instruction.

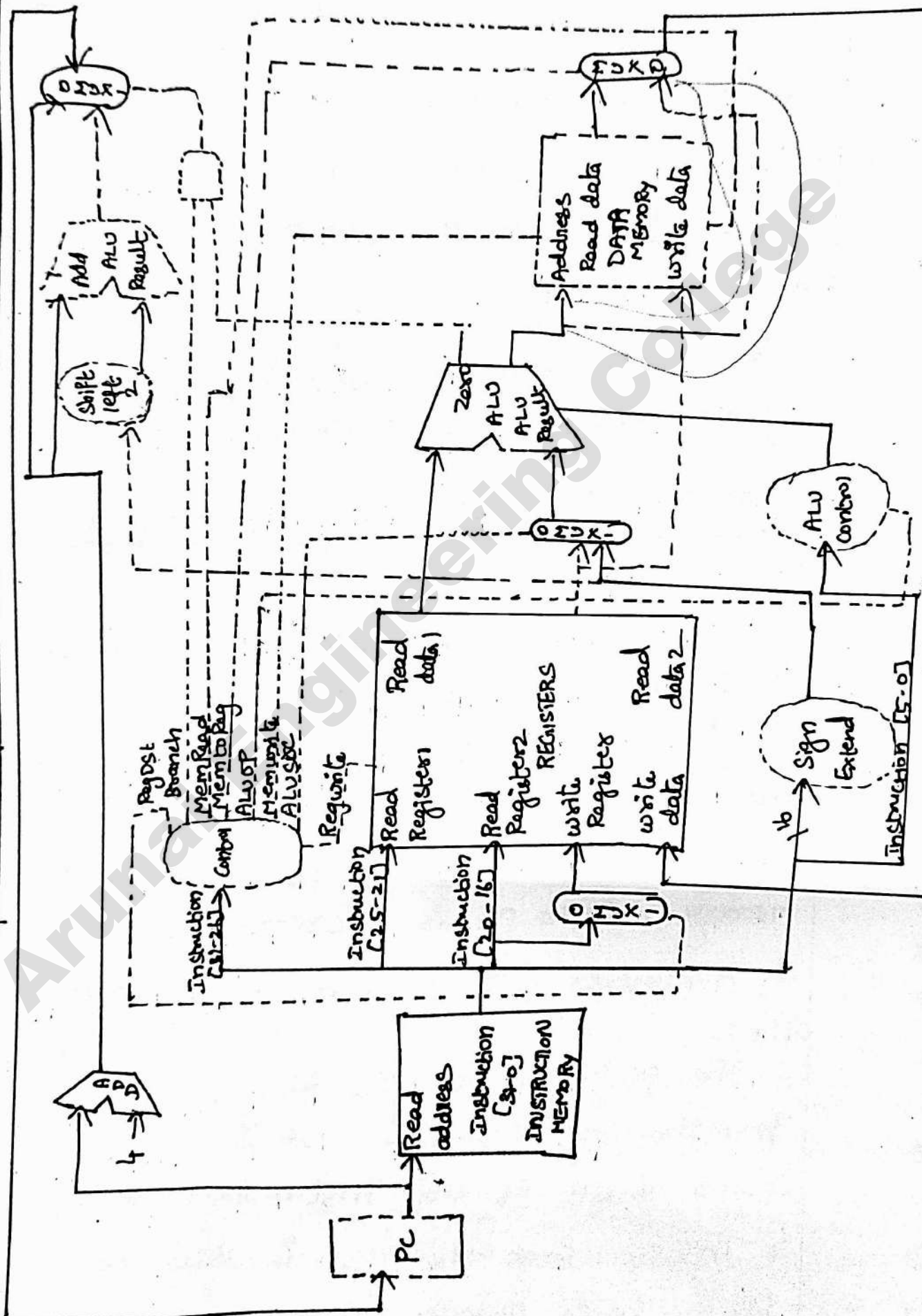
1. An instruction is fetched from the instruction memory and the PC is incremented.

2. A register ($\$t2$) value is read from the register file.

3. The ALU computes the sum of the value read from the register file and the sign-extended, lower 16 bits of the instruction (offset).

4. The sum from the ALU is used as the address for the data memory.

Data path operation for load instruction



5. The data from the memory unit is written into the register file in the destination register (\$t1)

3. operation of the datapath for branch-on-equal instruction.

* Instructions such as `beq $t1, $t2, offset` comes under this category. Four steps are needed to execute the instruction. They are

1. An instruction is fetched from the instruction memory and the PC is incremented.

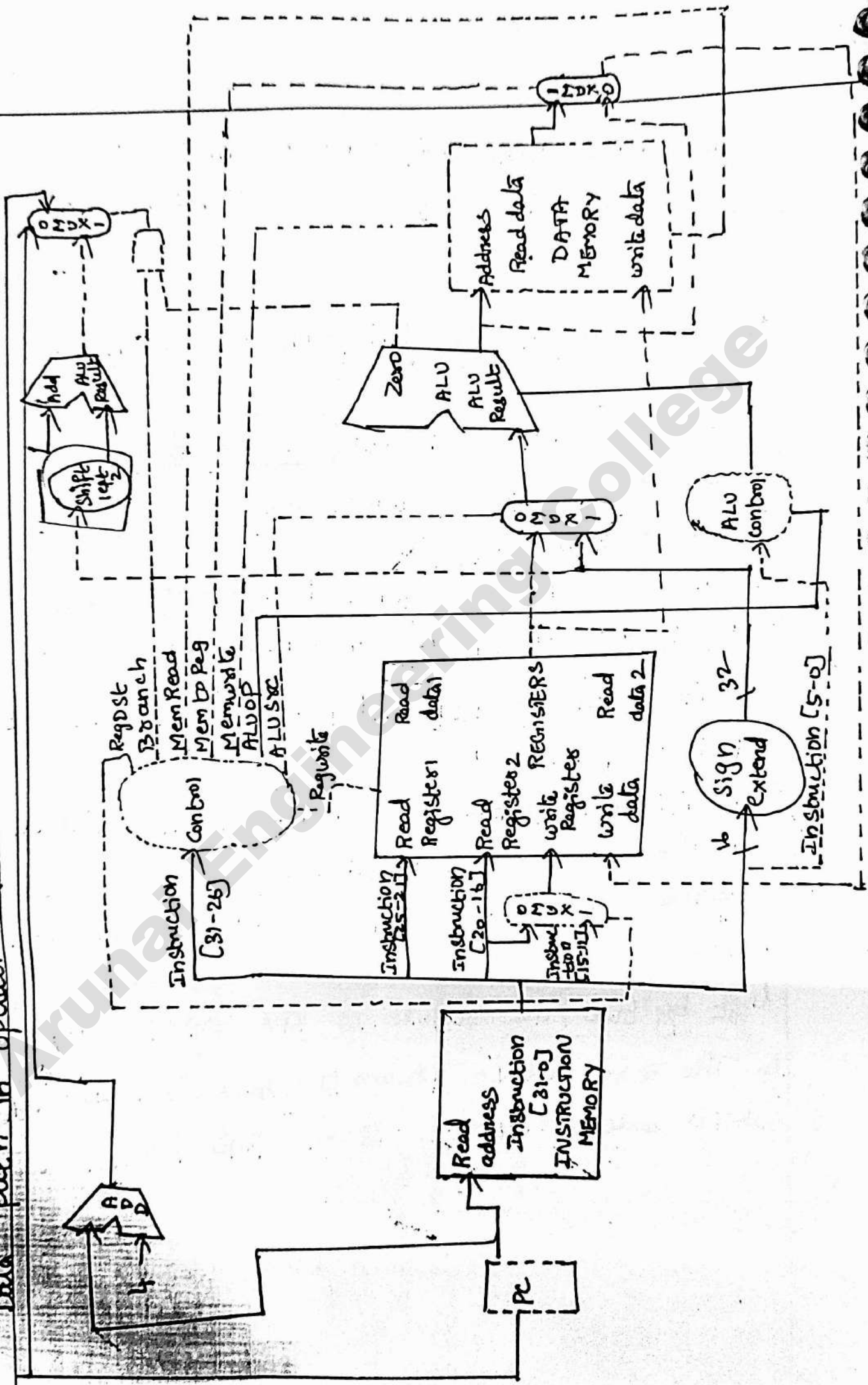
2. Two registers, \$t1 and \$t2 are read from the register file.

3. The ALU performs a subtract on the data values read from the register file.

* The value of PC+4 is added to the sign-extended, lower 16 bits of the instruction (offset) shifted left by two; the result is the branch target address.

4. The zero result from the ALU is used to decide which address result to store into the PC.

Data path in operation for branch on equal instruction



7) The following sequence of instructions are executed in the basic 5-stage pipelined processor:

[Apr/May-18]

O₀ r₁, r₂, r₃

O₀ r₂, r₁, r₄

O₀ r₁, r₁, r₂

- Indicate dependences and their type.
- Assume there is no forwarding in this pipelined processor. Indicate hazards and add NOP instructions to eliminate them.
- Assume there is full forwarding. Indicate hazards and add NOP instructions to eliminate them.

Answer:

a) Indicate dependences and their type

Instructions:

I₁: O₀ r₁, r₂, r₃

I₂: O₀ r₂, r₁, r₄

I₃: O₀ r₁, r₁, r₂

i) True Data dependency:

* Instruction I₂ and I₁ have true data dependency because I₂ can be executed only when I₁ is completed (i.e.) I₂ needs r₁ which is available only when I₁ is completed.

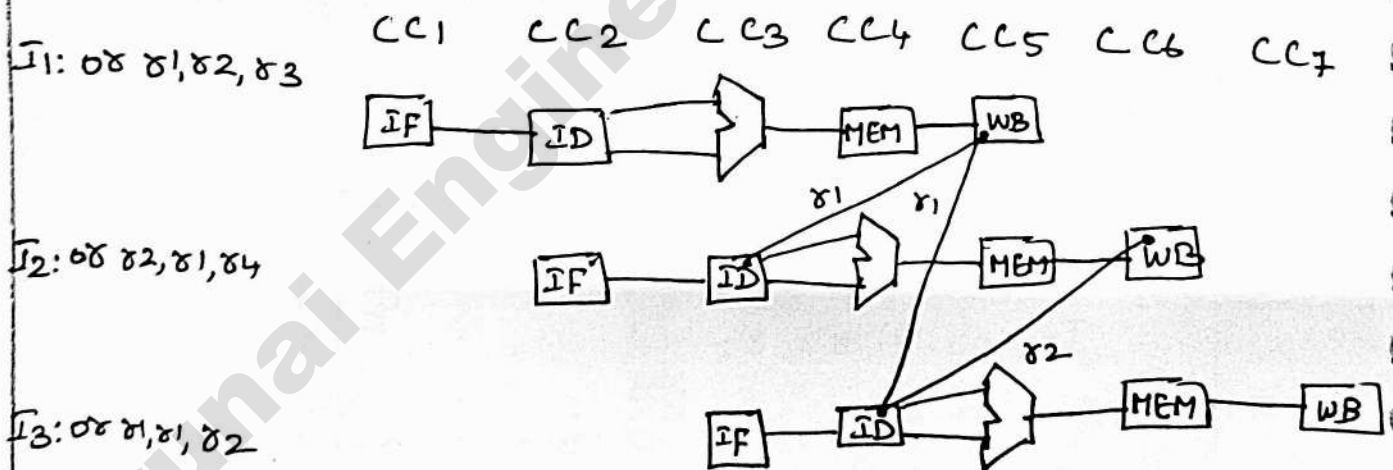
* Similar kind of data dependency exists between I_1 and I_3 , I_2 and I_3 .

2. output dependency:

* Instruction I_1 and I_3 are having output dependency, because I_3 cannot be executed before I_1 .

* No forwarding mean data hazard occurs.

b) Assume there is no forwarding in this pipelined processor. Indicate hazards and add Nop instructions to eliminate them.



* As in above figure, r_1 is available only at CC5, but instructions I_3 need r_1 at CC3 and CC4. This is a situation of data hazard.

* Nop instruction can be inserted to avoid data hazard.

* In instruction I_2 , after instruction fetch (IF) ID stage is delayed until CC5, so that r_1 will be available at CC5.

* Similar case exist between I_2 and I_3

I_1 : $OR\ r_1, r_2, r_3$

I_2 : NOP

I_3 : NOP

I_4 : $OR\ r_2, r_1, r_4$

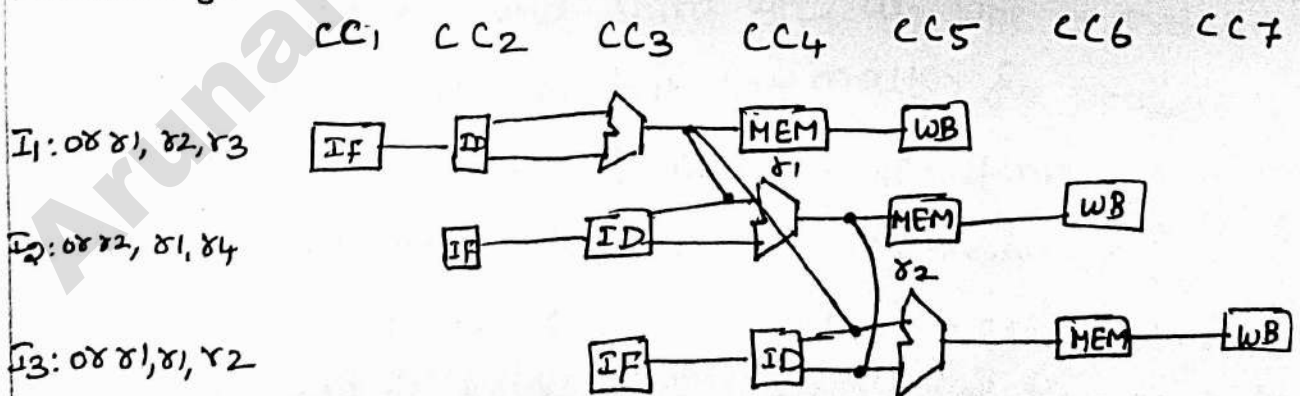
I_5 : NOP

I_6 : NOP

I_7 : $OR\ r_1, r_2, r_2$

c) Assume there is full forwarding, indicate hazards and add NOP instructions to eliminate them.

Forwarding:-



* when forwarding is applied, there is no need of NOP instruction.

* R_1 is forwarded to instructions I_2 and I_3 at CC_3 .

* Similarly, r_2 is forwarded to I_3 at CC_4 .

Arunai Engineering College

8) The following sequence of instructions are executed in the basic 5-stage pipelined processor. [NOV/DEC-18]

lw \$1, 40(\$6)

add \$6, \$2, \$2

sw \$6, 50(\$1) Indicate dependencies and their type. Assuming there is no forwarding in this pipelined processor, indicate hazards and add nop instructions to eliminate them.

Answer:

* Dependencies involve writing registers. There are three possibilities:

i) RAW - read after write

ii) WAR - write after Read

iii) WAW - write after write

The second and third type are not a problem for simple pipelines, but can matter for processors which permit out-of-order execution

Instruction sequence

I₁: lw \$1, 40(\$6)

I₂: add \$6, \$2, \$2

I₃: sw \$6, 50(\$1)

Dependencies

RAW on \$1 from I₁ to I₃

RAW on \$6 from I₂ to I₃

WAR on \$6 from I₁ to I₂ and I₃

* In the basic five-stage pipeline WAR and WAW dependences do not cause any hazards.

* Without forwarding, any RAW dependence between an instruction and the next two instructions (if the register read happens in the second half of the clock cycle and the register write happens in the first half)

* The code that eliminates those hazards by inserting nop instruction is

Instruction sequence

I₁: lw \$1, 40(\$6)

I₂: add \$6, \$2, \$2
nop

I₃: sw \$6, 50(\$1)

Delay I₃ to avoid RAW hazard on \$1 from I₁

A processor has five individual stages, namely IF, ID, EX, MEM and WB and their latencies are 250ps, 350ps, 150ps, 300ps and 200ps respectively. The frequency of the instructions executed by the processor are as follows; ALU: 40%, Branch 25%, Load: 20% and Store: 15%. What is the clock cycle time in a pipelined and non-pipelined processor? If you can split one stage of the pipelined datapath into two new stages, each with half the latency of the original stage, which stage would you split and what is the new clock cycle time of the processor? Assuming there are no stalls or hazards, what is the utilization of the data memory? Assuming there are no stalls or hazards, what is the utilization of the write-register ports of the "Registers" unit?

i) For a pipelined processor, the clock cycle is the time of the pipeline element with the largest latency $T_{clk} = 350\text{ps}$

For a non-pipelined processor, the clock cycle is the sum of the latencies of all the pipeline elements $T_{clk} = 250 + 350 + 150 + 300 + 200$
 $T_{clk} = 1250\text{ps}$

ii) Split the longest stage: ID

The next longest stage is: MEM at 300ps.

$$T_{clk} = 300\text{ps}$$

UNIT - IV

PART-A

1) What is Flynn's Classification? (Nov/Dec-2014)

* Flynn's Classification divides parallel hardware into four groups based on the number of instruction streams and the number of data streams

1. Single Instruction Stream Single Data Stream (SISD)
2. Single Instruction Stream Multiple Data Stream (SIMD)
3. Multiple Instruction Stream Single Data Stream (MISD)
4. Multiple Instruction Stream Multiple Data Stream (MIMD)

2. Brief about Multithreading (Nov/Dec-2014)

* The term multithreading implies that there are multiple threads of control in each processor.

* Multithreading offers an effective mechanism for hiding long latency in building large scale multiprocessors.

3. What is ILP? [Nov/Dec-2015]

Instruction level parallelism is a kind of parallelism it executes the instruction in parallel way

4. Define a super scalar processor [Nov/Dec-2015]

Super scalar processor is a dynamic multiple issue processor. It is an advanced technique that enables the processor to execute more than one instruction per clock cycle by selecting them during execution

5. What is speculation?

An approach that allows the compiler or the process to "guess" the outcome of an instruction to remove its as a dependence in executing other instruction is called speculation

6. What is the need for speculation [Nov/Dec-2015]

It is the most important methods for finding and exploiting more ILP

Differentiate between Strong Scaling and weak Scaling. [APR/MAY-2015]

Strong Scaling	Weak Scaling
<p>1. In this method speed up is achieved on a multiprocessor without increasing the size of the problem</p> <p>* It means measuring speed up while keeping the problem size fixed.</p>	<p>1. In this method speed up is achieved on a multiprocessor while increasing the size of the problem, proportionally to the increase in the number of processors.</p>

8. Compare UMA and NUMA Multiprocessors [May/June-2015]

Uniform Memory Access	Non-Uniform Memory Access
1. Programming challenges are easy	1. Programming challenges are hard
2. UMA machines can scale small size	2. NUMA machines can scale to larger sizes
3. It has higher latency	3. It has lower latency to nearby memory.

9. Define Synchronization

Synchronization is the process of coordinating the behavior of two or more processes which may be running on different processors

10. What is fine grained Multithreading? (May/June 2017)

Fine grained Multithreading is a version of hardware multithreading that implies switching between threads after every instruction.

11. What is Coarse grained Multithreading?

Coarse grained Multithreading is a version of hardware multithreading that implies switching between thread only after significant events such as last level cache miss

12. What is meant by thread?

Thread is a light-weight process which includes the program counter, the registers state and stack. It shares a single address space.

13. What is Multithreading? [NOV/DEC-2016]
[NOV/DEC-2014]

* Multithreading Allows multiple threads to share the functional units of a single processor in an overlapping fashion in order to utilize the hardware resources effectively.

14. What is Instruction level parallelism?
[APR/MAY-2017]
[NOV/DEC-2016]
[MAY/JUNE-2016]

* Instruction level parallelism (ILP) is a measure of number of instructions that can be performed simultaneously during a single clock cycle.

* The potential overlap among instructions is called as Instruction level parallelism.

* There are two primary methods for increasing the potential amount of instruction level parallelism

1. Increasing the depth of the pipeline to overlap more instructions
2. Multiple issue.

15 Define strong scaling and weak scaling.

[NOV/DEC-2017]

Strong scaling:

* In these methods speed up achieved on a multiprocessor without increasing the size of the problem.

" Strong scaling means measuring speed up while keeping the problem size, fixed".

Weak scaling:

* In this method speed up is achieved on a multiprocessor while increasing the size of the problem proportionally to the increase in the number of processors.

16 Difference between Fine-grained multithreading and coarse-grained multithreading. [NOV/DEC-2017]

Fine-grained Multithreading

1. Fine-grained Multithreading is a version of hardware multithreading that implies switching between threads after every instruction.

Coarse-grained Multithreading

Coarse-grained multithreading is a version of hardware multithreading that implies switching between thread only after significant events such as last level Cache miss.

17

Distinguish Implicit multithreading and Explicit Multithreading. [APR/MAY-2017]

Implicit Multithreading	Explicit Multithreading
<p>1. Implicit Multithreading is concurrent Execution of multiple threads extracted from single sequential program.</p> <p>2. Implicit threads defined statically by compiler or dynamically by hardware</p>	<p>1. Explicit Multithreading is concurrent Execution of instructions from different explicit threads.</p> <p>* Either by interleaving instructions from different threads on shared pipelines or</p> <p>* by parallel execution on parallel pipelines</p>

Arunai Engineering

18

Web server is to be enhanced with a new CPU which is 10 times faster on computation than old CPU. The original CPU spent 40% of its time processing and 60% of its time waiting for I/O. What will be the overall speed up? [Nov/Dec-18]

$$\text{Speedup} = \frac{\text{Execution time for task without enhancement}}{\text{Execution time for task with enhancement}}$$

$$\text{Speedup} = \frac{\text{Execution time old}}{\text{Execution time new}} = \frac{1}{(1 - \text{fraction}_{\text{enhanced}}) + \frac{\text{fraction}_{\text{enhanced}}}{\text{Speedup}_{\text{enhanced}}}}$$

$$\text{fraction}_{\text{enhanced}} = 40\% = 0.4$$

$$\text{Speedup}_{\text{enhanced}} = 10$$

Using formula (1)

$$\frac{1}{(1 - 0.4) + \frac{0.4}{10}} = \frac{1}{0.64} = 1.56$$

19

Classify shared memory multiprocessors based on the memory access latency. [NOV/DEC-18]

* Shared memory multiprocessors based on the memory access latency

i) Uniform Memory Access (UMA)

* The same time to access main memory no matter which processor requests it and no matter which word is requested, such machines are called UMA multiprocessors.

ii) non-uniform memory access (NUMA)

* Some memory access are much faster than others, depending on which processor asks for which word. such machines are called NUMA multiprocessors.

20 Give example for each class in Flynn's Classification.
[Apr/May -18]

i) Single Instruction Single Data streams (SISD)

Eg Intel Pentium 4

ii) Single Instruction Multiple Data streams (SIMD)

Eg SSE of Instructions of X86

iii) Multiple Instruction Single Data streams (MISD)

NO Examples

iv) Multiple Instruction Multiple Data streams (MIMD)

Eg: Intel Core I7

21 Protein string matching code has 4 days execution time on current machine doing integer instructions in 20% of time, doing I/O in 35% of time and other operations in the remaining time. which is the better trade off among the following two proposals?
First: Compiler optimization that reduces number of integer instructions by 25% (assume each integer instruction takes the same amount of time); Second: Hardware optimization that reduces the latency of each I/O operations from bus to 5us.

Ans:

Second proposal: Hardware optimization that reduces the latency of each I/O operations from 6 μ s to 5 μ s is the best, because 35% of the time is consumed by I/O operations whereas integer instructions takes only 20% of the time.

Further time take for executing integer instructions is very minimum when compared to time taken for an I/O operation.

Hence, the second is the best proposal.

PART-B

Explain in detail about Flynn's classification of Parallel hardware. [Nov/Dec, 2016, 2015]

Discuss about SISD, MIMD, SIMD, SPMD and VECTOR System [May - 15] [Apr/May - 17]

Explain with diagrammatic illustration Flynn's classification [Nov/Dec 17]

Flynn's classification:

* Parallel processing can be classified in many ways. It can be classified according to internal organization of processors, according to interconnection structure used between processors or according to flow of information through the system.

* In 1966, Michael J. Flynn, classified processor parallelism based on the number of simultaneous instruction and data streams seen by the processor during program execution.

* Flynn's classification divides computers into four major groups No. of I stream No. of data stream

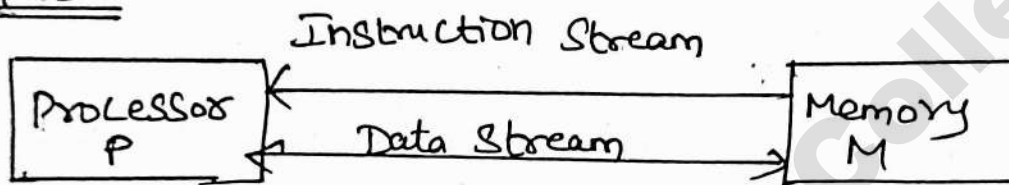
- i) Single Instruction Stream - Single Data Stream (SISD)
- ii) Single Instruction Stream - Multiple Data Stream (SIMD)

iii) Multiple Instruction Streams - Single Data Stream (MISD)

iv) Multiple Instruction Streams - Multiple Data Streams (MIMD)

Instruction and data streams in a sequential

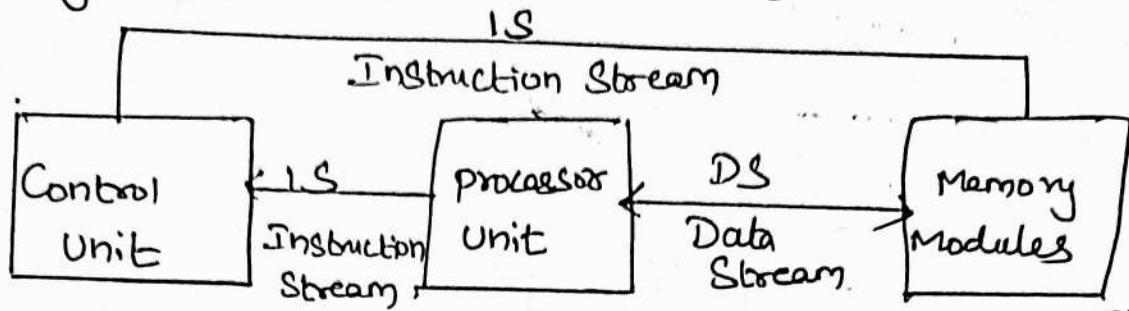
Computer



* A processor operates by fetching instructions and operands from main memory or cache, executing the instructions and placing the final results in memory.

* The instructions form an instruction stream flowing from memory to the processor, while the operands from another stream, data stream flowing to and from the processor.

i. Single Instruction Stream Single Data Stream (SISD)



SISD Computer. Von Neumann computers.
uniprocessors

* Most conventional machines with one CPU containing a single arithmetic-logic unit capable only of scalar arithmetic fall into the category.

* SISD and sequential computers are thus synonymous.

* In SISD computers instructions are executed sequentially but overlap in their execution stages (pipelining).

* They may have more than one functional unit, but all functional units are control by a single control unit.

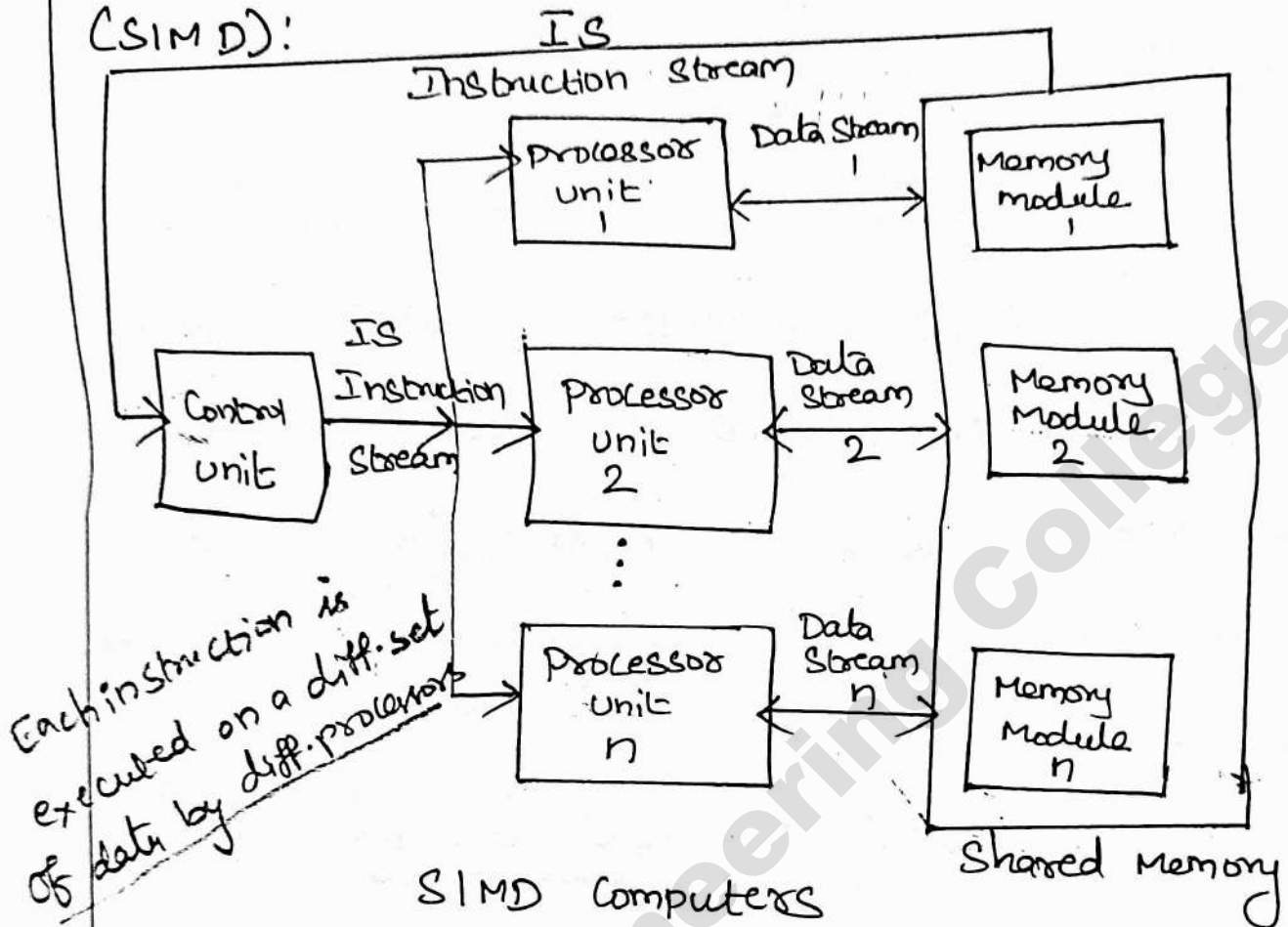
Examples of SISD machines includes

> CDC 6600 which is unpipelined but has multiple function units

> Cray-1 which supports vector processing.

ii) Single Instruction Stream Multiple Data Streams

(SIMD):



* This category corresponds to array processors.

They have multiple processing/execution units and one control unit.

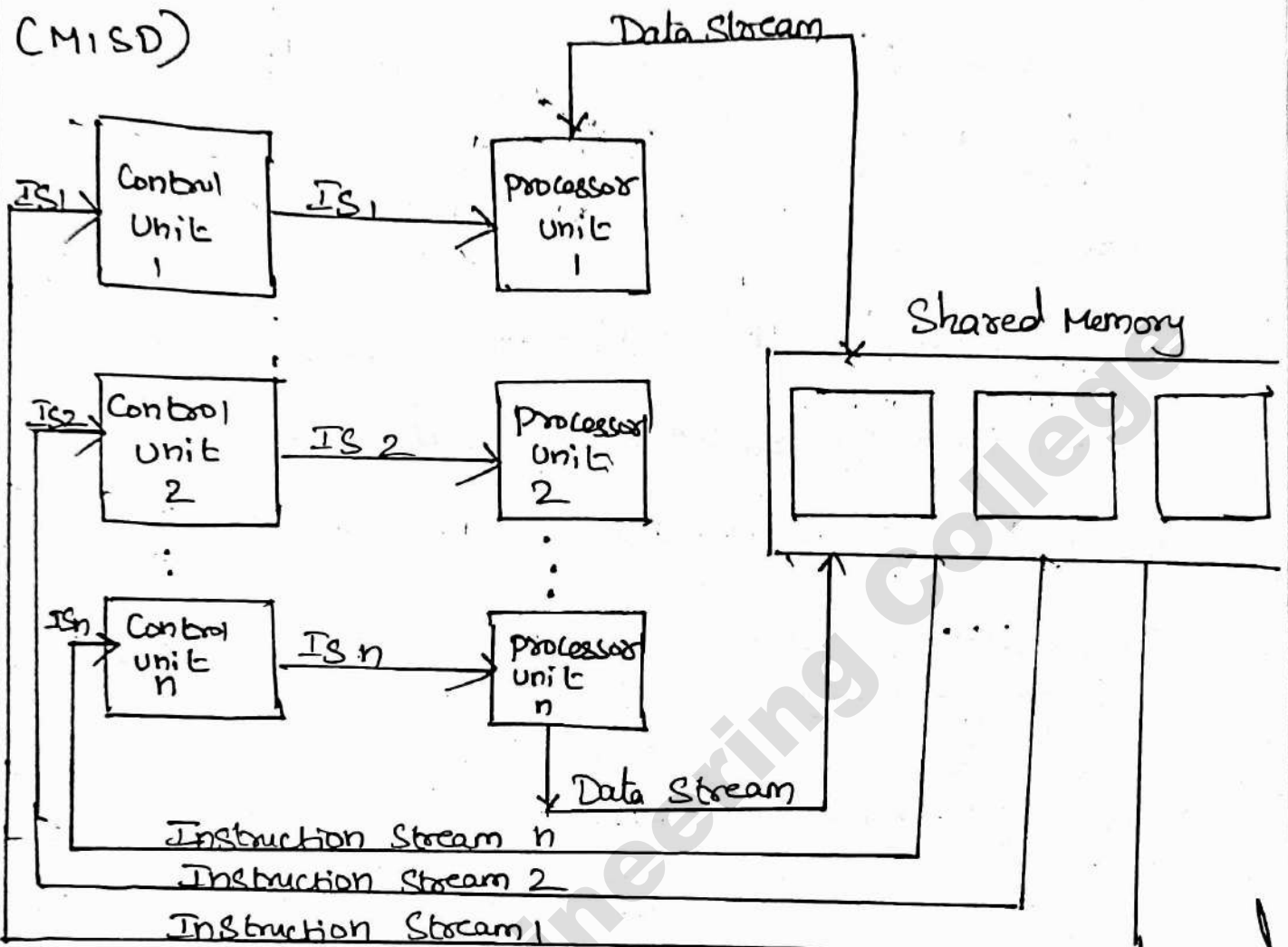
* Therefore, all processing/execution units are supervised by the single control unit.

* Here all processing elements received same instruction from control unit but operate on different data sets from distinct data streams.

* This category is also known as Single program Multiple Data Stream (SPMD)
eg. model G1 pu

iii) Multiple Instruction Streams & single Data Stream

(MISD)



MISD computer *This kind of organization has never been used*

* Not many parallel processors fit well into this category.

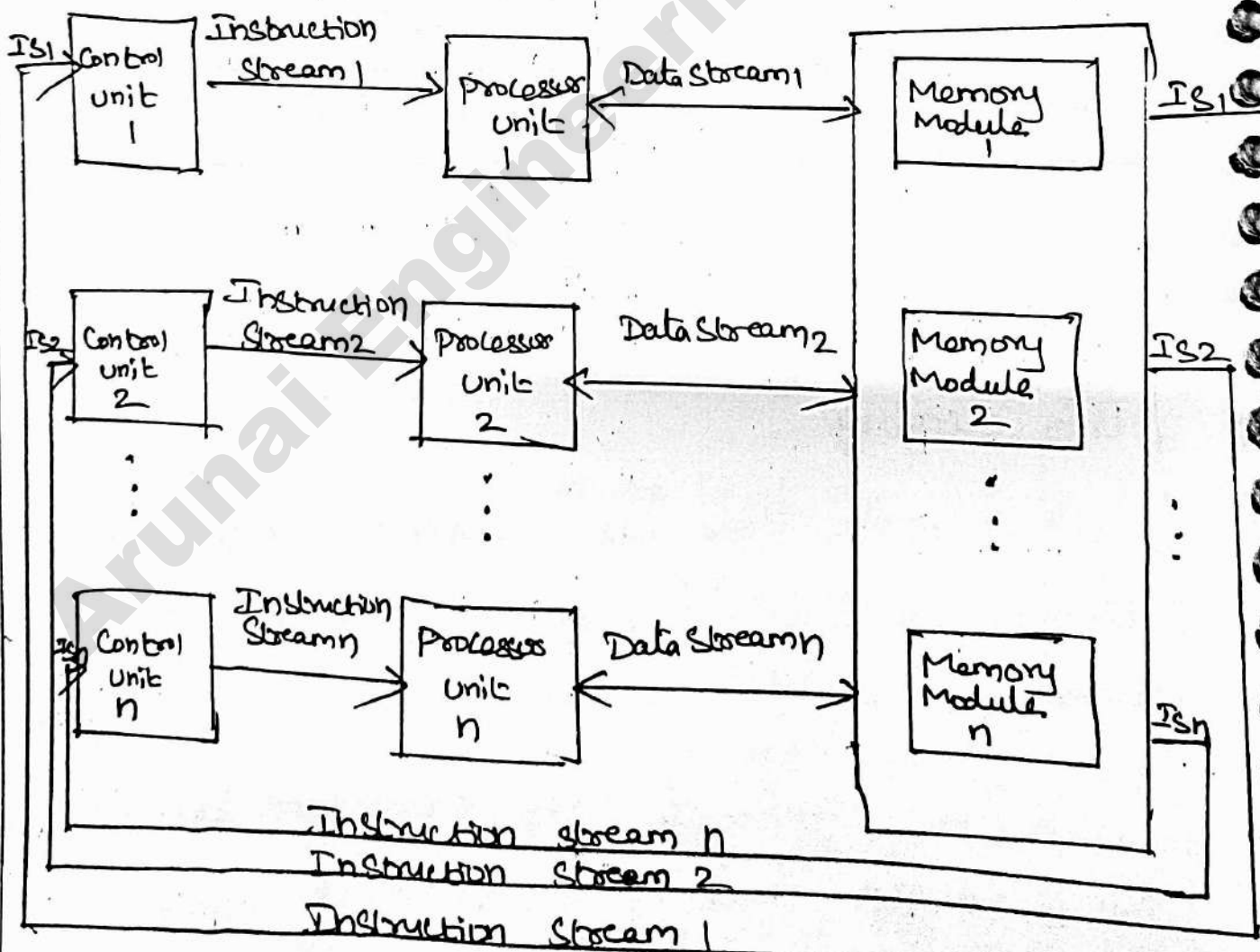
* In MISD, there are n Processor units. Each receiving distinct instructions operating over the same data stream and its derivatives.

* The results of one processor become the input to the next processor in the micropipe.

* The fault-tolerant computers where several processing units process the same data using different programs belongs to the MISD class.

* The results of such apparently redundant computations can be compared and used to detect and eliminate faulty results.

iv) Multiple Instruction Streams Multiple Data Streams (MIMD):



MIMD Computer

* Most multiprocessors system and multiple computers system can be classified in this category.

* In MIMD, there are more than one processors unit having the ability to execute several program simultaneously. 1) Each processor has a separate pgm. 2) An instruction stream is generated from each program.

* MIMD computer implies interactions among the multiple processors because all memory streams are derived from the same data space shared by all processors. 3) Each instruction operates on diff. data.

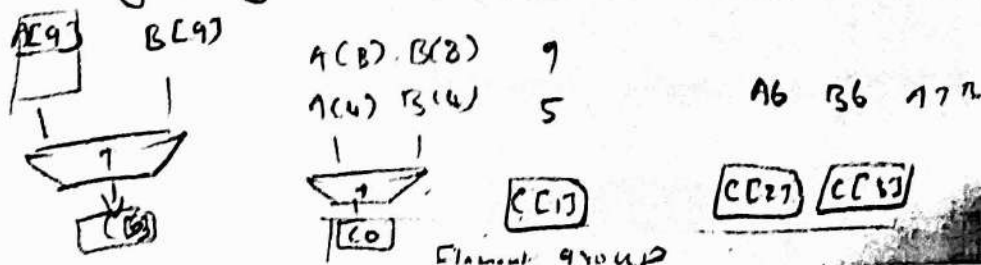
* If the n data streams are derived from disjointed sub spaces of the shared memories then we would have the so-called Multiple SISD (MSISD) operation.

scalar - single data

Vector systems:

* SIMD computers operate on vectors of data and it uses vector architecture.

* Vector architectures pipelined the ALU to get good performance at lower cost.



* The basic philosophy of vector architecture is to collect data elements from memory put them in order into a large set of registers, operate on them sequentially in registers using pipelined execution units and then write the results back to memory.

* Vector architecture has a set of 32 vector registers and each with 64 bit elements

* Vector elements are ^{pipelining parallelism can work well} independent and it can be operated on in parallel,

* All modern vector computers have vector functional units with multiple parallel pipelines called vector lanes

* Vector functional units with multiple parallel pipelines; Pipelines produce two or more results per clock cycle.

Vector - Inst-set contains instruc. that operates on one-dimensional array of data called vector.

Adv: Highly memory access, Reduce instruction fetch bandwidth

Draw: memory can easily bottleneck

What is hardware Multithreading? Compare and Contrast Fine grained and coarse grained multithreading. [May-2015] [May-17]

Explain in detail about hardware Multithreading

Describe Simultaneous Multithreading (SMT) with an example [Nov/Dec-2014] [Nov/Dec-2015]

Hardware Multithreading:-

* Hardware multithreading allows multiple threads to share the functional units of single processors in an overlapping fashion try to utilize the hardware resource efficiently.

* To permit this sharing, the processor must duplicate the independent state of each thread.

Example:

* Each thread would have a separate copy of the register file and the program counter

* The memory itself can be shared through the virtual memory mechanisms it already support multiprogramming.

* Hardware multithreading increase the utilization of a processor by switching to another thread when one thread is stalled.

* Hardware must support the ability to change to a different thread relatively quickly.

* Thread is a light weight process and threads share a single address space but processes don't share. Thread switch is more efficient than a process switch.

* Process includes one or more threads, the address space and the operating system state.

* Process switch invokes the operating system but not a thread switch.

↳ Hardware multithreading has two main approaches such as

1. Fine grained Multithreading
2. Coarse grained Multithreading

1. Fine Grained Multithreading: cycle-by-cycle.

* Fine-grained multithreading switches between threads on each instruction, resulting in interleaved execution of multiple threads.

* This interleaving is often done in a round-robin fashion, skipping any threads that are stalled at that time.

* To make fine-grained multithreading practical, the processor must be able to switch threads on every clock cycle.

Advantage: - no need for dependency checking b/w instructions

* It can hide the throughput losses that arise from both short and long stalls because instruction from other threads can be executed when one thread stalls.

Disadvantage

↳ It slows down the execution of the individual threads because thread that is ready to execute without stalls will be delayed by instructions from other threads.

2. Coarse Grained Multithreading: ^{switch to event} ^{or cache miss} ^{synchronizes events} ^{fault}

* Coarse-grained Multithreading was invented as an alternative to fine-grained multithreading

* Coarse-grained multithreading switches threads only on costly stalls, such as second-level cache misses.

* This change relieves the need to have thread switching be essentially free and is much less likely to slow down the execution of an individual thread, since instructions from other threads will only be issued when a thread encounters a costly stall.

Drawback:

* It is limited in its ability to overcome throughput losses especially from shorter stalls.

* This limitation arises from the pipeline start up costs of coarse grained multithreading issues instructions from a single thread, when a stall occurs the pipeline must be empty.

Fine & coarse grained multithreading can start execution of instructions from only a single thread at a given cycle.

Benefit:

* The new thread will begin executing after the stall must fill the pipeline before instructions will be able to complete.

* Due to this startup overhead, coarse grained multithreading is more useful for reducing the penalty of high cost stalls.

Simultaneous Multithreading (SMT)

Instruction from multiple threads executed concurrently in same cycle

* Simultaneous multithreading is a variation on hardware multithreading that uses the resources of a multiple issue, dynamically scheduled pipelined processor to exploit thread level parallelism at the same time it exploits instruction level parallelism.

* It has multiple processors and more functional unit parallelism available than most single threads can effectively use.

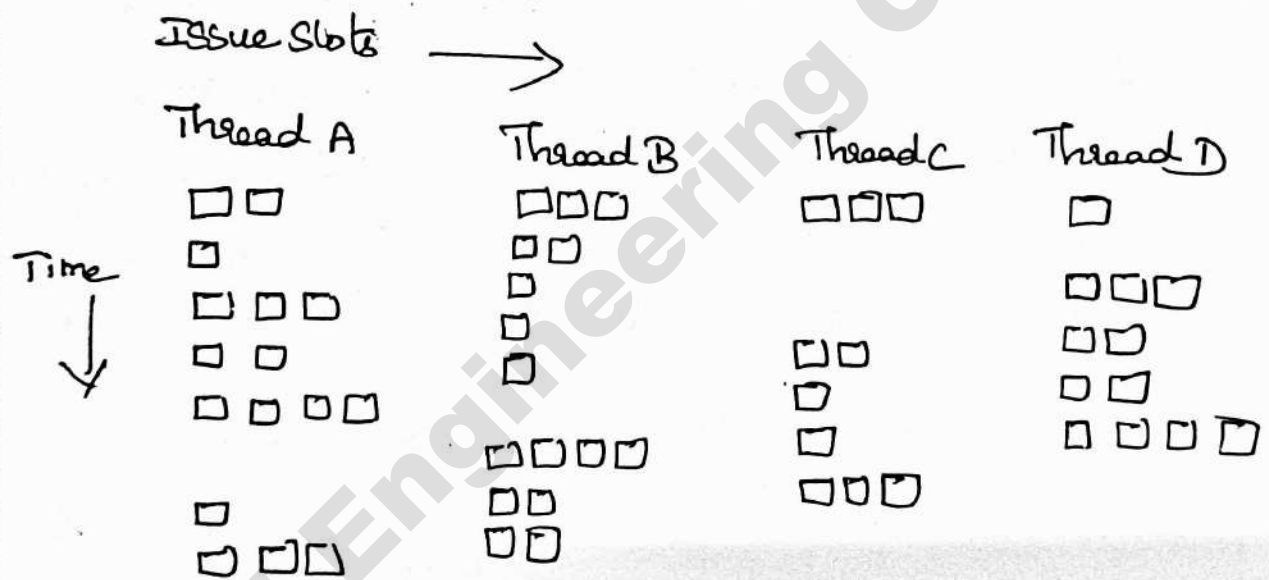
* It has register renaming and dynamic scheduling policy with these features the following tasks can be obtained.

* Multiple instructions from independent threads can be issued without regard to the dependences among them and the resolution of the dependences can be handled by the dynamic scheduling capability.

* SMT relies on the existing dynamic mechanisms. It does not switch resources every cycle.

* SMT is always executing instructions from multiple threads, leaving it up to the hardware to associate instructions slots and renamed registers with their proper threads.

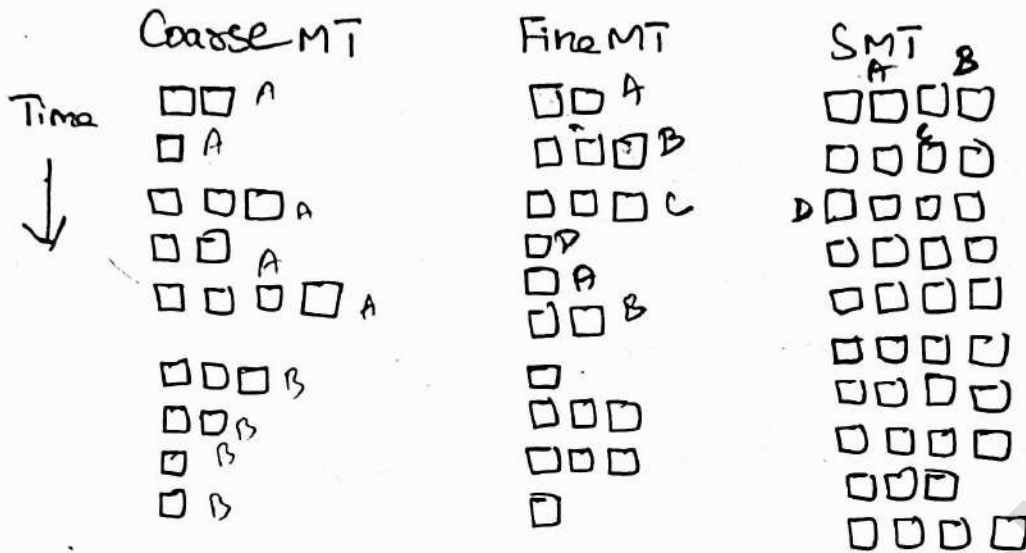
Four Threads Execute independently on a superscalar with no multithreading support.



Four Threads would be combined to execute on the processors more efficiently using these multithreading options.

1. A superscalar with coarse grained multithreading
2. A superscalar with fine grained multithreading
3. A superscalar with simultaneous multithreading

Issue slots



* The horizontal dimension represents the instruction issue capability in each clock cycle and vertical dimension represents a sequence of clock cycle.

Coarse-grained Multithreading:

* The long stalls are partially hidden by switching to another thread that uses the resources of the processors.

* It will reduce the numbers of completely idle clock cycle and the pipeline start up overhead still leads to idle cycles

Fine grained Multithreading:

The interleaving of threads mostly eliminates idle clock cycles. Because only a single thread issues instructions in a given clock cycle.

* Limitations in instruction level parallelism will lead to idle slots within some clock cycles.

Simultaneous Multithreading:

* In SMT thread level parallelism and instruction level parallelism both are exploited with multiple threads using the issue slots in a single clock cycle.

* Ideally the issue slot usage is limited by imbalances in the resources needs and resource availability over multiple threads, but in practice other factors can restrict how many slots are used.

3. Explain Instruction level parallel processing.

State the challenges of Parallel processing [Nov/Dec-2014]

*

Instruction Level Parallelism:

* pipelining is a technique that runs programs faster by overlapping the execution of instructions.

This is one example of Instruction level parallelism

* Pipelining exploits the potential parallelism among instructions. This parallelism is called Instruction-level parallelism (ILP).

* There are two primary methods for increasing the potential amount of instruction-level parallelism

1. Increasing the depth of the pipeline to overlap more instructions.

2. Replicate the internal components of the computers.

* Another approach is to replicate the internal components of the computers so that it can launch multiple instructions in every pipeline stage.

* The general name for this technique is multiple issue.

* Launching Multiple Instructions per Stage will allow the instruction execution rate to exceed the clock rate or the CPI to be less than 1.

* It is sometimes useful to flip the metric and use IPC or instructions per clock cycle.

* There are two major ways to implement a multiple issue processors such as

1. Static multiple issue

2. Dynamic multiple issue

* The major differences between these two kinds of issues are the division of work between the compiler and the hardware, because the division of work dictates whether decisions are made at compile time or during execution time.

* A combination of compiler based optimization and hardware techniques can be used to maximize Instruction level parallelism.

Multiple Issue Pipeline:

There are two primary and distinct responsibilities that must be dealt within a multiple issue pipeline such as

1. packaging Instructions into issue slots
2. Dealing with data and control hazards

Concept of speculation:

* one of the most important methods for finding and exploiting more instruction level parallelism is speculation.

* It is an approach that allows the compiler or the processor to guess about the properties of an instruction so to enable execution that begin for other instruction. may depend on the speculated instruction.

* speculation is an approach whereby the compiler or process, guesses the outcome of an instruction to remove its dependence in executing other instructions.

* speculation may be done in the compiler or by the hardware.

Eg The compiler can use speculation to reorder instructions, moving an instruction across a branch or load across a store.

Difficulty with speculation:

* Difficulty with speculation is it may be wrong. So any speculation mechanism must include a method to check if the guess was right and a method to undo or backout the effects of the instructions was executed speculatively. But it will add complexity.

1: Packaging Instructions into issue slots:

* Issue slots are the positions from which instructions could issue in a given clock cycle and it correspond to positions at the starting blocks for a sprint

* How does the processor determine how many instructions and which instructions can be issued in a given clock cycle?

* In static issue processors, this process is partially handled by the compiler and in case of dynamic issue processors it is dealt with runtime by the processor

2. Dealing with data and control hazards

* In static issue processors, the compiler handled some or all of the data and control hazards statically.

* In dynamic issue processors at least some classes of hazards using hardware techniques operating at execution time.

* The recovery mechanisms used for incorrect speculation is done in two ways

1. Hardware speculation
2. software speculation

1. Hardware speculation:

* In hardware speculation the processor usually buffers the speculative results until it knows they are no longer speculative

* If the speculation is correct, the instructions are completed by ^{allowing} the contents of the buffers to be written to the registers or memory.

* If the speculation is incorrect, the hardware flushes the buffers and re-executes the correct instruction sequence.

2. Software speculation:

* In software speculation the compiler usually inserts additional instructions that check the accuracy of the speculation and provide a fix up routine to use when the speculation is incorrect

* speculation can improve performance when it is done properly and decrease performance when it is done carelessly

Static Multiple Issue Processors:

* In static multiple issue, all processors use the compiler to assist with packaging instructions and handling hazards. It will use the Issue Packet.

* Issue Packet is set of instructions that issue together in one clock cycle and the packet may be determined statically by the compiler or dynamically by the processor.

* Static multiple issue processors usually restricts the mix of instruction that can be initiated in a given clock cycle.

* So issue packet is a single instruction that allows several operations in certain predefined fields. This approach is called Very Long Instruction word (VLIW)

* VLIW is a style of instruction set architecture that launches many operations that are defined to be independent in a single wide instruction and it has many separate opcode fields.

* In most static issue processors compiler takes some of the responsibility for handling data and control hazards

The responsibilities taken by the compiler are

1. Static branch Prediction
2. Code Scheduling to reduce or prevent all hazards.

An Example: Static multiple Issue with MIPS ISA

Static multiple Issue processor functions can be explained using simple two Issue MIPS processor.

↳ In two Issue MIPS processor one of the instruction is stored in integer ALU operation or branch

↳ Other instruction can be a load or store.

Issuing two instructions per cycle require fetching and decoding 64 bits of instructions.

* In many static multiple Issue processors and all VLIW processors issuing instructions simultaneously is restricted to simplify the decoding and instruction issue.

* So instruction must be paired and aligned on a 64 bit boundary with the ALU or branch position must appear first.

Static two issue Pipeline in operation

Instruction type	Pipe Stages							
	IF	ID	EX	MEM	WB			
ALU or branch Instruction	IF	ID	EX	MEM	WB			
Load or store Instruction	IF	ID	EX	MEM	WB			
ALU or branch Instruction		IF	ID	EX	MEM	WB		
Load or store Instruction		IF	ID	EX	MEM	WB		
ALU or branch Instruction			IF	ID	EX	MEM	WB	
Load or store Instruction			IF	ID	EX	MEM	WB	
ALU or branch Instruction				IF	ID	EX	MEM	WB
Load or store Instruction				IF	ID	EX	MEM	WB

* ALU and data transfer instructions are issued at the same time.

* Keeping the register writes at the end of the pipeline will simplify the handling of exceptions and maintenance of a precise exception model. But it is more difficult in multiple issue processors.

* In some static multiple issue processors the compiler takes full responsibility for removing all hazards, scheduling the code and inserting no-ops.

* For these kinds of design the code executes without any need for hazard detection or hardware generated stalls.

* Hardware - It can detect data hazards and generates stalls between two issue packets for that compiler avoids all dependences within an instruction pair.

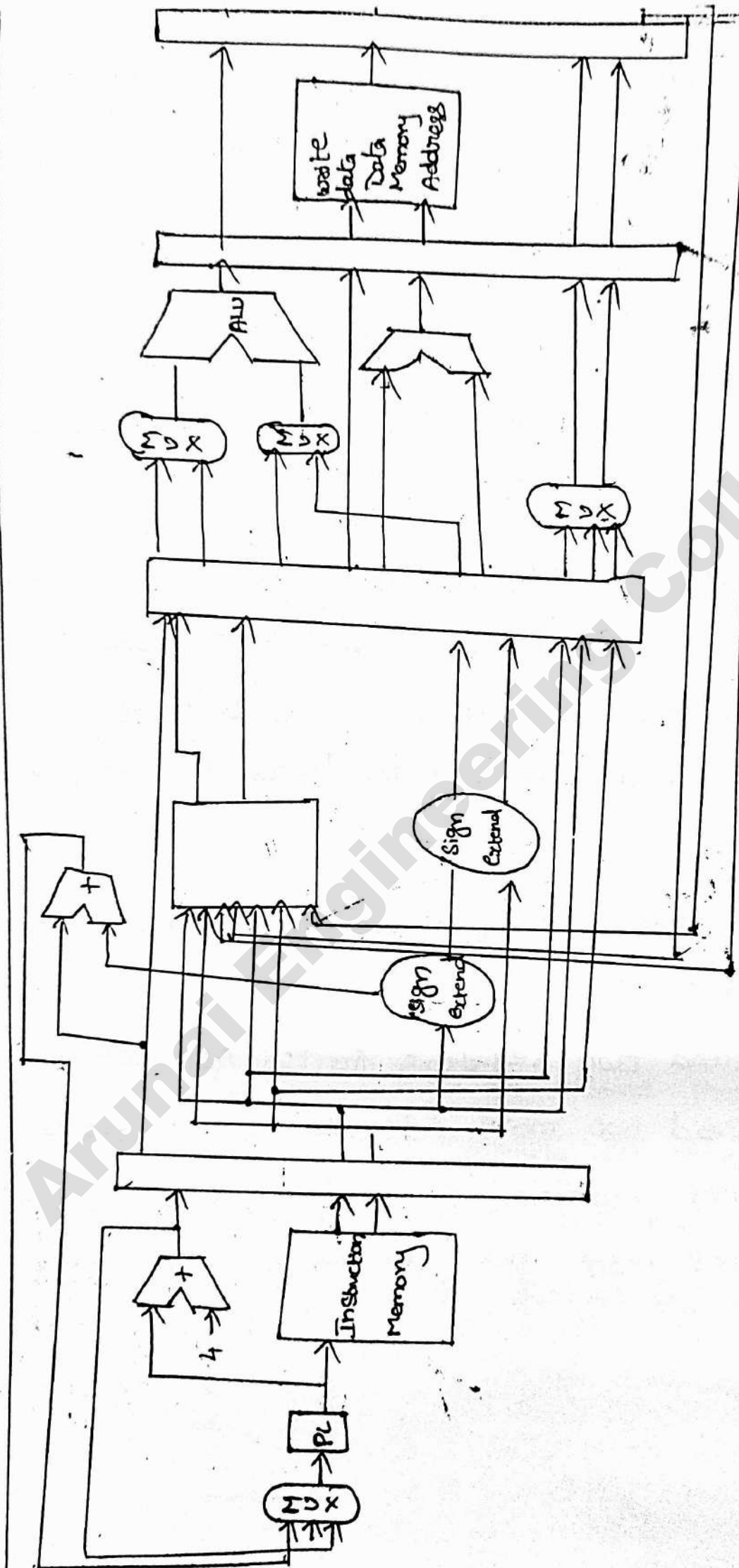
* Software - It must handle all hazards or only try to reduce the fraction of hazards between separate issue packets.

* To perform ALU and data transfer operation in parallel, first we need an additional hardware.

* With additional hardware units some other units also necessary such as hazard detection and stall logic.

In one clock cycle we need to perform the following tasks

1. Read two registers for the ALU operation
2. Read two more registers for a store.
3. one register for write port for a load
4. one write port for an ALU operation



A Static Two-Issue Datapath

* Two-issue processors can improve performance by up to a factor of 2. But many instructions are overlapped and additional overlap increases the relative performance loss from data and control hazards.

2. Dynamic Multiple Issue Processors:

* Dynamic multiple issue processors are also known as super scalar processors or simply superscalars.

* Super scalar is an advanced pipelining technique that enables the processor to execute more than one instruction per clock cycle by selecting them during execution.

* Achieving good performance on dynamic multiple issue processors it requires the compiler to schedule instructions to move dependences and improve the instruction issue rate.

* In super scalar processors, whether the code has scheduled or not it is given to the hardware to execute correctly.

* In super scalar processors, the compiled code can run always correctly independent of the issue rate or pipeline structure of the processor.

* Some static issue processors the code can run correctly across different implementations but it requires poorly compilation.

* Many superscalars extend the basic framework of dynamic issue decisions by including dynamic pipeline scheduling.

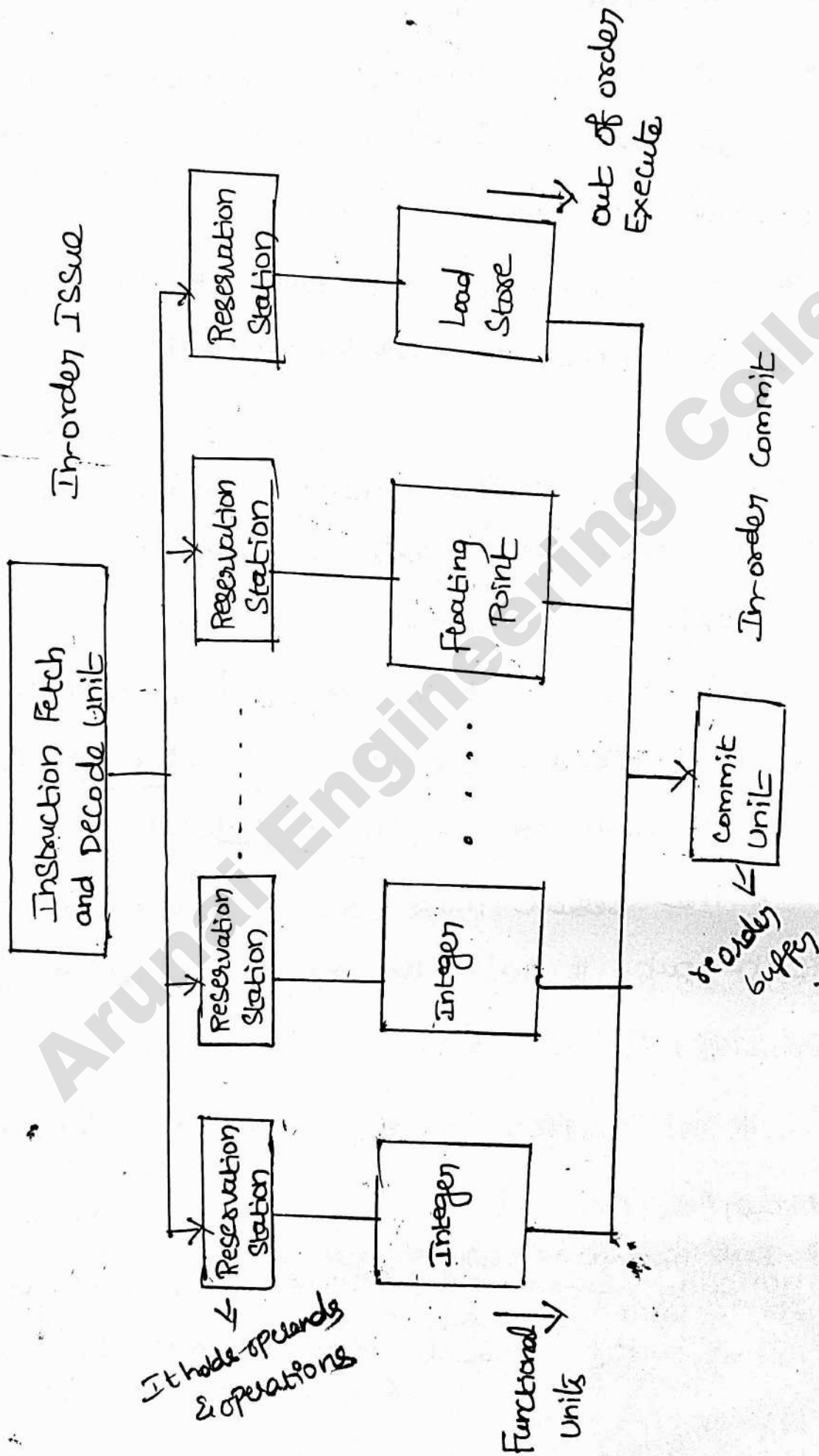
* Dynamic pipeline scheduling chooses which instructions to execute in a given clock cycle while trying to avoid hazards and stalls.

Dynamic pipeline scheduling:

* Dynamic pipeline scheduling chooses which instructions to execute next by reordering the instructions it avoid stalls.

* In such processors, the pipeline is divided into three major units

1. An instruction fetch and issue unit
2. Multiple functional units
3. Commit unit



The Three primary units of a dynamically scheduled pipeline

* The final step of Updating the state is also called retirement or graduation

* The first unit fetches the instructions, decodes them and sends each instruction to a corresponding functional unit for execution

* Each functional unit has a buffer that is called reservation stations. It holds the operands and operation.

* Once the buffer contains all its operands and the functional unit is ready to execute then the result is calculated.

* When the result is completed, it is sent to any reservation stations waiting for this particular result as well as to the commit unit

* The result must be buffered until it is safe to put it into the register file or for a store or into memory.

* The buffer in the commit unit is called the reorder buffer. It holds the results in a dynamically scheduled processor until it is safe to store the results into memory or a register

* once a register is committed to the register file it can be fetched directly from there just like a normal pipeline.

* Form of register renaming is obtained by combination of buffering operands in the reservation stations and results in the reorder buffer

Below two steps are explained how the renaming of register is obtained.

Step 1:

* When an instruction issues, it is copied to a reservation station for the appropriate functional unit.

* If any operands are available in the register file or reorder buffer are also immediately copied into the reservation station.

* The instruction is buffered in the reservation station until all the operands and the functional units are available.

* For issuing instruction, the register copy of the operand is no longer required and for a write to that register occurred means the value could be overwritten.

Step 2:

* If an operand is not in the register file or reorder buffer, it must be waiting to be produced by a functional unit.

* The name of the functional unit will produce the result is tracked, when the functional unit eventually produces the result it is copied directly into the waiting reservation station by passing the registers from function unit.

↳ These two steps use the reorder buffer and the reservation stations to implement register renaming.

out of order execution

* The processor executes the instruction in some order that preserves the data flow order of the program. This style of execution is called out of order execution.

In-order commit

* A commit in which the results of pipelined execution are written to the programmer visible state in the same order how the instructions are fetched.

Discuss the challenges in parallel processing with necessary eg: ~~Example~~
Parallel Processing Challenges: \uparrow performance - 17
 \downarrow time to execute

* Parallel processing will increase the performance of processors and it will reduce the utilization time to execute task. But obtaining the parallel processing is not an easy task

* The difficulty with parallelism is not in the hardware side it is in the software side.

* It is difficult to write software that uses multiple processors to complete one task faster and the problem gets worse as the number of processors increases.

* Developing the parallel processing programs are so harder than the sequential programs because of the following reasons.

1. It must get better performance or better energy efficiency from a parallel processing program on a multiprocessor.

2. If we would not have efficient energy then we have to use a sequential program on a uniprocessor because sequential programming is very simple.

* Using Uniprocessors we can solve the problem in developing parallel processing programs.

* Because uniprocessor design techniques such as superscalar and out of order Execution.

take advantage of instruction Level Parallelism

* Uniprocessor design techniques will reduce the demand for rewriting programs for multiprocessors.

* If number of processors level increased mean it is more difficult to write parallel processing programs.

* The following reasons we cannot get parallel processing programs as faster than sequential programs

1. Scheduling
2. Partitioning the work into parallel pieces
3. Balancing the load evenly between the workers
4. Time to Synchronize
5. Overhead for communication between the parties

1. Scheduling:

* Scheduling is the method by which threads, processes or data flows are given access to system resources

eg processor time, communications bandwidth

* Scheduling is done to load balance and share system resources effectively or achieve a target quality of service.

* Scheduling can be done in various fields among that process scheduling is more important, because in parallel processing we need to schedule the process correctly

↳ Process scheduling can be done in the following ways

1. Long term scheduling
2. Medium term scheduling
3. Short term scheduling
4. Dispatcher.

2. Partitioning the work

* The task must be broken into equal number of pieces because otherwise some task may be idle while waiting for the ones with larger pieces to finish

* To obtain Parallel processing task must be divided into equally to all the processors. Then only we can avoid the idle time of any processors

3 Balancing the load:

* Load balancing is the process of dividing the amount of work that a computer has to do between two or more processors. So that more work gets done in the same amount of time and in general all process get served faster.

* work load has to be distributed evenly between the processors to obtain Parallel processing task.

Time to Synchronize:

* Synchronization is the most important challenge in Parallel processing. Because all the processors have equal work load so it must complete the task within specific time period.

* For Parallel processing program it must have time to synchronization process, because if any process does not complete the task within specific time period then we cannot obtain Parallel processing

Two methods are found to increase the scale up,

Such methods are

1. Strong scaling
2. Weak scaling

1. Strong Scaling:

* In this method speed up is achieved on a multiprocessor without increasing the size of the problem.

* "Strong scaling means measuring speed up while keeping the problem size fixed".

2. Weak scaling:

* In this method speed up is achieved on a multiprocessor while increasing the size of the problem proportionally to the increase in the number of processors.

* "weak scaling means the problem size grows proportionally to the increase in the number of processors".

* Let's assume that the size of the problem is M is the working set in main memory and we have P processors.

* Then the memory per processor for strong scaling is approximately M/P and for weak scaling it is approximately M .

* By considering the memory hierarchy it is easy to know that weak scaling being easier than strong scaling.

* But if the weakly scaled dataset no longer fits in the last level cache of a multicore microprocessors then the resulting performance could be much worse than by using strong scaling.

* Depending on the application we can select the scaling approach.

4. Discuss shared memory multiprocessors with a neat diagram. [Nov/Dec-2016]

Explain the term Multicore processor [Nov/Dec-14]

Multicore processors:

* Hardware multithreading improved the efficiency of processors but it has big challenge to deliver on the performance potential of Moore's law by efficiently programming the increasing number of processors per chip.

* Rewriting old programs to run well on parallel hardware is more difficult. Computer designers must do something to overcome these difficulties.

* Solution for above program is using a single physical address space for all processors so that programs need not concern themselves with where their data is and programs may be executed in parallel.

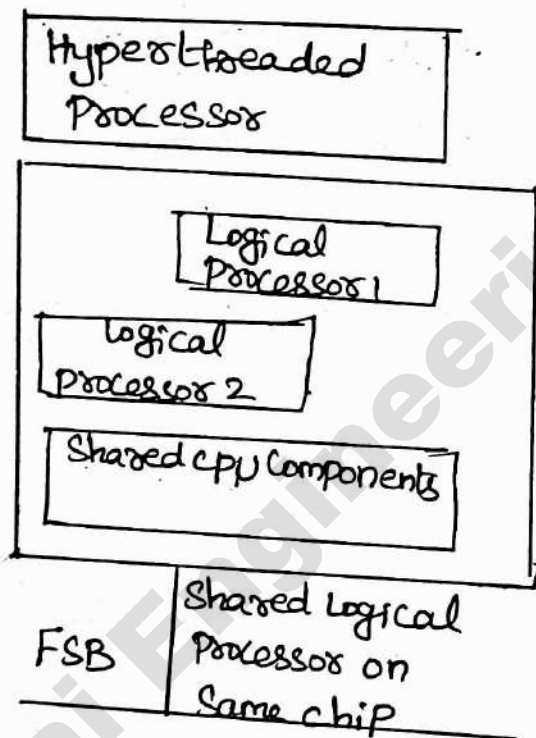
↳ Multicore architectures are classified according to

i) Number of processors: Two processors (dual core), four processors (Quad core) and eight processors (octa-core) configuration.

- 2) Approaches to processor-to-processor communication
- 3) cache and memory implementations
- 4) Implementation of I/O bus and the Front Side Bus (FSB)

Three common configurations that support multiprocessing.

Type 1:



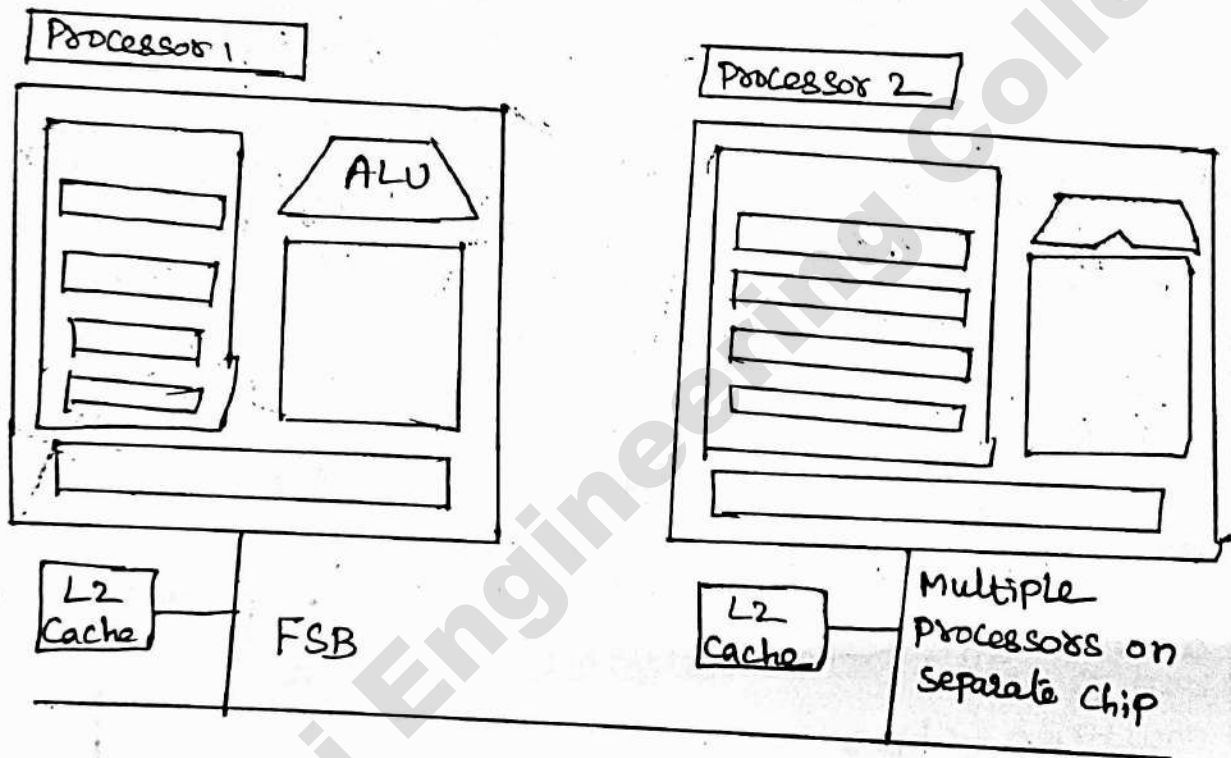
* It uses hyperthreading technology. It allows a single processor to act like two separate processors to the operating system and the application programs that use it.

* However, there is a single processor running multiple threads.

* Thus, in hyperthreaded technology the multiple processors are logical instead of physical.

* In hyperthreaded technology has some duplication of hardware but not enough to qualify a separate physical processors.

Type 2:

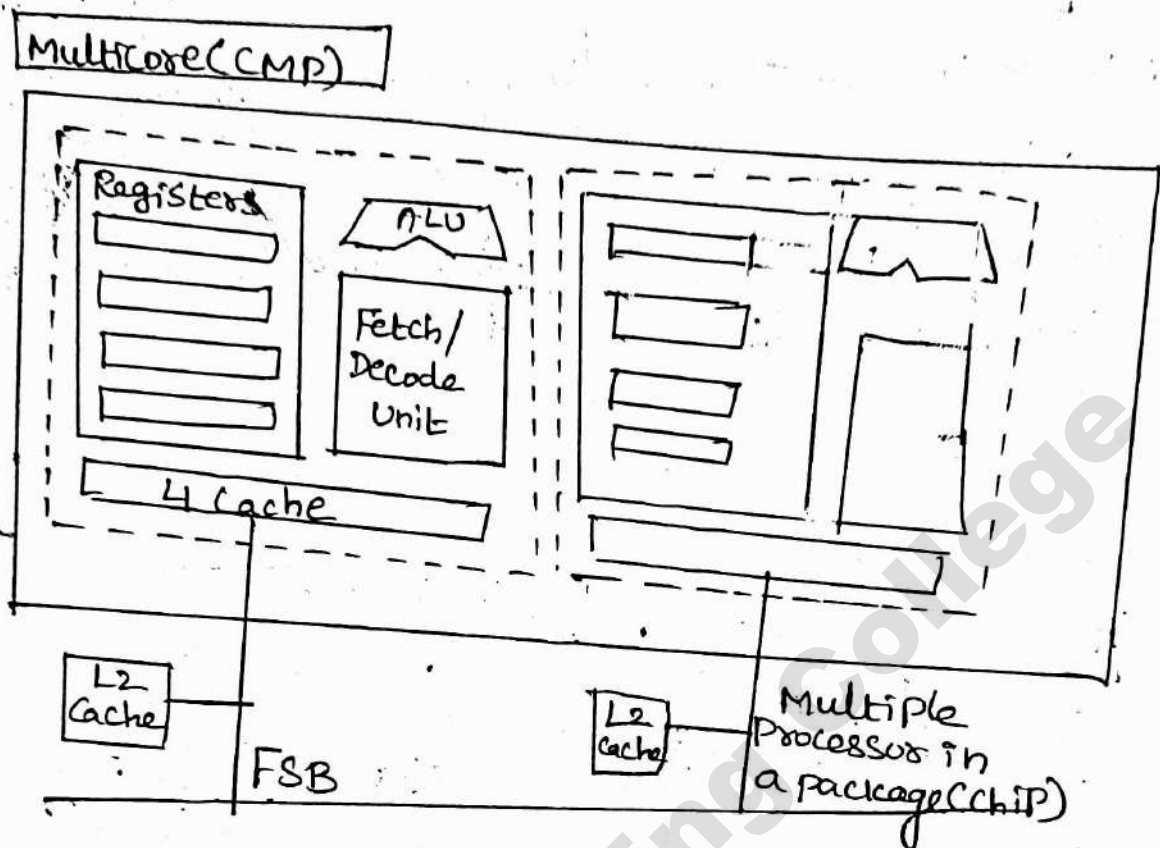


* It is the classic multiprocessor. A multiprocessor is a tightly coupled computer system having two or more processing units (multiple processors) on a separate chip with its own hardware.

Type 3:

* It represents multicore system. It provides two or more complete processors on a single chip.

Type 3:

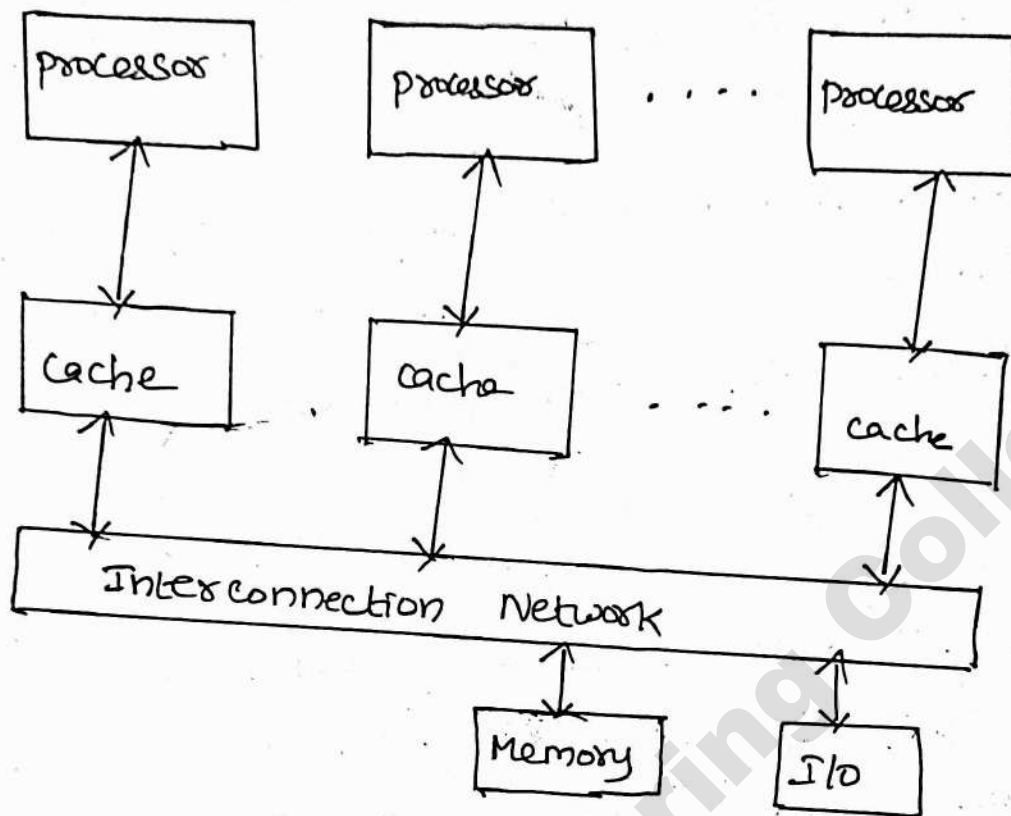


* A Shared Memory Multiprocessor (SMP) allows programmer to use single physical address space across all processors which is always suitable for multicore chips

* In shared address multiprocessor processors can communicate through variables in memory with all processors capable of accessing any memory location via loads and - stores.

* The classic organization of a shared memory multiprocessor. It can run independent jobs in their own virtual address spaces even if they all share a physical address space.

Types of single address



Classic organization of a shared memory multiprocessor

Single address space multiprocessor comes in two styles such as

1. Uniform Memory Access (UMA)
2. Non uniform Memory Access (NUMA)

* UMA is a multiprocessor in which latency to any word in main-memory is same no matter which processor requests the access.

* NUMA is a type of single address space multiprocessor in which some memory access are much faster than others depending on which processor asks for which word

Uniform Memory Access	Nonuniform Memory Access
1. Programming challenges are easy 2. UMA machines can scale small sizes 3. It has higher latency	1. Programming challenges are hard 2. NMA machines can scale to larger sizes 3. It has lower latency to nearby memory.

Synchronization:

* As processors operating in parallel will normally share data, they also need to coordinate with operating on shared data otherwise one processor can start working on data before another is finished with it. This coordination is called Synchronization.

* Synchronization is the process of coordinating the behavior of two or more processes which may be running on different processors

* When sharing is supported with a single address space there must be a separate mechanism for synchronization such as called Lock

* Lock is a synchronization device that allows access to data to only one processor at a time and other processors interested in shared data must wait until the original processor unlocks the variable.

5.

A pipelined processor uses delayed branch technique. Recommend any one of the following possibilities for the design of the processor. In the first possibility the processor has a 4-stage pipeline and one delay slot. In the second possibility, it has a 6-stage pipeline and two delay slots. Compare the performance of these two alternatives, taking only the branch penalty into account. Assume that 20% of the instructions are branch instructions and that an optimizing compiler has an 80% success rate in filling in the single delay slot. For the second alternative the compiler is able to fill the second slot 25% of the time. [APP/MAY-2017]

Solution:

$$\text{Speed up factor} = \frac{\text{Pipeline depth}}{1 + \text{Pipeline stall per instruction}}$$

$$\text{Pipeline stall per instruction} = \text{Branch instruction} \times \text{Number of delays} - \text{Branch instruction} \times \text{Success rate in 1st delay} - \text{Branch instruction} \times \text{Success rate in 2nd delay}$$

$$\text{Branch Instruction} = 20\% = 0.2$$

$$\text{Success rate in 1st delay} = 80\% = 0.80$$

$$\text{Success rate in 2nd delay} = 25\% = 0.25$$

4-stage pipeline:

$$* \text{ Pipeline depth} = 4$$

$$\text{Number of delay} = 1$$

$$\begin{aligned} \text{Pipeline Stall per instruction} &= 0.2 \times 1 - 0.2 \times 0.8 \\ &= 0.2 - 0.16 \\ &= 0.04 \end{aligned}$$

$$\text{Speed up factor} = \frac{4}{1+0.04} = \frac{4}{1.04} = 3.846$$

6-stage pipeline:

$$* \text{ Pipeline depth} = 6$$

$$\text{Number of delays} = 2$$

$$\begin{aligned} \text{Pipeline Stall per instruction} &= 0.2 \times 2 - 0.2 \times 0.8 - 0.2 \times 0.25 \\ &= 0.4 - 0.16 - 0.05 \\ &= 0.4 - 0.21 \\ &= 0.19 \end{aligned}$$

$$\text{Speed up factor} = \frac{6}{1+0.19} = \frac{6}{1.19} = 5.042$$

\therefore Speed up factor of 6-stage pipeline is high, So it will be recommended.

6. Suppose you want to achieve a speed-up of 90 times faster with 100 processors. What percentage of the original computation can be sequential?

[Nov/Dec-17]

Solution:

$$\text{Speedup} = \frac{1}{(1 - F_e) + \frac{F_e}{S_e}} \rightarrow (i)$$

Where

$$\text{Speedup} = 90$$

F_e = Fraction enhanced, is always less than or equal to 1

S_e = Speedup, is always greater than 1

(i.e) 100 (Number of processors)

Using the values in equation (i) we will get,

$$90 = \frac{1}{(1 - F_e) + \frac{F_e}{100}}$$

$$90 \left[(1 - F_e) + \frac{F_e}{100} \right] = 1$$

$$90(1 - F_e) + \frac{90F_e}{100} = 1$$

$$90 - 90F_e + 0.9F_e = 1$$

$$90 - 1 = 90f_e - 0.9f_e$$

$$89 = (90 - 0.9)f_e$$

$$89 = 89.1f_e$$

or

$$89.1f_e = 89$$

$$\therefore f_e = \frac{89}{89.1} = 0.999$$

$$f_e = 0.999$$

Thus to achieve a speed-up of 90 from hundred processors, the sequential percentage can only be 0.1%. (i.e., 0.001)

ii)

Suppose you want to perform two sums; one is a sum of 10 scalar variables, and one is a matrix sum of a pair of two-dimensional arrays, with dimensions 10 by 10. For now let's assume only the matrix sum is parallelizable; we'll see how to parallelize scalar sums. What speed-up do you get with 10 versus 40 processors? Next, calculate the speed-ups assuming the matrices grow to 20 by 20.

[Nov/Dec-17]

Solution:

If we assume performance is a function of the time for an addition, t , then there are 10 additions that do not benefit from parallel processors and 100 additions that do. If the time for a single processor is $110t$, the execution time for 10 processors is

Execution time after improvement =

$$\frac{\text{Execution time affected by improvement}}{\text{Amount of improvement}} + \text{Execution time unaffected}$$

$$\begin{aligned} \text{Execution time after improvement} &= \frac{100t}{10} + 10t \\ &= 20t \end{aligned}$$

So, the speed-up with 10 processors is $110t/20t = 5.5$
The execution time for 40 processors is

$$\begin{aligned} \text{Execution time after improvement} &= \frac{100t}{40} + 10t \\ &= 12.5t \end{aligned}$$

So, the speed-up with 40 processors is $110t/12.5t = 8.8$

* Thus, for this problem size, we get about 55% of the potential speed-up with 10 processors, but only 22% with 40.

* Look what happens when we increase the matrix. The sequential program now takes $10t + 400t = 410t$. The execution time for 10 processors is

$$\begin{aligned} \text{Execution time after improvement} &= \frac{400t}{10} + 10t \\ &= 50t \end{aligned}$$

* So the Speed-up with 10 processors is $\frac{410t}{50t}$
 $= 8.2$

The Execution time for 40 processors is

$$\begin{aligned} \text{Execution time after improvement} &= \frac{40t}{40} + 10t \\ &= 20t \end{aligned}$$

* So the Speed-up with 40 processors is $\frac{410t}{20t}$
 $= 20.5$

* Thus, for this larger problem size, we get 82% of the potential Speed-up with 10 processors and 51% with 40.

* These problem show that getting good Speed-up on a multiprocessor while keeping the problem size fixed is harder than getting good Speed-up by increasing the size of the problem.

UNIT-V

PART-A

1. What is DMA? [Nov/Dec-2014]

The CPU and I/O controller interact only when the CPU must yield control of the memory bus to the I/O controller in response to requests from the latter. This level of I/O control is called Direct Memory Access (DMA)

2. Differentiate programmed I/O and Interrupt I/O [Nov-2014]

Programmed I/O	Interrupt I/O
The time that the CPU spends for testing I/O device status and executing I/O data transfer is high	The time that CPU spends for testing I/O device status and executing I/O data transfer is low

dot

3. What is the purpose of dirty/modified bit in cache memory? (Nov-Dec 2014)

A dirty bit or modified bit is a bit that is associated with block of computer's memory and it indicates whether or not the corresponding block of memory has been modified. This bit is set when the processor writes to this memory.

4. What is the need to implement memory as a hierarchy?

- * Decreasing cost per bit (May-2015)

- * Increasing capacity

- * Increasing access time

- * Decreasing frequency of access of the memory by the processors.

5. Point out how DMA can improve I/O speed.

[May/June-2015]

- * The I/O device is connected to the system bus via a special interface circuit called DMA controller. In this data counter stores the number of words that remain to be transferred.

- * It is automatically incremented or decremented after each transfer and tested for zero, when DC reaches

6. What are the various memory technologies? (Nov/Dec-2015)

Computer memory is classified as primary memory and secondary memory.

i) Primary memory technologies in use:

Static RAM, Dynamic RAM, ROM and its types (PROM, EPROM, EEPROM) etc.

ii) Secondary memory technologies in use:

Optical disk, flash-based removal memory card, hard-disk drives, magnetic tapes etc.

7. Define Hit ratio. (Nov-Dec-2015)

The percentage of accesses where the processor finds the code or data word it needs in the cache memory is called the hit rate or hit ratio.

8. Define Memory hierarchy (May/June-2016)

A memory hierarchy is a structure of memory that uses multiple levels of memories as the distance from the processor increases; the size of the memories and the access time both increase.

9. State the advantage of virtual memory (may/june-2016)

1. Allocating memory is easy
 2. Eliminates external fragmentation
 3. Allows demand paging and pre-paging
 4. More efficient swapping
 5. It is the ability to load and execute a process that requires a larger amount of memory than what is available by loading the process in parts and then executing them.
-

10. Define locality of reference. What are its types?

* The program may contain a simple loop, nested loops, or a few procedures that repeatedly call each other. The point is that many instructions in localized area of the program are executed repeatedly during some time period and the remainder of the program is accessed relatively infrequently. This is referred as locality of reference.

11 Define Memory Interleaving [APR/MAY - 2017]

* In computing, interleaved memory is a design made to compensate for the relatively slow speed of dynamic random-access memory (DRAM) or core memory, by spreading memory addresses, evenly across memory banks

* Memory interleaving is a category of techniques for increasing memory speed.

Ex: with separate memory banks for odd and even addresses, the next byte of memory can be accessed while the current byte is being refreshed.

12. Summarize the sequence of events involved in handling (handling) an interrupt request from a single device. [APR/MAY-2017]

* An interrupt is an unscheduled event that disrupts program execution, used to detect overflow.

↳ Processor checks for interrupts

↳ If no interrupts, fetch the next instruction for the current program.

↳ If an interrupt is pending, suspend execution of the current program, and execute the interrupt handler.

Multiple Interrupts:

Two approaches to handle multiple interrupts

1. Define priorities for interrupts. Allow interrupts to interrupt interrupts.
2. Disable interrupt while an interrupt is being processed.

13 What is virtual memory? [NOVIDEC-2017]

* virtual memory is a technique that uses main memory as a "cache" for secondary storage.

Motivations for virtual memory are

1. To allow efficient and safe sharing of memory among multiple programs.

2. To remove the programming burdens of a small, limited amount of main memory.

14 How many total bits are required for a direct-mapped cache with 16KB of data and 4-word blocks, assuming a 32-bit address? [NOVIDEC-2017]

* We know that 16KB is 4K words, which is 2^{12} words and with a block size of 4 words (2^2), 2^{10} blocks. Each block has 4×32 or 128 bits of data plus a tag, which is $32 - 10 - 2 - 2$ bits, plus a valid bit.

Total cache size is

$$2^{10} \times (128 + (32 - 10 - 2 - 2) + 1) = 2^{10} \times 147 = 147 \text{ Kbits}$$

147 kbits or 18.4 kB for a 16 kB cache.

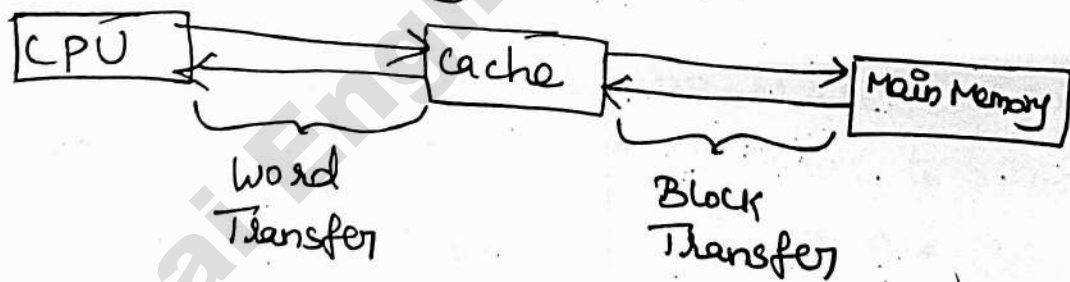
* For this cache, the total number of bits in the cache is about 1.15 times as many as needed just for the storage of the data.

Arunai Engineering College

15. What is meant by address mapping? [NOV/DEC-2016]

Addressing mapping, which is defined as the smallest unit of addressed data (from the persistent store) that can be mapped independently to an area of the virtual address space.

16. What is cache memory? [NOV/DEC-2016]



* Cache memory is a small-sized type of volatile computer memory that provides high-speed data access to a processor and stores frequently used computer programs, applications and data.

17 Distinguish SRAM and DRAM. [Apr/May-18]

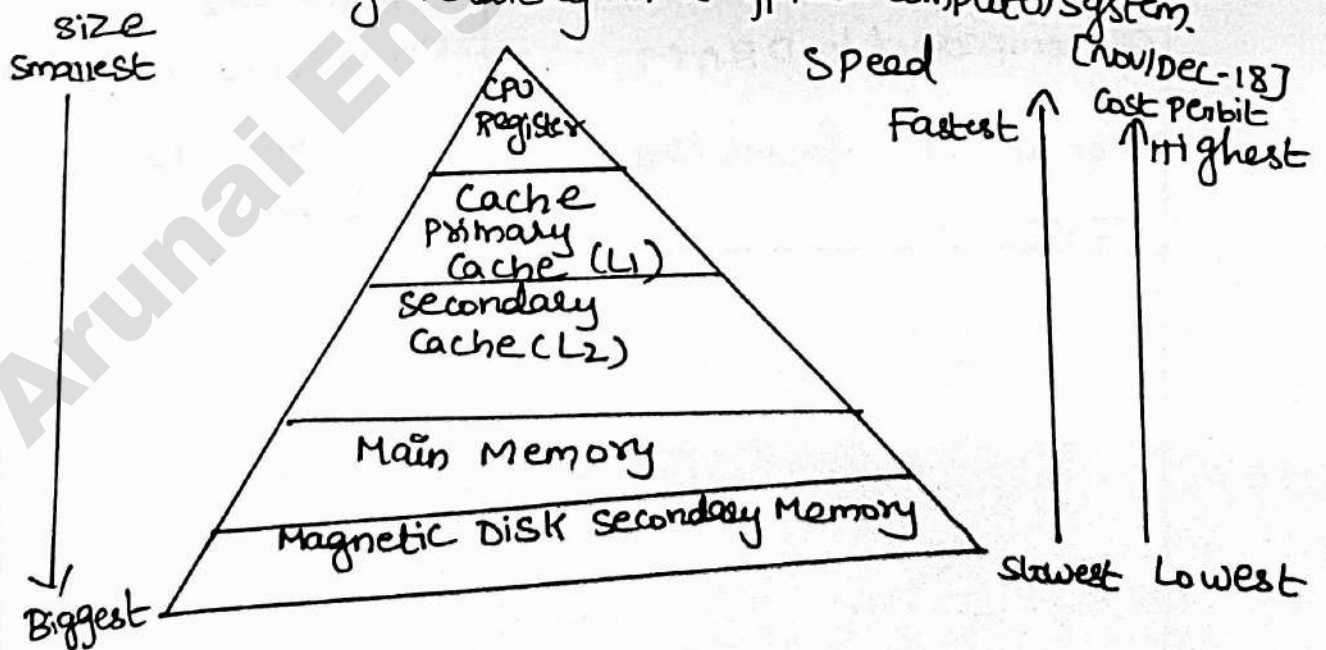
Static RAM (SRAM)	Dynamic RAM (DRAM)
<ul style="list-style-type: none">* Static Random Access Memory uses <u>transistors</u> to store a single bit of data.* SRAM does not need periodic refreshment to maintain data.* SRAM are used in Cache memory* SRAM are expensive as compared to DRAM* SRAM are faster than DRAM	<ul style="list-style-type: none">* DRAM Dynamic Random Access memory uses a separate <u>capacitor</u> to store each bit of data.* DRAM needs periodic refreshment to maintain data* DRAM are used in main memory* DRAM are less expensive as compared to SRAM* DRAM are slower than SRAM

18 What is the use of DMA Controller? [Apr/May-18]

* Direct Memory Access is used to transfer large blocks of data at high speed, a special unit is provided to allow transfer of a block of data directly between an external device and the main memory, without continuous intervention by the processor.

* The process is managed by a chip known as a DMA controller (DMAC)

19 Draw the memory hierarchy in a typical computer system.



20 What is meant by memory-mapped I/O? [Nov/Dec-18]

* If a part of the main-memory address space is assigned to I/O ports, then such systems are called as memory-mapped I/O systems

* A memory-referencing instruction becomes I/O instructions automatically, if the address is address of an I/O port.

* The usual memory load and store instructions are used to transfer data words to or from I/O ports. So, no special I/O instructions are needed.

PART-B

Discuss DMA Controller with block Diagram [Nov/Dec 2016]

Draw the typical block diagram of a DMA controller and explain how it is used for direct data transfer between memory and peripherals? [Nov/Dec-2015]

Explain about DMA Controller with the help of a block diagram [May/June-2016] [·

What is meant by DMA? Explain the use of DMA Controllers in a DMA controller: Computer System [May/June-17]

↳ The Programmed I/O Method has two limitations

1. The speed with which the CPU can test and service I/O device limits I/O data transfer rates.

2. The time that the CPU spends testing I/O device status and executing I/O data transfers can often be better spent on other tasks.

↳ The influence of the CPU on I/O transfer rates is two-fold.

1. A delay occurs while an I/O device needing service waits to be tested by the CPU. If there are many I/O devices in the system then each device may be tested infrequently.

2. Programmed I/O transmits data through the CPU rather than allowing it to be passed directly from main memory to the I/O device and vice versa.

*DMA and interrupt circuits use special control lines it has generic names as DMA REQUEST and INTERRUPT REQUEST, these two lines are connected the I/O devices to the CPU

*DMA and interrupt circuits increase the speed of I/O operations by eliminating most of the role played by the CPU.

*Signals on these lines cause the CPU to suspend its current activities at appropriate breakpoints and attend to the DMA or interrupt request.

*These special request lines eliminate the need for the CPU to execute routines that determine I/O device status.

*DMA also allow I/O data transfers to take place without the execution of I/O instructions by the CPU.

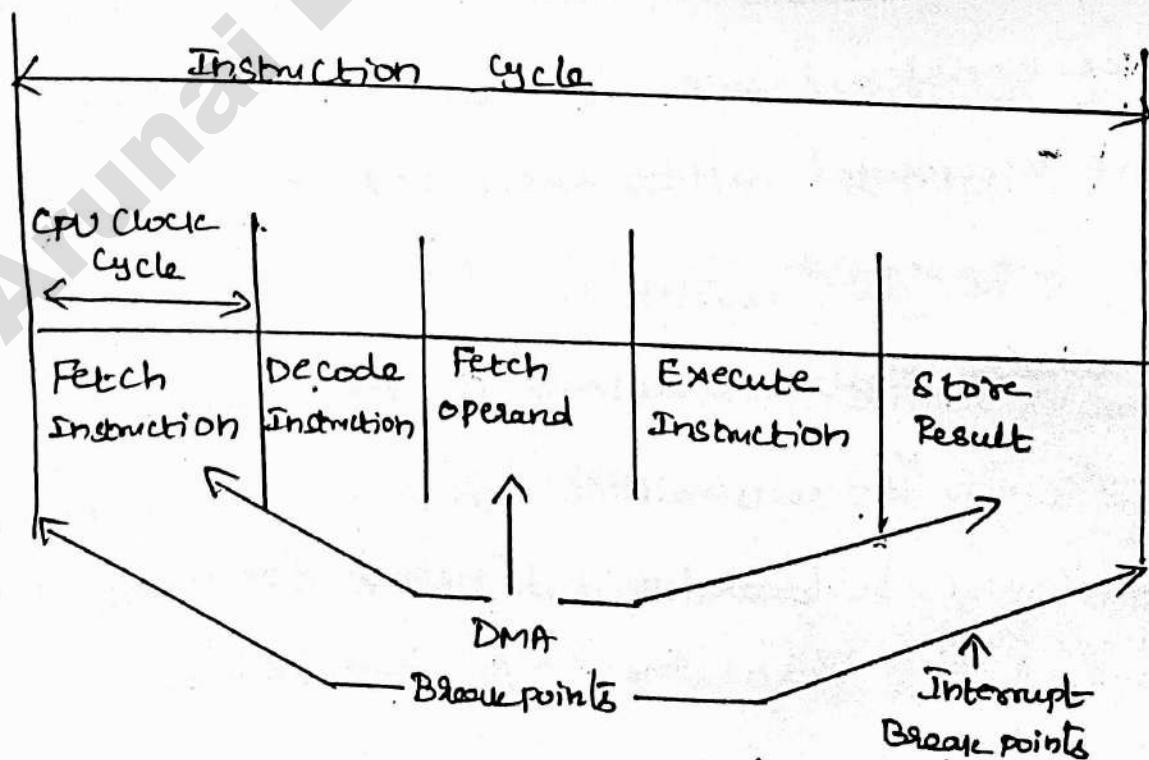
*DMA request by an I/O device only requires the CPU to grant control of the memory (system) bus to the requesting device

The CPU can yield control at the end of any transactions involving the use of system bus.

Typical sequence of CPU actions during execution of a single instruction.

↳ The instruction cycle is collection of number of CPU cycles and several cycles are required use of the system bus. During the instruction cycle there are five break points when the CPU can respond to a DMA request.

↳ When such a request is received by the CPU, it waits until the next breakpoint, release the system bus and signals the requesting I/O device by activating a DMA ACKNOWLEDGE control line.



Direct Memory Access (DMA)

The hardware needed to implement DMA

* The circuit assumes that all access to main memory is via a shared system bus.

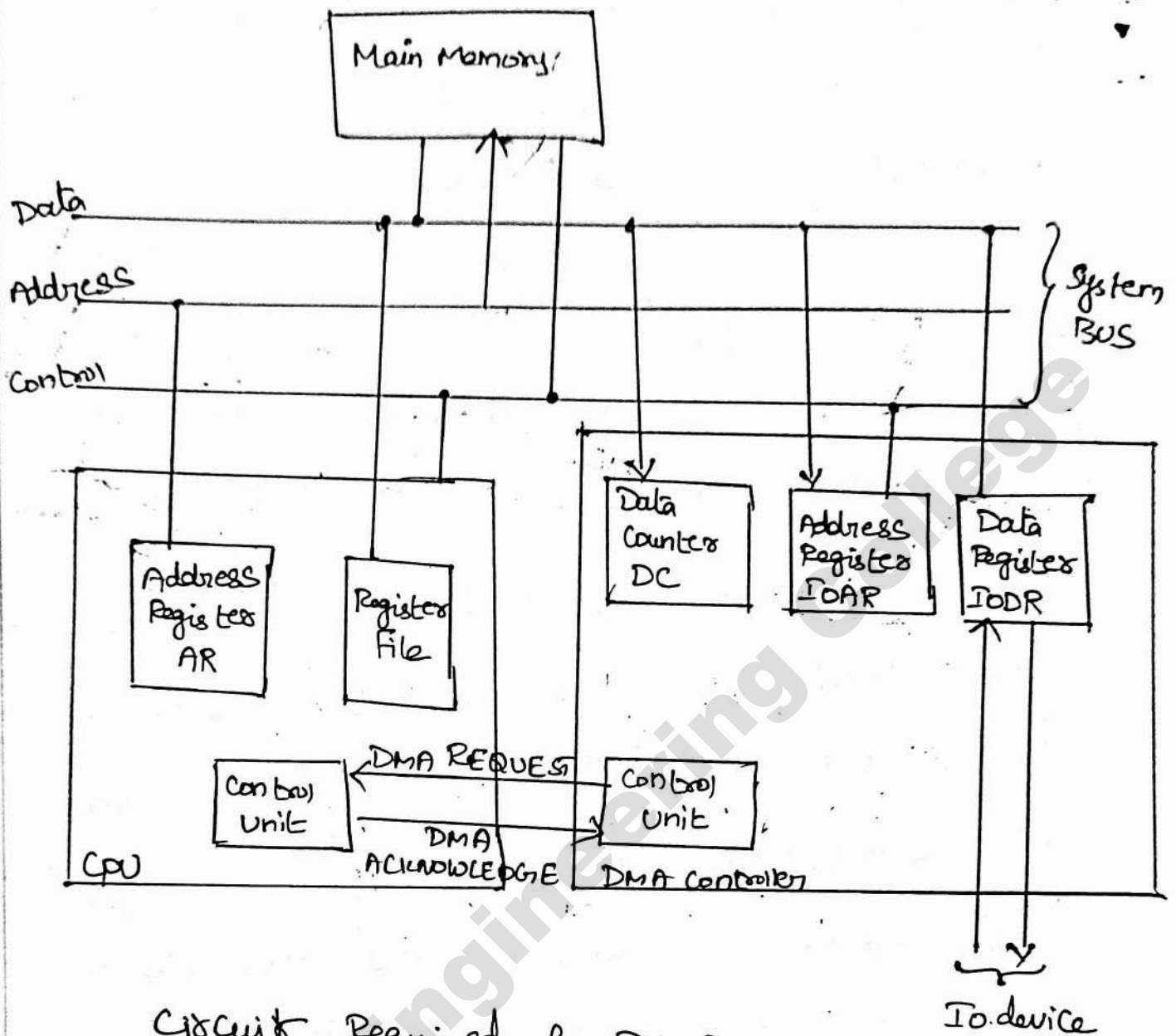
* The I/O device is connected to the system bus via a special interface circuit called DMA Controller. It contains a data buffer register IOBR, it also controls an address register IOAR and data count Register DC.

* These registers enable the DMA controller to transfer data to or from a contiguous region of memory.

* IOAR stores the address of the next word to be transferred. It is automatically incremented or decremented after each word transfer.

* The data counter DC stores the number of words that remain to be transferred.

It is automatically incremented or decremented after each transfer and tested for zero, when DC reaches zero the DMA transfer halts.



Circuit Required for DMA

* The DMA Controller is normally provided with interrupt capabilities to send an interrupt to the CPU to signal the end of the IO data transfer.

* The logic needed to control DMA can be placed in a single IC with other IO control circuits.

* A DMA Controller can be designed to supervise DMA transfers involving several IO devices, each

5-5

device has a different priority of access to the system bus.

* Under DMA control data can be transferred in several different ways. In a DMA block transfer a data word is a sequence of arbitrary length that is transferred in a single burst while the DMA controller is master of the memory bus.

* This DMA mode needs the secondary memories like disk drives, where data transmission cannot be stopped or slowed without loss of data.

* Block DMA transfer supports the fastest I/O data transfer rates but it can make the CPU inactive for relatively long periods by tying up the system bus.

* An alternative technique called cycle stealing, it allows the DMA controller to use the system bus to transfer one word data after it must return control of the bus to the CPU.

* Cycle Stealing reduces the maximum I/O transfer rate and also reduces the interference by the DMA Controller in the CPU's memory access.

* It is possible to eliminate this interference completely by designing the DMA interfaces.

* CPU is not actually using the system bus. This is called transparent DMA.

DMA circuit has the following DMA transfer process:

1. The CPU executes two I/O instructions, which load the DMA registers IOAR and IC with their initial values

a) IOAR register should contain the base address of the memory region to be used in the data transfer.

b) IC register should contain the numbers of words to be transferred to or from the region.

2. When the DMA controller is ready to transmit or receive data, it activates the DMA REQUEST line of the CPU.

a) The CPU waits for the next DMA breakpoint

5-4

b) DMA controller relinquishes control of the data and address lines and activates DMA ACKNOWLEDGE

c) DMA REQUEST and DMA ACKNOWLEDGE are essentially used in BUS REQUEST and BUS GRANT lines for control of the system bus

d) DMA requests from several DMA controllers are resolved by bus priority control techniques.

3. Now DMA controller transfers data directly to or from main memory. After transferring a word, ICAR and DC registers are updated.

4. If DC has not yet reached zero but the I/O device is not ready to send or receive the next batch of data then the DMA controller will release the system bus to the CPU by deactivating the DMA REQUEST line

a) The CPU responds by deactivating DMA ACKNOWLEDGE and resuming control of the system bus.

5. If DC register is decremented to zero then the DMA controller again relinquishes control of the system bus. It may also send an interrupt request signal to the CPU. The CPU responds by halting the I/O device or by initiating a new DMA transfer

* DMA can be obtained under a general method for system

2) What is virtual memory? Explain the steps involved in virtual memory address translation. [MAY-2015]

What is virtual memory? Explain in detail about how virtual memory is implemented with neat diagram [NOV/DEC-2015]

Virtual Memory: (Storage allocation scheme) [NOV/DEC-2016]

* Virtual Memory is a technique that uses ^{May-17} main memory as a "cache" for secondary storage.

Motivations for virtual memory are

Motivations for virtual memory are

1. TO allow efficient and safe sharing of memory among multiple programs.

2. TO remove the programming burdens of a small, limited amount of main memory.

* TO allow multiple virtual machines to share the same memory we must protect the virtual machines from each other.

* We need to ensure that a program can only read and write the portions of main memory that has been assigned to virtual machines.

* Main memory need to contain only the active portions of the many virtual machines and virtual memory allow us to efficiently share the processor as well as the main memory.

* We cannot know which virtual machines will share the memory with other virtual machines when we compile them.

* Because the virtual machines sharing the memory that can change dynamically while the virtual machines are running.

* We need to compile each program into its own address space.

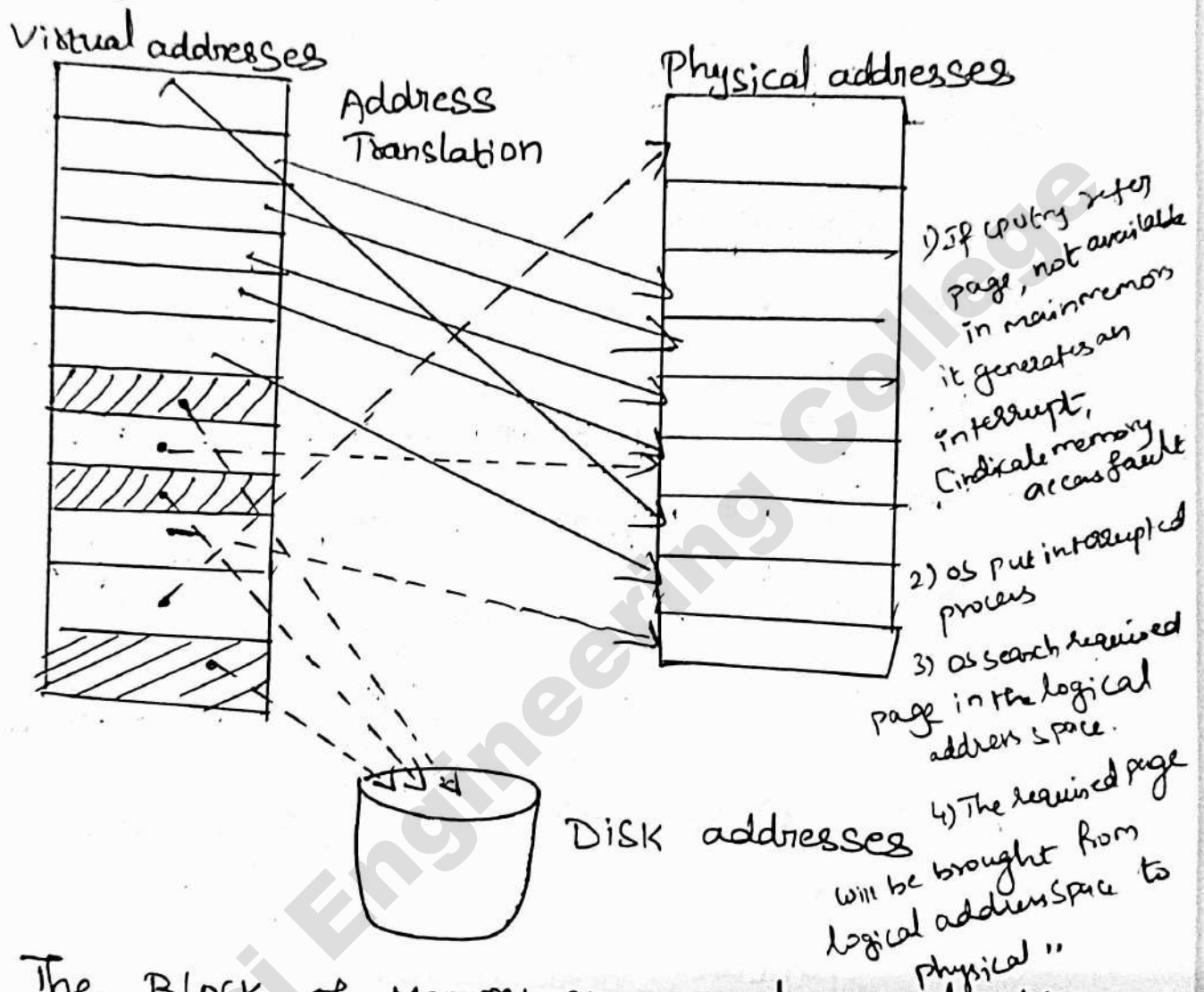
* Virtual memory implements the translation of a program address space to physical address.

* This translation process enforces protection of a program's address space from other virtual machines.

* In a virtual memory block is called as page and a virtual memory miss is called as page fault. With virtual memory, the processor produces a virtual address.

* Virtual address is an address that corresponds to a location in virtual space and is translated by address mapping to a physical address when the memory is accessed.

* The virtually addressed memory with pages mapped to main memory. This process is called address mapping or address translation.



The Block of Memory are mapped virtual addresses to physical addresses.

* The processor generates virtual addresses while the memory is accessed using physical addresses. Both virtual memory and physical memory are broken into Pages.

5) page replacement algm used for decision making of replacing the page in physical address
 ↳ signal will be sent to the CPU to continue the program

* Not necessary all pages or segments are present in main memory. Required pages need to be loaded into memory, virtual memory is implemented by Demand Paging or Demand Segmentation.

* A virtual page is mapped to physical page. It is also possible for a virtual page to be absent from main memory and not be mapped to a physical address, in the case the page resides on disk.

* Virtual memory can be used to control two memory hierarchy levels such as

1. DRAM's and flash memory in personal mobile devices
2. DRAM's and magnetic disks in servers.

* Virtual memory also simplifies loading the program for execution by providing relocation method

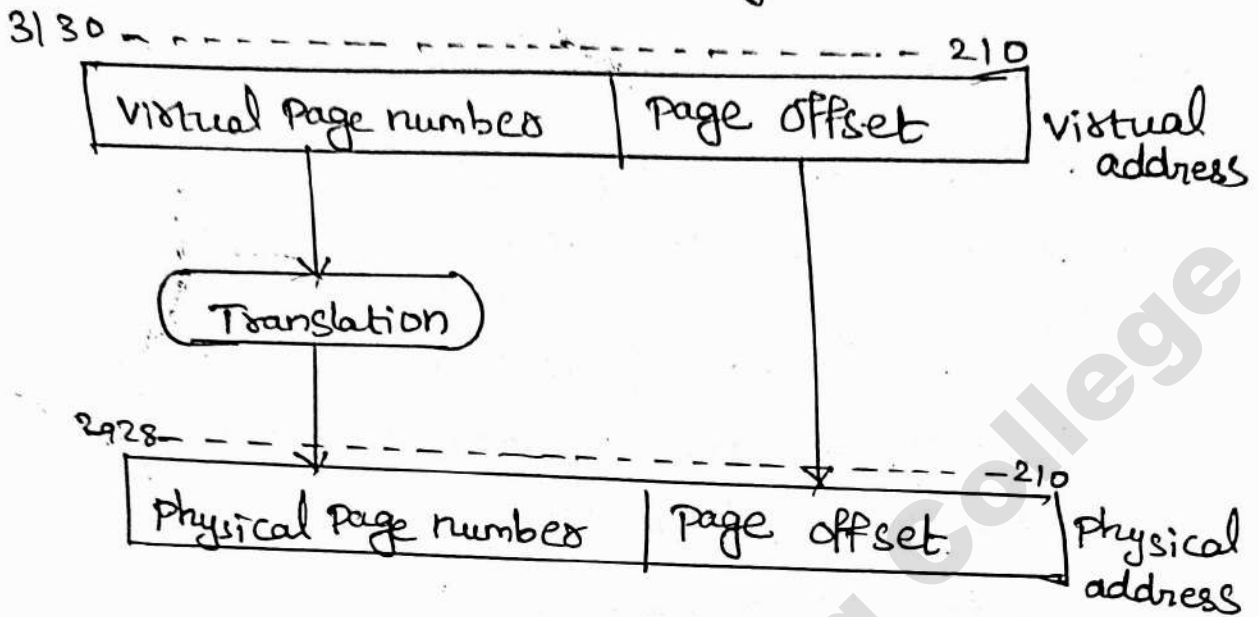
* In virtual memory, the address is broken into two parts

1) virtual page number

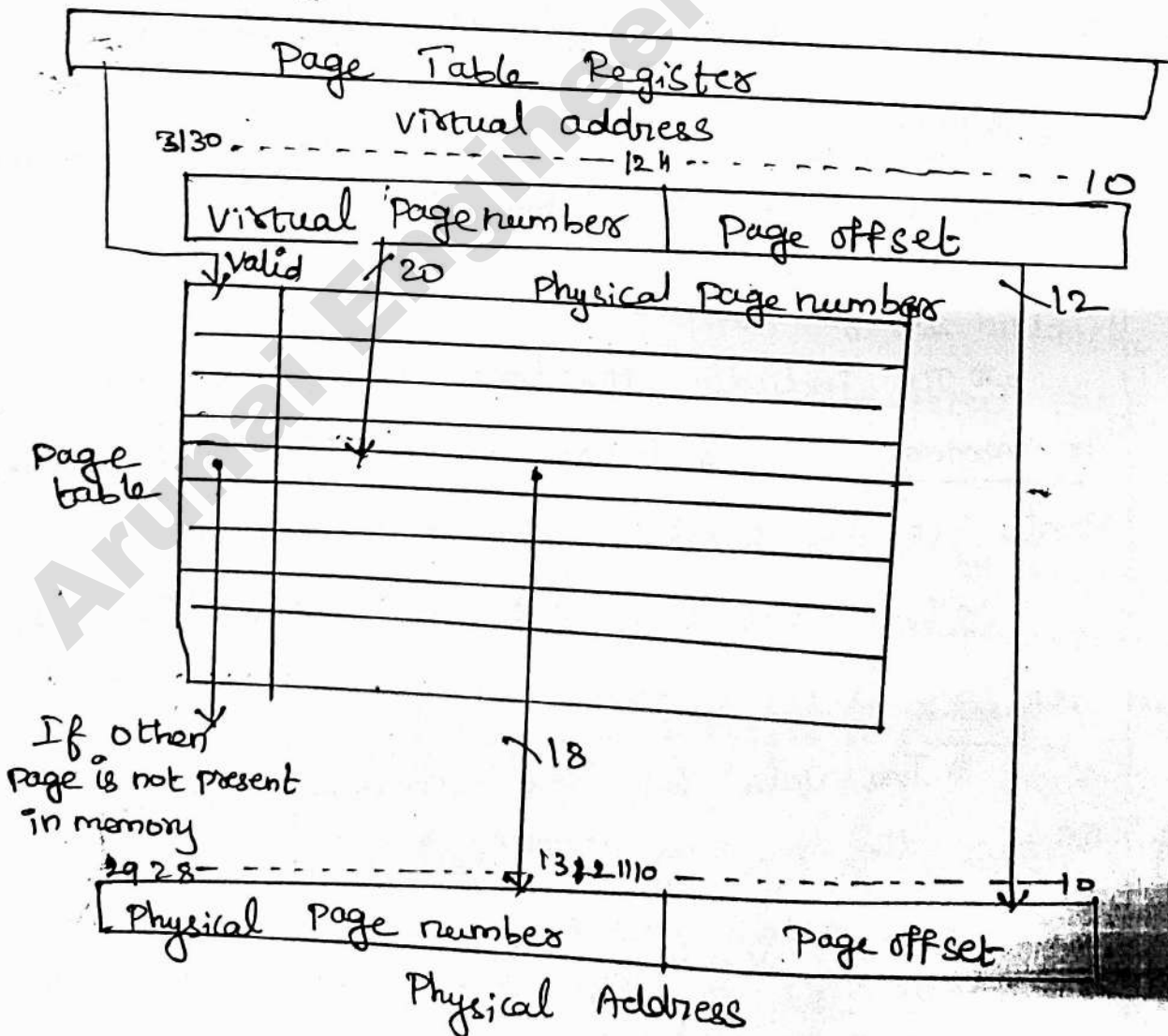
2) page offset

* The number of bits in the page offset field determines the page size and the number of pages addressable with the virtual address need not match the number of pages addressable with the physical address.

Mapping from a virtual to physical address



Placing a page:



* To reduce the page fault the designer must optimize page placement in more attractive way. If virtual page can be mapped to any physical page then when a page fault occurs the operating system can choose to replace any page it wants.

* In virtual memory systems, we can locate pages by using a table that indexes the memory and this structure is called a Page table.

* Page table resides in memory and its indexed with the page number from the virtual address to discover the corresponding physical page numbers.

* Each program has its own page table, which maps the virtual address space of that program to main memory.

* To indicate the location of the page table in memory, the hardware includes a register that points to the start of the page table.

* These registers are called as Page table register.

* The valid bit for each entry indicates whether the mapping is legal or not.

↳ If valid bit is off. - Then the page is not present in memory.

↳ If valid bit is on - The page is in memory.

* The Entry contains the physical page number because the page table contains a mapping for every possible virtual page.

Page Faults:

* If the valid bit for a virtual page is off it causes a page fault. If page fault occurs then the control has given to operating system.

* Once the operating system gets control, it must find the page in the next level of the hierarchy and decide where to place the requested page in main memory.

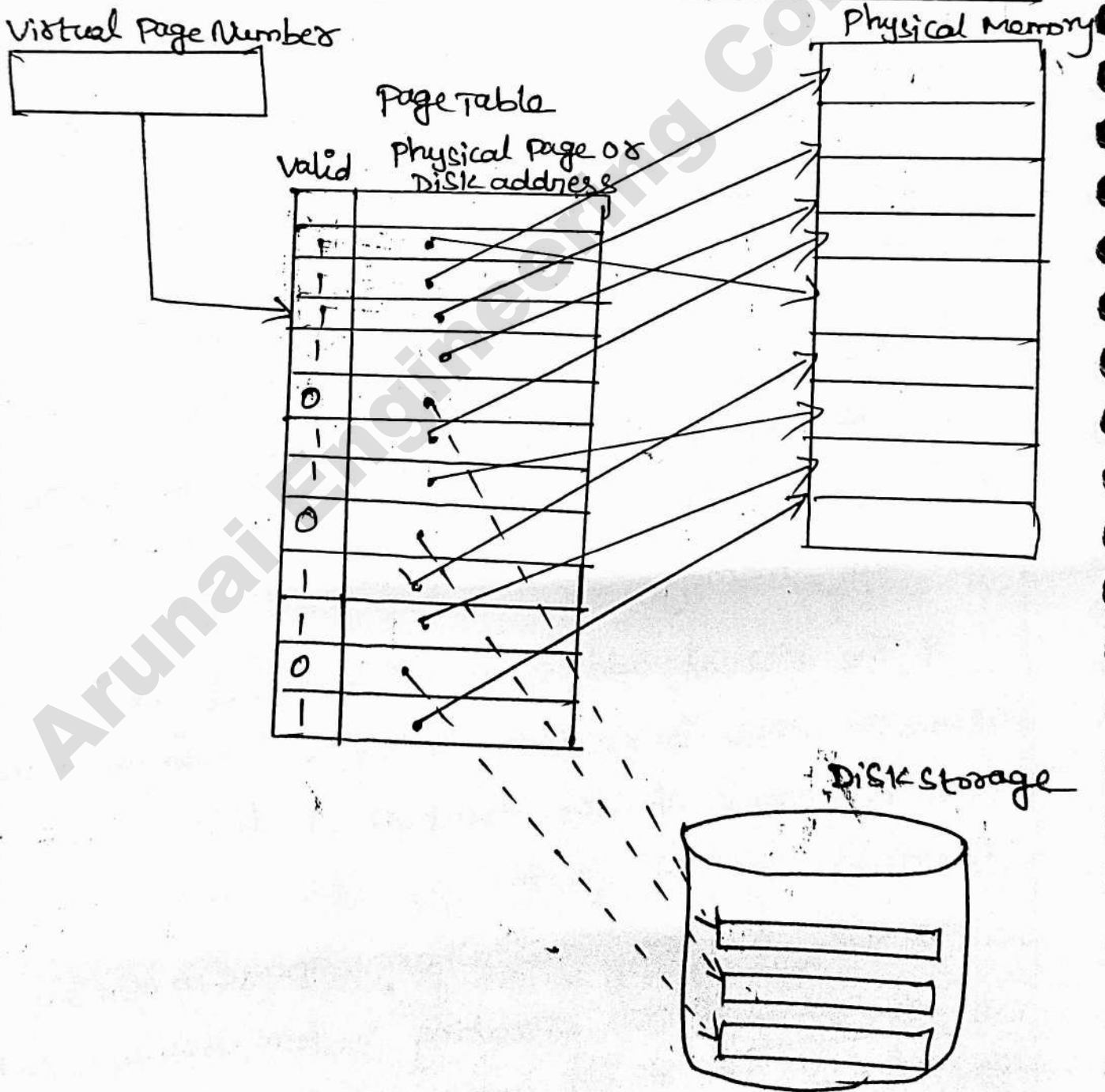
* The virtual address alone does not tell us where the page is on disk, so virtual memory system, must keep track of the location on disk of each page in virtual address space.

* We do not know when a page in memory will be replaced. So operating system usually creates the space on flash memory or disk for all the

Pages of a process when it creates the process. This space is called the swap space.

Swap space - is a space on the disk or flash memory reserved for the full virtual memory space of a process.

Page table maps each page in virtual memory to either a page in main memory or a page stored on disk



* The virtual page number is used to index the page table. If the valid bit is on, the page table supplies the physical page number corresponding to the virtual page.

* If the valid bit is off, the page currently resides only on disk at a specified disk address. Pages in the main memory and the pages on disk are the same size.

* Operating system also creates a data structure that tracks which processes and which virtual addresses use each physical page.

* When a fault occurs, if all the pages in main memory are in use then the operating system must choose a page to replace.

* OS follow LRU (Least Recently Used) replacement scheme. The replaced pages are written to swap space on the disk.

Writes in virtual memory system:

Virtual memory system must use write back it performing the individual writes into the page in memory and copying the page back to disk when it is replaced in the memory.

Advantage:

* write back scheme has major advantage in virtual memory system.

* Because the disk transfer time is small compared with its access time and copying back an entire page is much more efficient than writing individual words back to the disk.

Disadvantage:

* It is costly

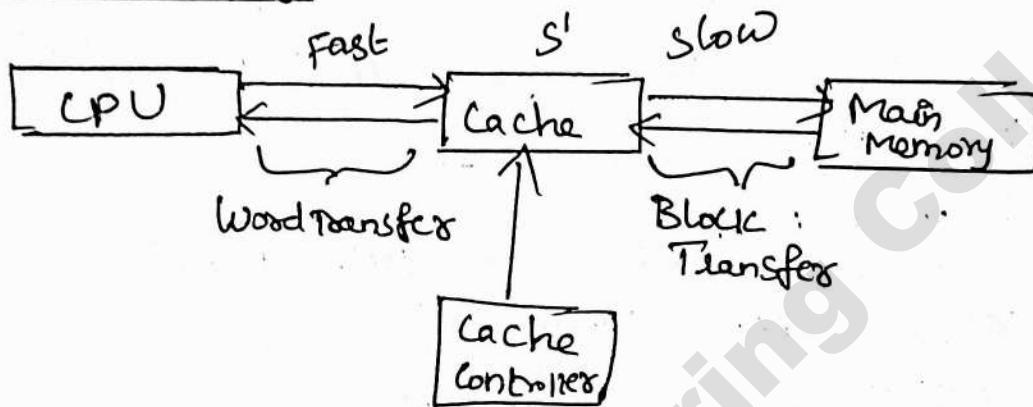
* Dirty bit is set when any word in a page is written.

* If the OS chooses to replace the page, the dirty bit indicates whether the page needs to be written out before its location in memory can be given to another page.

3. Define Cache Memory? Explain the various Mapping Techniques associated with each memories [May/June-2016]

Explain mapping functions in cache memory to determine how memory blocks are placed in cache [Nov/Dec-2014]
[Apr/May-17]

Cache Memory:



Cache memory System

* The part of program (code) and data that work at a particular time is usually accessed from the SRAM memory. This is accomplished by loading the active part of code and data from main memory to SRAM memory.

* This small section of SRAM memory added between processor and main memory to speed up execution process is called Cache memory

Mapping Techniques:

* Usually, the cache memory can store a reasonable number of memory blocks at any given time, but this number is small compared to the total number of blocks in the main memory.

* The correspondence between the main memory blocks and those in the cache is specified by a mapping function.

The mapping techniques are classified as

1. Direct-mapping technique
2. Associative-mapping technique.

> Fully-associative

> Set-associative techniques

To discuss these techniques of cache mapping we consider a cache consists of 128 blocks of 16 words each, for a total of 2048 (2K) words and assume that the main memory has 64K words. This 64K words of main memory is addressable by a 16-bit address and it can be viewed as 4K blocks of 16 words each. The group of 128 blocks of 16 words each in main memory form a page.

1. Direct-Mapping:

* It is the simplest mapping techniques.

* In this technique, each block from the main memory has only one possible location in the cache organization.

Eg: The block i of the main memory maps onto the block j ($j = i \text{ modulo } 128$) of the cache.

\therefore one of the main memory blocks $0, 128, 256, \dots$ is loaded in the cache, it is stored in cache block 0.

Blocks $1, 129, 257, \dots$ are stored in cache block 1

In general, the mapping expression is

$$j = i \text{ modulo } m$$

where $i =$ main memory block number

$j =$ cache block (line) number

$m =$ number of blocks (lines) in the cache

The address is divided into three fields

i) word field

ii) Block field

iii) Tag field

Tag field: - Identifies which memory block is in the slot

* It is used to store the high-order 5-bits of memory address of the block. These 5-bit (tag bits) are used to identify which of the 32 blocks (pages) that are mapped into the cache.

↳ When memory is accessed, the 7-bit cache block field of each address generated by CPU points to a particular block location in the cache.

↳ The high-order 5-bits of the address are compared with the tag bits associated with that cache location.

↳ If they match, then the desired word is in that block of the cache.

↳ If there is no match, then the block containing the required word must first be read from the main memory and loaded into the cache.

Q. Associative-Mapping (Fully-Associative Mapping) -

* In this technique, a main memory block can be placed into any cache block position.

* There is no fix block, the memory address has only two fields

i) word ii) Tag

* This technique is also referred to as fully-associative mem...

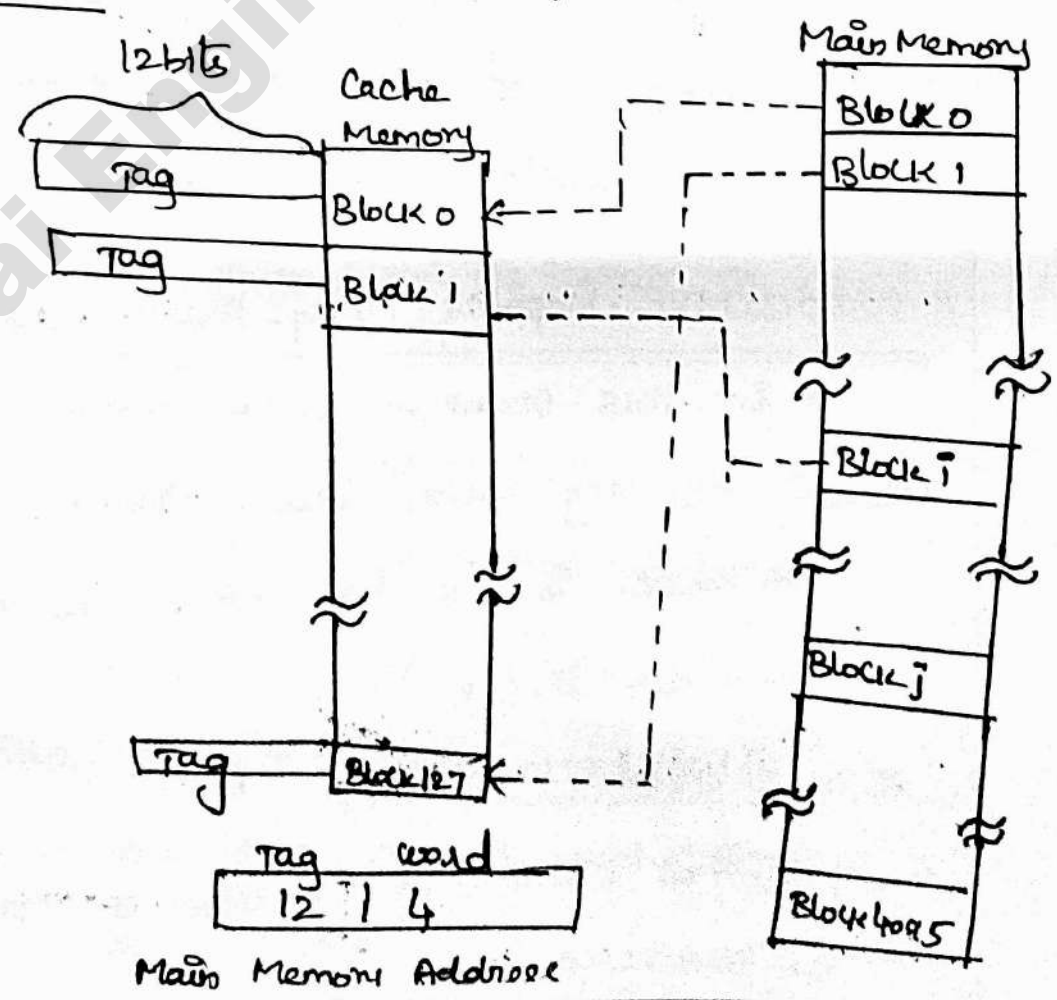
* The 12-tag bits are required to identify a memory block when it is resident in the cache.

* The high-order 12-bits of an address received from the CPU are compared to the tag bits of each block of the cache to see if the desired block is present

* once the desired block is present, the 4-bit word is used to identify the necessary word from the cache.

* This technique gives complete freedom in choosing the cache location in which to place the memory block. Thus, the memory space in the cache can be used more efficiently.

* A new block that has to be loaded into the cache has to replace (or remove) an existing block only if the cache is full.



Disadvantage:

↳ In associative-mapped cache, it is necessary to compare the high-order bits of address of the main memory with all 128 tag corresponding to each block to determine whether a given block is in the cache.

3) Set-Associative Mapping: - Cache is divided into number of sets
Each set contains a number of lines
* The Set-associative mapping is a both direct and associative mapping.

* It contains several groups of direct-mapped blocks that operate as several direct-mapped caches in parallel.

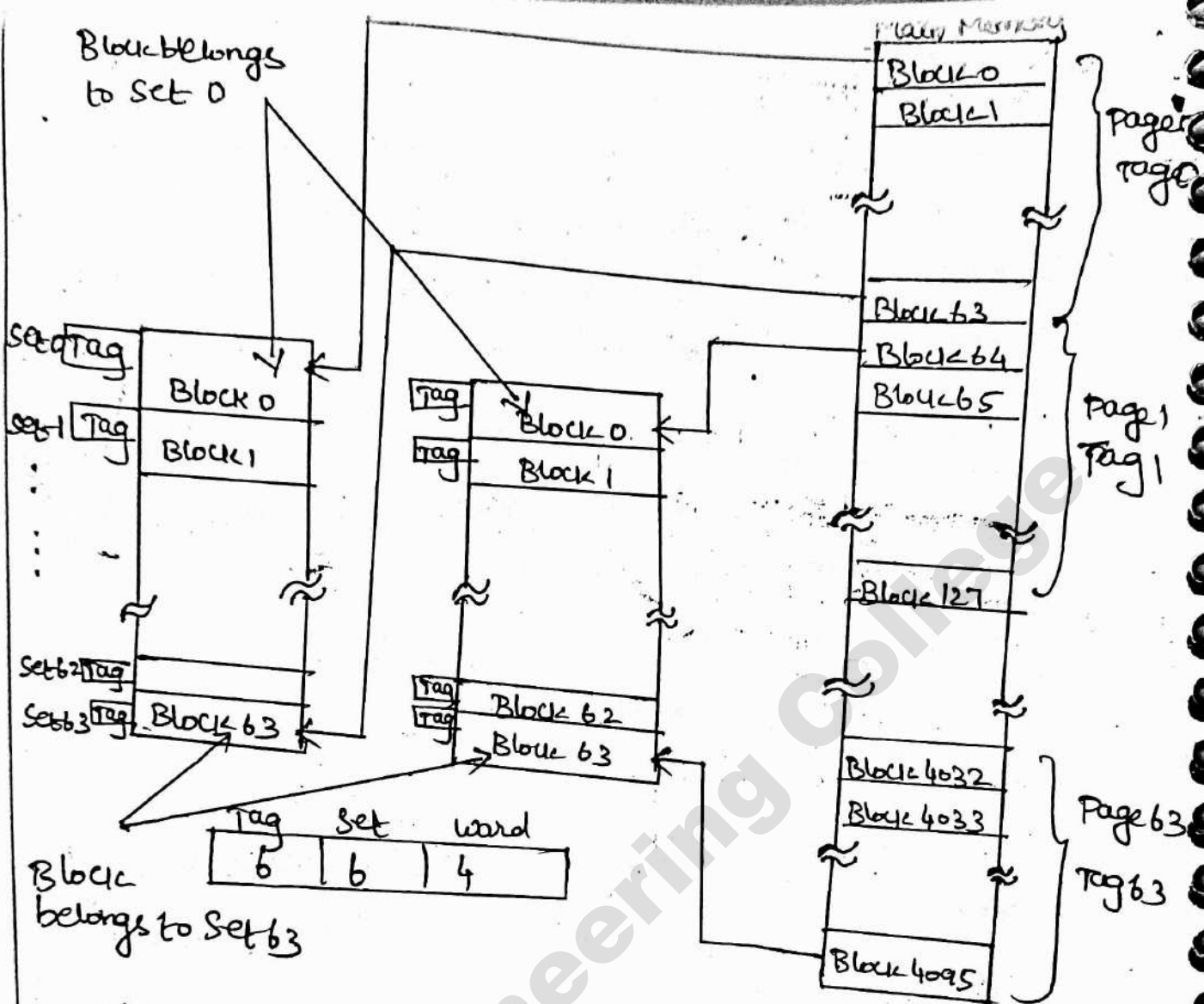
* A block of data from any page in the main memory can go into a particular block location of any direct-mapped cache.

* The required address comparisons depend on the number of direct-mapped caches in the cache systems.

Two-way Set-associative Cache

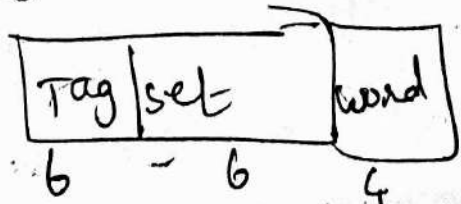
* Each page in the main memory is organized in such a way that the size of each page is same as the size of one directly mapped cache.

* It is called two-way set-associative cache because each block from main memory has two choices for block placements.



* In this technique, block 0, 64, 128, ... 4032 of main memory can map into any of the two (block 0) blocks of set 0, block 1, 65, 129, ... 4033 of main memory can map into any of the two (block 1) blocks of set 1 and so on.

Memory address



Tag: $2^6 = 64$ different memory lines in each of the $2^6 = 64$ different set values. When the number of lines per set is n , the mapping is called n -way associative.

word (4-bits):- The lower order 4-bits are used to select one of the 16 words in a block.

Set (6-bits):- The 6-bits set field determines which set the cache contains the desired block.

Direct Mapping:

Advantage:

- * It is simple and easy to implement.

Disadvantage:

- * It is not flexible.
- * Since more than one memory block is mapped onto a given cache block position, contention may arise for that position even when the cache is not full.

Associative mapping:

Advantages:

- * It gives complete freedom in choosing a cache location to load a memory block in cache.
- * Space in the cache can be used efficiently.
- * A new block replaces an existing block only if the cache is full.

Dis advantages:

- * Cost of associative mapping cache is higher than the cost of direct mapping cache.
- * All the 128 tags are searched to find whether a given block is in the cache or not

Set - Associative Mapping:

Advantages:

- * The contention method of direct mapping technique is solved by having a few choices for block replacement.
- * The hardware cost is reduced by decreasing the size of associative search.
- * It is simple to implement.

Explain in detail about the Bus arbitration techniques

In DMA [NOV/DEC-2014] [May-17]

A bus is common pathway connecting two or more bus devices
CPU, I/O, Main memory are interconnected by common bus
← address bus
← control bus
← data bus

Bus Arbitration: → Deciding which devices get to use the shared bus data

* The device that is allowed to initiate data transfers on the bus at any given time is called the busmaster.

4) Two parts
← Issuing Command Bus Master
← Transferring the data Bus Slave

* Bus arbitration is the process by which the next device to become the bus master is selected and bus mastership is transferred to it. The selection of bus master is usually done on the priority basis.

Approaches to bus arbitration:

1) In Centralized bus arbitration, a single bus arbiter performs the required arbitration. The bus arbiter may be the processor or a separate controller connected to the bus.

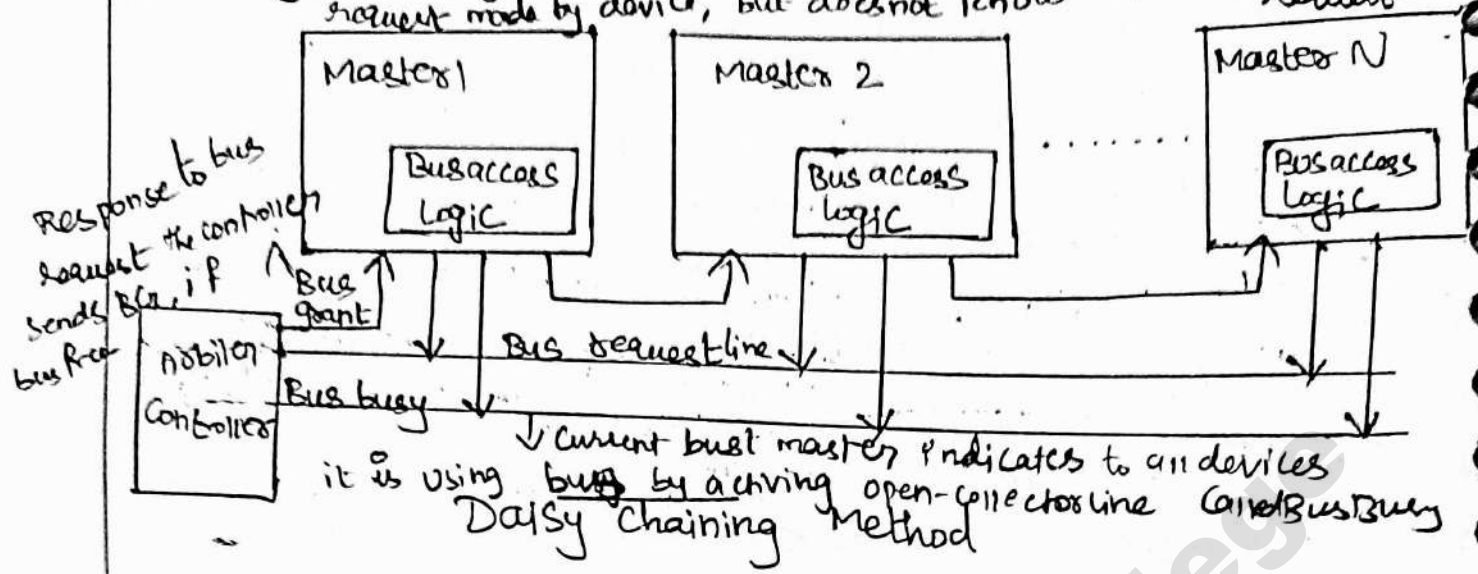
2) In Distributed bus arbitration, all devices participate in the selection of the next bus master.

1) Centralized Bus Arbitration Schemes: 1) higher reliability
2) operation of bus is not dependent on any single device

There are three different arbitration schemes that use the centralized bus arbitration approach.

- a) Daisy chaining b) Polling method c) Independent request

a) Daisy Chaining! BRL wired OR line, controller only request made by device, but does not know which device made the request



* It is a simple and cheaper method

* All masters make use of the same line for bus request.

* In response to a bus request the controller sends a bus grant if the bus is free.

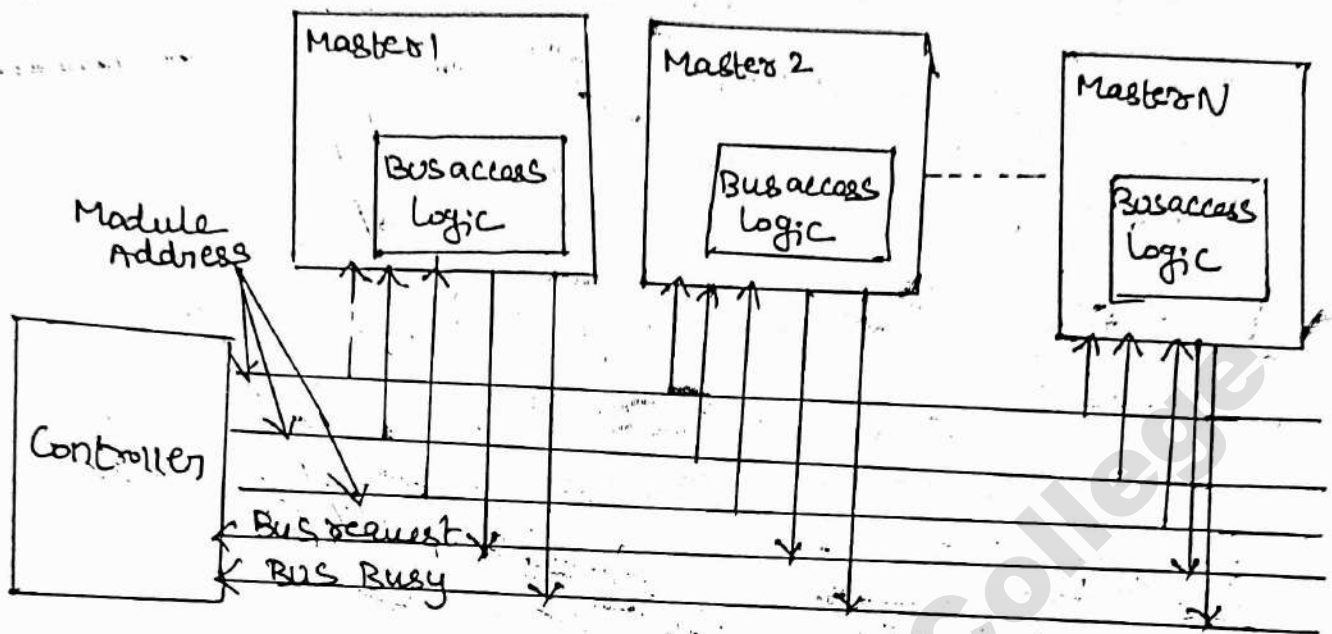
Advantages:

1. It is a simple and cheaper method
2. It requires the least number of lines and this number is independent of the number of masters in the system

Disadvantages:

1. The priority of the master is fixed by its physical location
2. Failure of any one master causes the whole system to fail.

b) Polling Method: - All masters use same line for Bus request.



* In this the Controller ^{assign addresses for master} is used to generate the addresses for the masters. (All master connected a local address bus).

* Number of address lines required depends on the number of masters connected in the system.

Eg 8 masters connected in the system, at least three address lines are required.

Advantages: * In response to bus request, the controller places the address of the bus master on address bus.

1. The Priority can be changed by altering the Polling Sequence stored in the controller.

2. If the one module fails entire system does not fail.

c) Independent request:

In this scheme each master has a separate pair of bus request and bus grant lines and

each pair has a priority assigned to it.

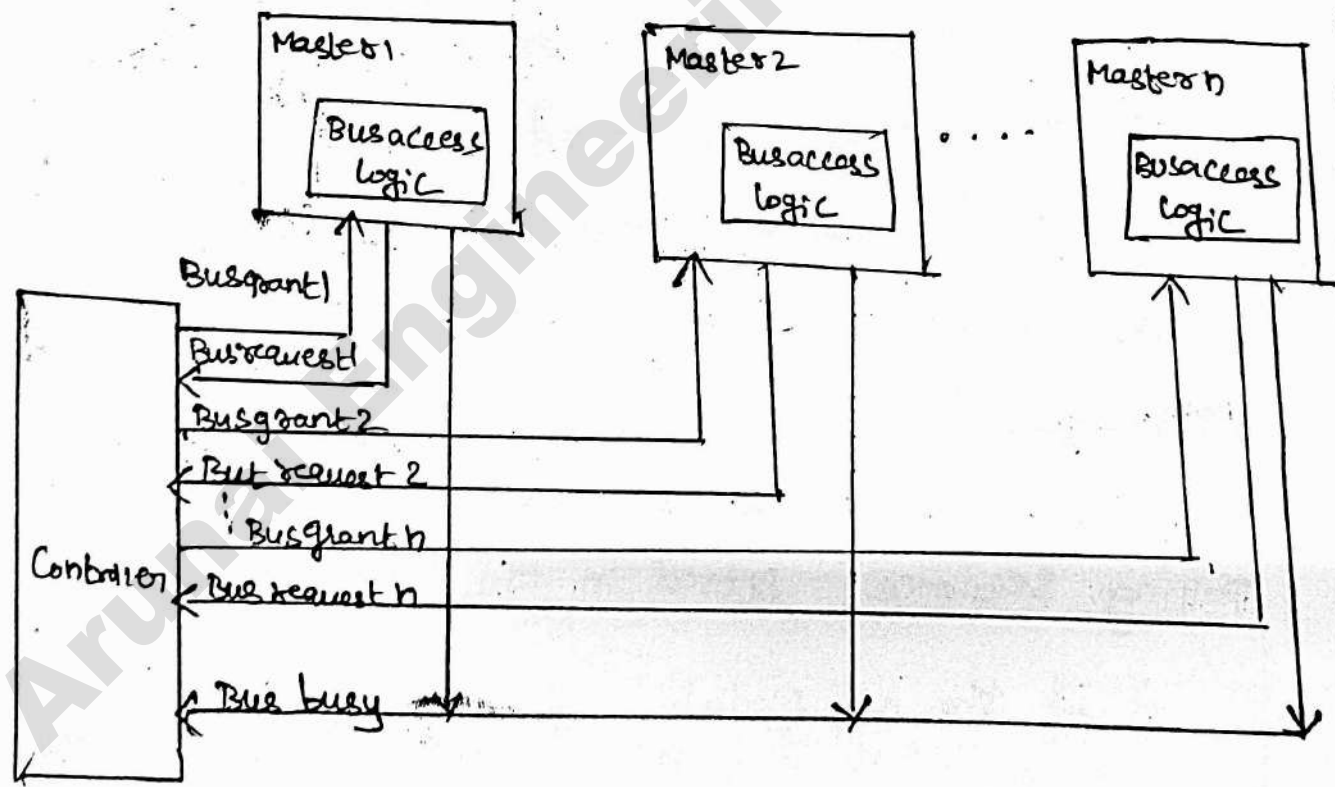
* The built in priority decoder within the controller selects the highest priority request and asserts the corresponding bus grant signal.

Advantage:

1. Due to separate pairs of bus request and bus grant signals, arbitration is fast and is independent of the number of masters in the system.

Disadvantage:

1. It requires more bus request and grant signals $(2 \times n)$ signals for n modules.



Independent request method

5 Elaborate on the various Memory technologies and its relevance. [NPRIMAY-2015]

Memory Technologies

There are four primary Technologies used in memory hierarchies such as

1. SRAM technologies
2. DRAM technologies
3. Flash Memory technologies
4. Disk memory technologies

1. SRAM Technology. low density, high power, expensive fast

* Static Random Access memory is a part of the Random Access Memory (RAM). It is main memory and located very close to the processor.

* SRAM'S are simply integrated circuits that are memory array with a single access port and it can provide either a read or a write.

* It has a fixed access time to any data through the read and write access times may differ.

* SRAM'S don't need to refresh so the access time is very close to the cycle time.

* Transistors hold the data as long as the power supply is not cut off

It hold only one bit data
* It typically use six to eight transistors per bit to prevent the information from being disturbed when read.

* It need only minimal power to obtain the charge in standby mode. * Address not divided
* use for caches

Fault
Accord

* MOST PCs and server systems used separate SRAM chips for either their primary, secondary or even caches

2-DRAM Technology: - high density, low power, cheap, slow

* In a SRAM, as long as power is applied the value can be kept indefinitely.

* In Dynamic RAM (DRAM) the value kept in a cell is stored as a charge in a capacitor

* DRAM has single transistor used to access this stored charge either to read the value or to overwrite the charge stored there.

Big
memory

* Because it uses only a single transistor per bit of storage they are much denser and cheaper than SRAM.

* DRAM'S store all the charge on a capacitor
So it cannot be kept indefinitely and must be periodically refreshed.

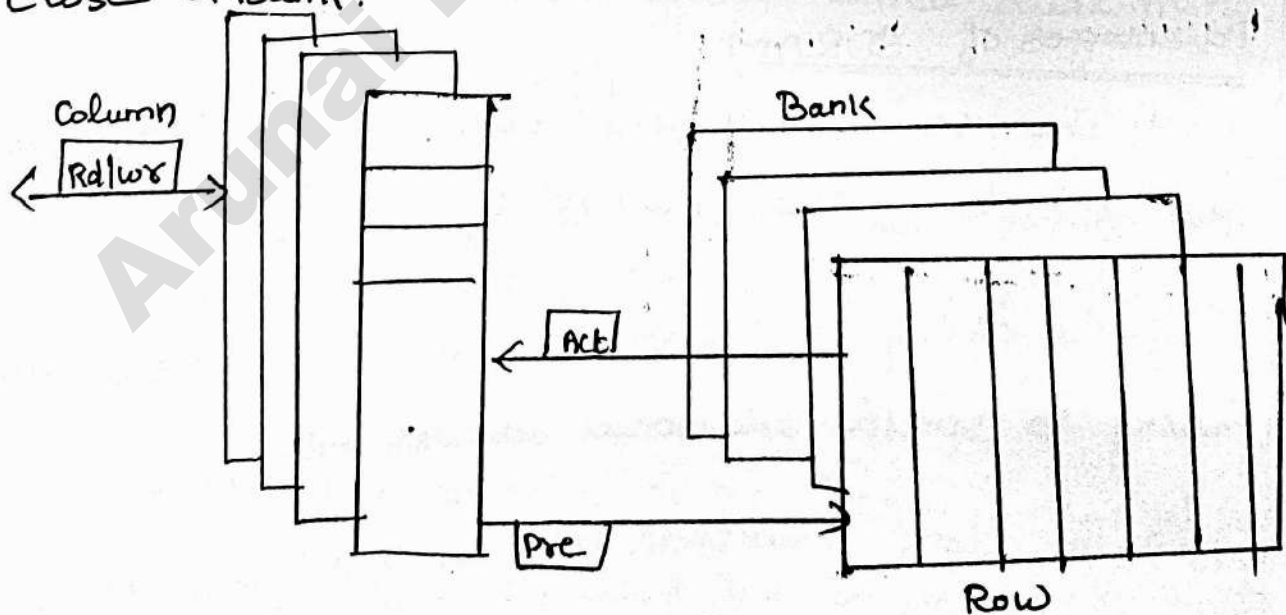
* If every bit had to be read out of the DRAM and then written back individually, so we constantly be refreshing ^{data} the DRAM by leaving no time for accessing it. [Two halves of address Row/Coln]
US for main memory

* DRAM uses a two level decoding structure and it will allow us to refresh the entire row with a read cycle followed immediately by a write cycle

* DRAM'S are organized in banks, typically four for DDR3. Each bank consists of a series of rows.

* sending a PRE (pre charge) command opens or

close a bank.



Internal organization of a DRAM

* A row address is sent with an Act (activate) which causes the row to transfer to a buffer, when the row is in the buffer it can be transferred by successive column address at whatever the width of the DRAM is or by specifying a block transfer and the starting address.

* Each command as well as block transfers are synchronized with a clock.

* To improve performance DRAM's buffer rows for repeated access.

* To improve the interface to processor's DRAM's added clocks with its and are called synchronous DRAM's or SDRAM's - syn. with bus

Advantages of SDRAM

1. The use of a clock eliminates the time for the memory and processor to synchronize.
2. The ability to transfer the bits in the burst without having to specify additional address bits.
3. The clock transfers the successive bits in a burst.

DDR- Double Data Rate

* The fastest version of SDRAM is called Double Data Rate.

* It means data transfers on both the rising and falling edge of the clock, thereby getting twice as much bandwidth as we expected based on the clock rate and the data width.

* The latest version of this technology is called DDR4.

* A DDR4-3200 DRAM can do 3200 millions transfers per second which means it has a 1600 MHz clock.

* Personal mobile devices like the iPad use individual DRAM's and memory for servers are commonly sold on small boards called dual inline memory modules (DIMM's)

* DIMM contains 4-16 DRAM and they are normally organized to be 8 bytes wide for server system.

* DIMM can have so many DRAM chips and only a portion of them are used for a particular transfer.

3. Flash memory Technology:

* Flash memory is a type of Electrically Erasable Programmable Read only Memory (EEPROM)

* The EEPROM technologies the writes can wear out flash memory bits. Most flash products include a controller to spread the writes by remapping blocks that have been written many times to less troubled blocks.

* This technique is called wear leveling

* In Digital cameras, flash memory is used to store picture image data.

* In MP3 players, flash memory is used to store the sound data.

* Flash memory are available in modules. These modules are implemented in two types

i) Flash cards

ii) Flash Drives

4. Disk Memory Technology:

* To read and write information on a hard disk, a movable arm containing a small electromagnetic coil called a read-write head is located just above each surface.

* The entire drive is permanently sealed to control the environment inside the drive, which in turn allows the disk head to be much closer to the drive surface.

* Each disk surface is divided into concentric circle called tracks.

* Each track is in turn divided into sectors that contains the information.

* Each track may have thousands of sectors and it has 512 to 4096 bytes in size.

* Sector is one of the segments that make up a track on a magnetic disk and it is the smallest amount of information that is read or written on a disk.

↳ To access the data, the OS must direct the disk through a three stage process

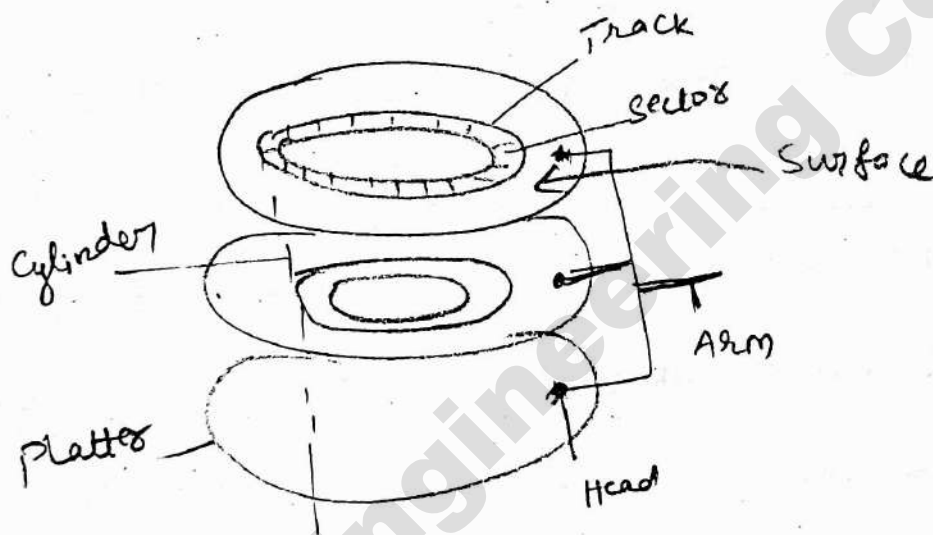
1. To position the head over the proper track. This operation is called seek and time to move the head to the desired track is called seek time.

2. Once the head has reached the correct track, we must wait for the desired sector to rotate.

under the read/write head. This time is called the rotational latency or rotational delay

3. Disk access is transfer time. It is the time to transfer a block of bits.

The transfer time is a function of the sector size, the rotation speed and the recording density of a track.



6. Assume the miss rate of an instruction cache is 2% and the miss rate of the data cache is 4%. If a processor has a CPI of 2 without any memory stalls and the miss penalty is 100 cycles for all misses determine how much faster a processor would run with a perfect cache that never missed. Assume the frequency of all loads and stores is 36%. [NOV/DEC-17]

Solution:

* The number of memory miss cycles for instructions in terms of the instruction count (I) is

$$2\% / 100 = 0.02$$

$$0.02 \times 100 = 2.00$$

$$\text{Instruction miss cycles} = I \times 2\% \times 100 = 2.00 \times I$$

* As the frequency of all loads and stores is 36%, we can find the number of memory miss cycles for data references:

$$\left. \begin{array}{l} 36\% / 100 = 0.36 \\ 4\% / 100 = 0.04 \end{array} \right\} 0.36 \times 0.04 = 0.0144 \times 100 = 1.44$$

$$\text{Data miss cycles} = I \times 36\% \times 4\% \times 100 = 1.44 \times I$$

* The total number of memory stall cycles is

$$2.00I + 1.44I = 3.44I$$

* This is more than three cycles of memory stall per instruction.

* Accordingly, the total CPI including memory stalls is $2 + 3.44 = 5.44$.

* Since there is no change in instruction count or clock rate, the ratio of the CPU execution time is

$$\frac{\text{CPU time with stalls}}{\text{CPU time with perfect cache}} = \frac{I \times \text{CPI}_{\text{stall}} \times \text{clock cycle}}{I \times \text{CPI}_{\text{perfect}} \times \text{clock cycle}}$$
$$= \frac{\text{CPI}_{\text{stall}}}{\text{CPI}_{\text{perfect}}}$$
$$= \frac{5.44}{2}$$

* The performance with the perfect cache is better by $\left. \vphantom{\begin{matrix} * \\ * \end{matrix}} \right\} = \frac{5.44}{2} = 2.72$.

Suppose, speed-up the computer in the example by reducing its CPI from 2 to 1 without changing the clock rate, which might be done with an improved pipeline.

* The system with cache misses would then have a CPI of $1 + 3.44 = 4.44$.

* The system with the perfect cache would be

$$\frac{4.44}{1} = 4.44 \text{ times as fast.}$$

* The amount of execution time spent on memory stalls would have risen from

$$\frac{3.44}{5.44} = 63\%$$

to

$$\frac{3.44}{4.44} = 77\%$$

Similarly, increasing the clock rate without changing the memory system also increases the performance loss due to cache misses.

7 Discuss about Programmed I/Os associated with Computers. [APR/May-18]

Programmed IO:

* If IO operations are completely controlled by the CPU, the computer is said to be using Programmed IO.

* In this case, the CPU executes programs that initiate, direct and terminate the IO operations.

* This type of control can be implemented easily with less hardware.

* It is more useful in small, low-speed systems where hardware cost must be minimized.

* Data transfer is between two programmable registers. One is CPU register and the other is attached to the IO device.

* The IO device does not have direct access to main memory. A data transfer from an IO device to memory requires the CPU to execute several

Instructions. It includes an input instruction to transfer a word from the I/O device to the CPU and a store instruction to transfer the word from CPU to M.

Memory-Mapped IO

* In programmed IO systems, the CPU, memory and I/O devices communicate via system bus.

* The address lines of the system bus that are used to select memory locations can also be used to select I/O devices.

* An I/O device is connected to the bus via an I/O port.

* In CPU's perspective, an I/O device is an addressable data register.

* If a part of the main-memory address space is assigned to I/O ports, then such systems are called as Memory-Mapped IO systems.

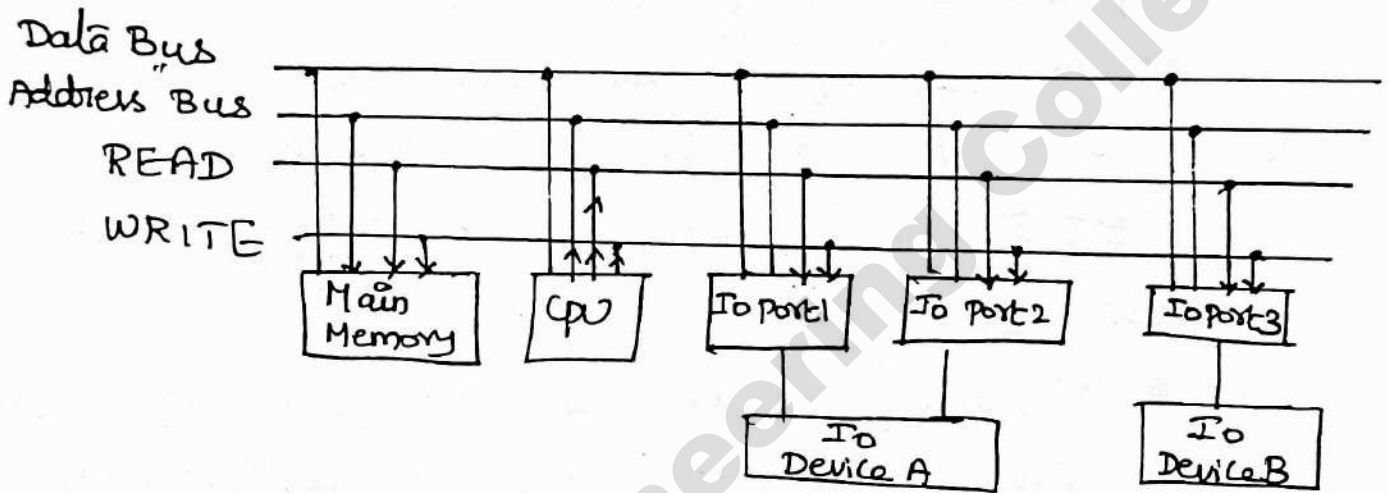
* A memory-referencing instructions becomes I/O instructions automatically, if the address is address of an I/O port.

* The usual memory load and store instructions are used to transfer data words to or from I/O ports. so no special I/O instructions are needed.

Disadvantage:

* CPU spend a lot of time in performing IO related functions.

Eg: IO-related function is testing the status of IO devices



Programmed IO with Memory-Mapped IO

* The control lines READ and WRITE are activated by the CPU for processing memory reference instructions and IO instructions.

* This technique is used in Motorola 680x0 Series

IO-Mapped IO:

* In IO-mapped IO systems, the memory and IO address space are separate.

* Similarly, the control lines used for activating memory and IO devices are also different.

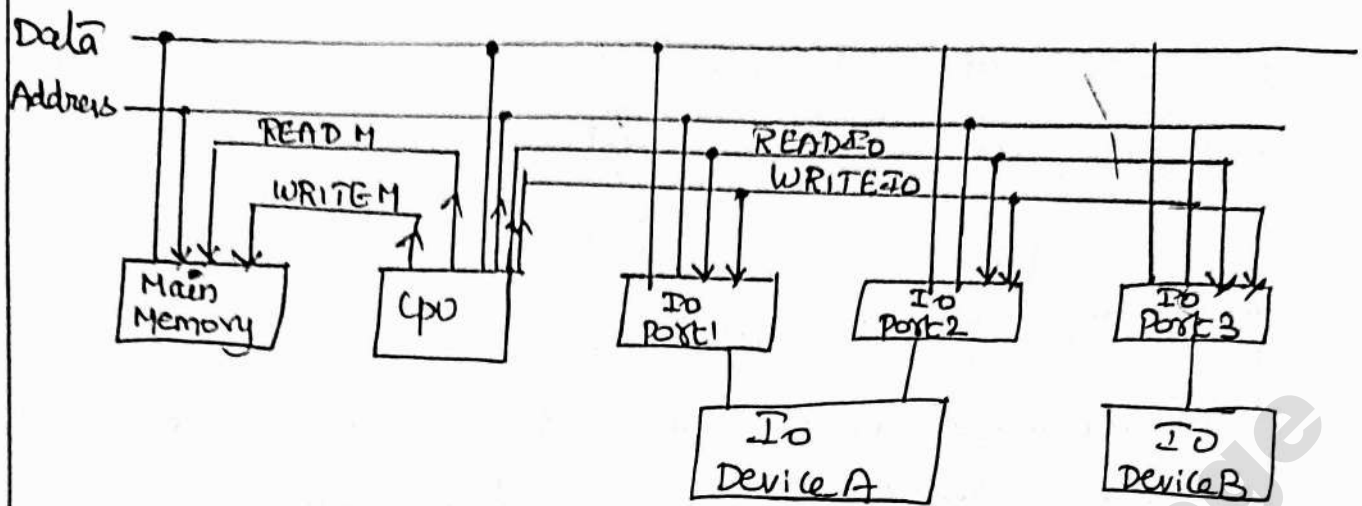
* Two sets of control lines are available READ M and WRITE M are related with memory and READ IO and WRITE IO are related with IO devices.

* A memory referencing instruction activates either the READ M or WRITE M control line, but these signals do not affect the IO devices.

* The CPU executes separate IO instructions to activate the READ IO and WRITE IO control lines.

* It causes a word to be transferred between the addressed port and the CPU.

* This technique is used in Intel 80x86 microprocessor series.



Programmed IO with IO-Mapped IO

IO Instructions:

* Two ^{IO} instructions are used to implement Programmed IO.

* Intel 80x86 series of microprocessors series have two IO instructions called IN and OUT.

IN: The instruction IN X causes a word to be transferred from IO port X to the accumulator register A.

OUT: The instruction OUT X transfer a word from the accumulator register A to the IO port X.

IO Data transfer:

* When the CPU executes an IO instruction such as IN or OUT, the addressed IO port must be ready to respond to the instruction

* For that, the CPU must know the IO device's status so that the data transfer is carried out only when the device is ready.

* In programmed IO, the CPU can be programmed to test the IO device's status before initiating the IO data transfer. The status is specified by a single bit information.

* The CPU must perform the following steps to determine the status of an IO device:

1. Read the IO device's status bit.
2. Test the status bit to determine if the device is ready to transfer data.
3. If ready, proceed with data transfer, otherwise go to step 1.

Limitations of programmed IO:

The programmed IO method has two limitations

1. The speed of the CPU is reduced due to low speed IO devices. The speed with which the CPU can test and transfers data between IO devices is limited due to low transfer rate of IO devices.

2. Most of the CPU time is wasted. The time that CPU spends testing IO device status and executing IO data transfers is too long. That can be spent on other tasks.

- 8) Write the sequence of operations carried out by a processor when interrupted by a peripheral device connected to it. [Apr/May-18] [Nov/Dec-18]

Interrupts:

* An exceptional event that causes the CPU to temporarily transfer control from its current program to another program is called an Interrupt.

* Interrupts are the primary means by which IO devices obtain the services of CPU.

* In Programmed IO, the CPU is continuously checking the status of the IO device still that is ready to transfer the data.

* But, interrupt frees the CPU from the need for checking the status of its IO devices.

* There are various sources of Interrupts both Internal and external to the CPU.

* IO Interrupts are external to CPU.

* Interrupts can also be categorized as

- i) Hardware Interrupts
- ii) Software Interrupts

Eg: Software generated interrupts are, attempt by an instruction to divided by zero, to execute a privileged instruction.

Interrupt Handling Mechanism:

* The basic method of interrupting the CPU is activating a special control line called INTERRUPT REQUEST which connects the interrupt source to the CPU.

* An Interrupt indicator is then stored in a CPU register that can be checked by the CPU periodically.

* CPU does this checking at the end of every instruction cycle.

* If an interrupt is identified, the CPU execute a specific interrupt-handling program.

* Each Interrupt source has its own interrupt service routines.

* So, the CPU must determine the address of the interrupt program to be executed.

* If more than one interrupts requests at the same time, then priorities must be assigned and the interrupt with highest priority is selected for handling.

The following steps are taken while handling the Interrupts:

1. I/O device enables the INTERRUPT REQUEST control line.

2. Interrupt indicator is enabled in the CPU register.

3. The CPU identifies the source of Interrupt.

4. The CPU obtains the memory address of the required Interrupt handler.

* Mostly this address is provided by the Interrupting device along with its interrupt request.

5. The Program Counter(PC) and other processor registers are saved in stack as a subroutine call.

6. The PC is loaded with the address of the interrupt handler. Execution proceeds until return instruction is reached. Immediately, the control is transferred back to the interrupted program.

9

Define Translation Lookaside Buffer (TLB). what is its use? [NOV/DEC-18]

Translation-lookaside buffer (TLB): Making Address Translation fast

* For every read and write access, the page table information is used by memory management unit. (MMU)

* So, if the page table is kept in MMU, it will be easy to access, But, the page table is too large to fit in MMU.

* It is impossible to include a complete page table as the part of processor chip.

* Therefore, the page table is kept in the main memory.

* But still, a copy of small portion of the page table can be accommodated within MMU.

* This small cache is called a Translation Lookaside Buffer (TLB).

* TLB contains the entries that correspond to most recently accessed pages.

* The structure of TLB is slightly modified from the page table in main memory.

* In addition to the information in the Page table entry, TLB includes the virtual addresses of the entry (virtual page number)

Virtual Page Number

TLB

valid	Dirty	Ref	Tag	Physical page address
1	0	1		
1	1	1		
1	1	1		
1	0	1		
0	0	0		
1	0	1		

Page table

valid	Dirty	Ref	physical page or disk Address
1	0	1	
1	0	0	
1	0	0	
1	0	1	
0	0	0	
1	0	1	
1	0	1	
0	0	0	
1	1	1	
1	1	1	
0	0	0	
1	1	1	

Physical memory

Disk storage

Arum College

* Each tag entry in the TLB holds a portion of the virtual page numbers

* Each data entry of the TLB holds a physical page number.

* Because, we access the TLB instead of the page table on every reference, the TLB will need to include other status bits, such as the dirty and the reference bits.

* on every reference, we look up the virtual page number in the TLB.

* If we get a TLB hit, the physical page number is used to form the address and the corresponding reference bit is turned on.

* If the processor is performing a write, the dirty bit is also turned on.

* If a miss in the TLB occurs, must determine whether it is a page fault or merely a TLB miss

i) The page is present in memory, then the TLB miss indicates only that the translation is missing.

In such cases, the processor can handle the TLB miss by loading the translation from the page table into the TLB and then trying the reference again.

ii) The page is not present in the memory, then the TLB miss indicates a true page fault.

* In this case, the processor invokes the operating system using an exception.

Arunai Engineering College