

STORED PROGRAM ORGANIZATION

one processor register: Accumulator.

instruction code format: with 2 parts.

First part: operation to be performed: opcode

Second part: specifies address of operand.

Memory Address tells control where to find operand in memory.

This operand is read from memory and used as the data to be operated on together with the data stored in the processor register.

2^{12} words in memory — 4096

If we store each instruction code in one 16 bit memory word, we have available 4 bits for the OPERATION CODE (opcode) to specify one out of 16 possible operations 12 bits to specify the address of an operand.

The control reads a 16 bit instruction from the PROGRAM portion of memory. It uses the 12-bit address part of the instruction to read a 16 bit operand from the data portion of

memory ??



It then executes the operation specified by the opcode.

ACCUMULATOR: AC

If an operation in an instruction code does not need an operand from memory, Rest of the bits in instruction can be used for other purposes.

Examples of instructions that don't need an operand:

LOAD AC & INCREMENT AC

For these instructions 1 types of operations, the second part of the instruction code (bits '0' through '11') is not needed for specifying a MEMORY address and can be used to specify other operations for the computer.

→ INDIRECT ADDRESS

→ It is sometimes convenient to use the address bits of an instruction code, not as address, but as the actual operand. When the second part of an instruction code specifies an operand, the INSTRUCTION is said to HAVE AN IMMEDIATE OPERAND.

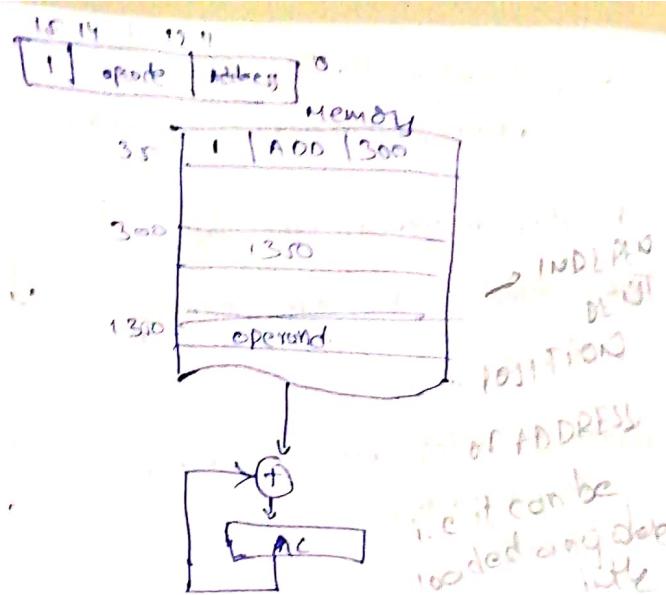
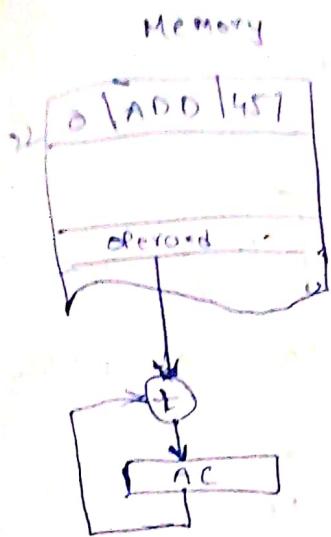
→ When the second part specifies the address of an operand the INSTRUCTION is said TO HAVE A DIRECT ADDRESS.

→ A THIRD possibility called INDIRECT ADDRESS where the bits in the second part of the instruction designate an address of a memory word in which the address of the operand is found.

→ one bit of the instruction code can be used to distinguish between a direct address & indirect address.

→ Consider INSTRUCTION CODE format shown in Fig(a)

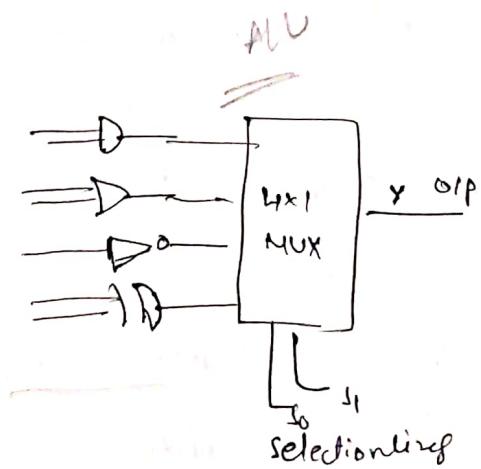
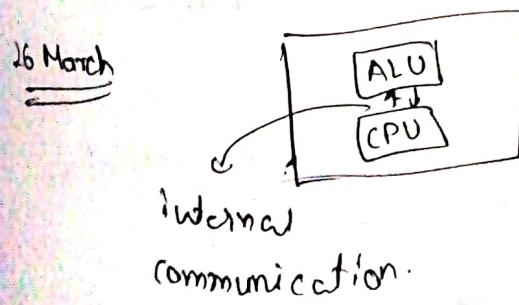
It consists of 3 bit opcode, a 12-bit address and an indirect address mode bit designated by I. The mode bit is '0' for direct & '1' for indirect address. it is placed in address section.



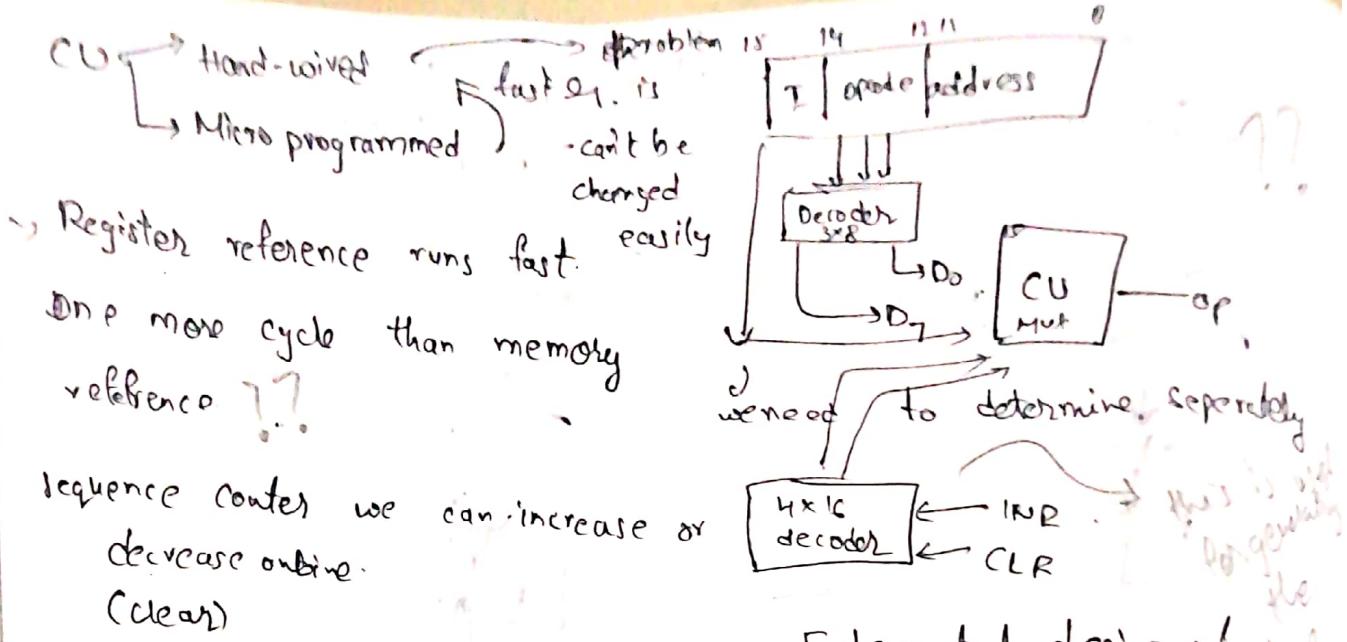
bit is '0' . So the instruction is recognised as a memory direct address instruction. The opcode specifies ADD instruction and the address part is the binary equivalent of 457. The control finds the operand in memory at address 457 and add it to the content of accumulator.

→ The INSTRUCTION in address 35 has mode bit I=1. Therefore it is recognized INDIRECT ADDRESS instruction. The address part is binary equivalent of contents of address 300. The control goes to address 300 to find the address of the operand. So the address of the operand in this case is ³⁰⁰1350. The operand found in address 1350 is then added to content of accumulator.

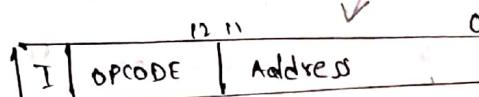
EFFECTIVE ADDRESS



→ Timing & Control signal to ALU



MP → program which sends timing & seqⁿ to ALU.



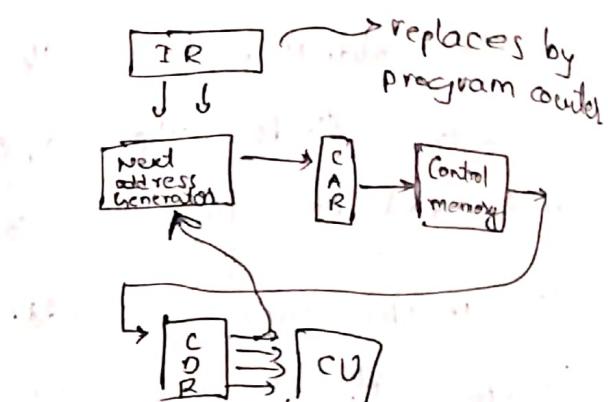
→ If we make control program modified easily control unit is changed.

→ Instruction decode, Instruction fetch.

27 March

COMPUTER REGISTERS:

- Computer instructions: stored in consecutive memory locations and are EXECUTED SEQUENTIALLY one at a time
- The control reads an instruction from a specific address in memory and executes it
- It then continues by reading the next instruction in sequence and executes it and so-on
- This type of instruction sequencing needs a COUNTER to calculate the address of the next instruction after execution of current instruction



pipelining → parallel execution of unrelated exp

It is also necessary to provide a register in the control unit for storing the INSTRUCTION CODE after it is read from memory.

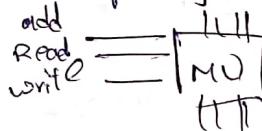
The computer needs PROCESSOR REGISTERS for manipulating data.

- A register for holding the memory address

- These requirements dictate the REGISTER CONFIGURATION in figure.

- The registers are listed in Table 5.1 together with a brief of their function and the # of bits they contain.

- The memory unit has a CAPACITY of 4096 words and each word contains 16 bits. Twelve bits of an INSTRUCTION WORD are needed to specify the address of an OPERAND. This leaves 3 bits for operation part (opcode) of instruction and a bit to specify a direct or indirect address.



- The DATA REGISTER (DR) holds the operand read from memory

- The ACCUMULATOR register is a general purpose processing register.

- The instruction read from memory is placed in the INSTRUCTION REGISTER (IR).

- The temporary register (TR) is used for holding temporary data during processes. → Indirect is used when it is relatively unknown.

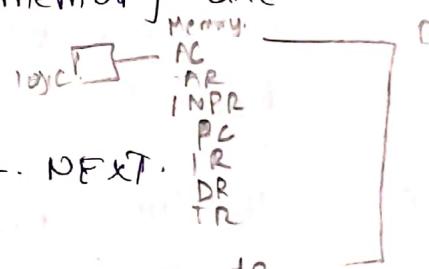
- The MEMORY ADDRESS REGISTER (MAR) known has 12 bits. Since this is the width of a memory address.

- The PROGRAM COUNTER (PC) also has 12 bits and it holds address of the NEXT INSTRUCTION to be read from memory after the current instruction previously stored in memory

- Instruction words are read and executed in a sequence unless BRANCH INSTRUCTION is ENCOUNTERED.

- A BRANCH INSTRUCTION calls for a transfer to a non consecutive instruction in the program. The address part of a branch instruction is transferred to PROGRAM COUNTER to become the address of next instruction.

- To read an instruction, the content of PC is taken as the address from the memory and a memory read cycle is initiated



March 29 COMMON BUS SYSTEM

- The BASIC COMPUTER under discussion;

Has 8 registers, a MEMORY unit, & a CONTROL UNIT

- Paths to transfer information from one register to another and between Memory & registers

- The number of wires will be excessive if connections are made between the outputs of each register and the inputs of the other registers

- A more efficient scheme for transferring information in a system with many registers is to a

the connection of the registers and MEMORY of the BASIC computer to a common bus shown in figure. the outputs of 7 registers and memory are connected to the COMMON BUS. The specific output that is selected for the bus lines at any given time, is determined from the binary value of the selection variables S_2, S_1 & S_0 . The number along each output shows the DECIMAL equivalent of the required binary selection. For example the number along the output of DR is '3'. The 16 bit outputs of DR are placed on the bus lines when $S_2, S_1, S_0 = 011$. Since this is the binary value of decimal '3'.

- The lines from common bus are connected to the inputs of each register and the data inputs of the memory.

- the particular register whose LD (load) input is ENABLED receives the data from the bus during the next clock pulse transition
LD \rightarrow performing the operation

- The memory receives the contents of the bus when its "write input is activated". The memory places its 16 bit output onto the bus when the READ input is activated.

- The memory places its 16 bit output onto the bus when the READ input is activated and $S_2, S_1, S_0 = 111$.

→ 4 registers, DR, AC, IR & TR have 16 bits each. Two registers AR & PC have 12 bits each since they hold memory address.

→ When the contents of AR or PC are applied

to the 16 bit common bus, the 4 most significant bits are set to 0's??

→ when AR or PC receives information from the bus, only 12 least significant bits are transferred into the register.

→ The INPUT REGISTER, INPR and the output register OUTR have 8 bits each and communicate with the LEAST SIGNIFICANT BITS in the bus. INPR is connected to provide information from the bus. This is because INPR receives a character from

an input device which is then transferred to AC. OUTR receives a character from AC and delivers it to an output device. There is no transfer from OUTR to any of the other registers.

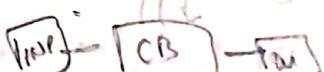
→ The 16 lines of the common bus receive information from 6 registers and the memory

→ The bus lines are connected to the inputs of 6 registers and the memory

→ 5 registers have 3 control inputs.

3 control inputs: LD (load), INR (increment) & CLR (clear)

MEMORY ADDRESS:



→ The input data and output data of the memory are connected to the common bus, but the memory address is connected to AR. Therefore AR must always be used to specify a memory address.

By using a single register [MAR] for address we eliminate the need for an ADDRESS BUS that would have been needed otherwise.

The 16 bits of AC come from an adder & logic circuit. This circuit has 3 sets of inputs. One set of 16-bit inputs come from the outputs of AC. These 16 inputs are used to implement register microoperations such as complement AC & shift AC.

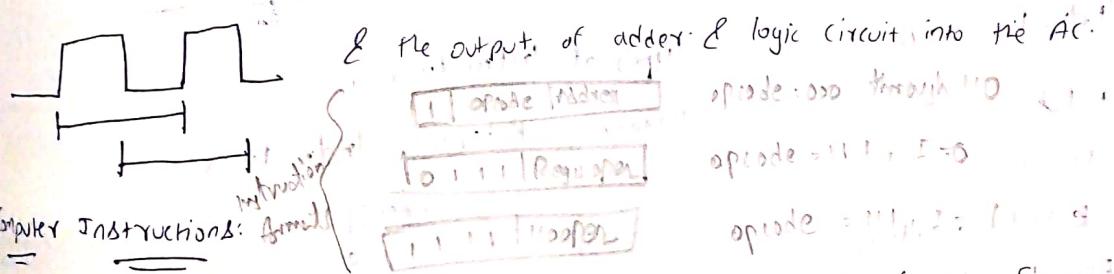
Another set 16-bit inputs come from DATA REGISTER (DR). The inputs from DR register & AC are used for Arithmetic operations & logic microoperations such as ADD DR to AC (or) AND DR to AC.

The result of an addition is transferred to an accumulator and the end carry of the addition is transferred to flip-flop E (extended AC bit).

The third set of 8-bit inputs come from the input register INPR.

Note: The content of any register can be applied on to the bus and an operation can be performed in the adder and logic circuit during the same clock cycle.

The clock transition at the end of the cycle transfers the content of the bus into the designated ^{destination} register.



Instruction format: The basic computer has 3 instruction code formats. Fig 5.5 each format has 16-bits.

The opcode part of the instruction contains 3 bits & the meaning of the remaining 13 bits depends on the opcode encountered.

A memory reference instruction uses 12 bits to specify an addressed one bit to specify the addressing mode. It has the value '0' for direct address

'1' for indirect address.

The Register reference instructions are recognized by the operation code 111.

With a '0' in the leftmost bit of instruction.

A register reference instruction specifies an operation on by a test of the

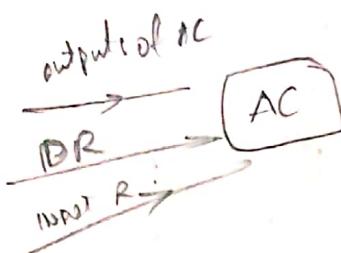
AC register. An operand from memory is not needed, therefore the other 12 bits are used

to specify an operation or test to be executed.

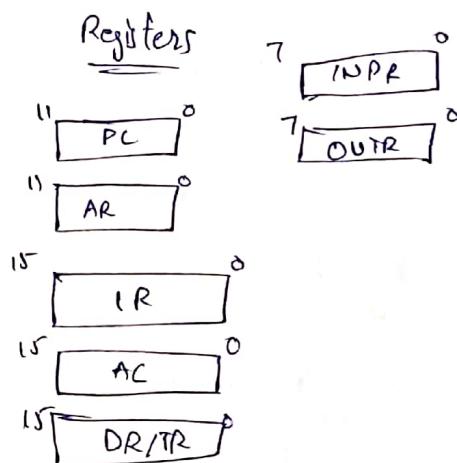
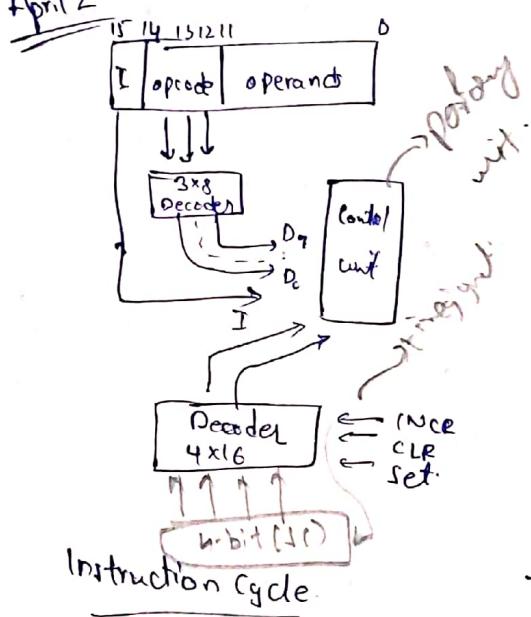
- Similarly, an I/O instruction does not need a reference to the memory & is recognized by the opcode 111 8 with leftmost bit as 1 of the instruction.

The remaining 12 bits are used to specify the type of I/O operation (or) test performed.

The type of instruction



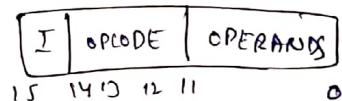
April 2



Types of Instructions

1. Fetch
2. Decode
3. Effective address
4. Execute

1. Memory- Reference Instructions



2. Register- Reference Instruction

12-14 - bits — 111 — I → O

3. I-O instruction

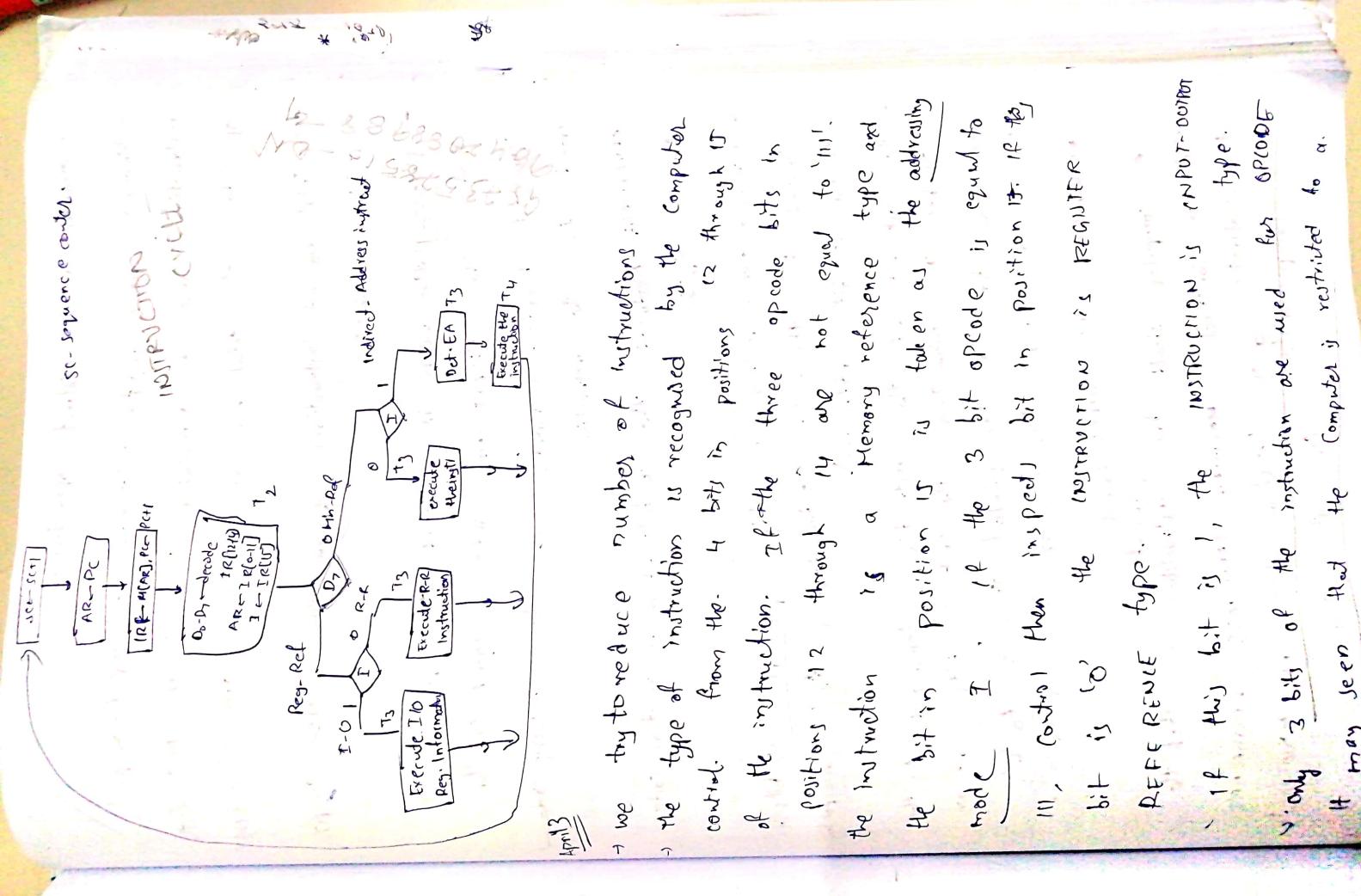
12-14 - bits — 111 — I → I

Fetch

T₀: AR ← PC
T₁: IR → M[AR],
PC ← PC + 1

Decode: I ← IR[15] JR[0-11]

T₂: D₁ : D₈ ← Decod[12-14], AR ← JR[0+1], I ← IR[15]



- we try to reduce number of instructions
- the type of instruction is recognized by the computer control from the 4 bits in positions 12 through 15 of the instruction. If the three opcode bits in positions 12 through 14 are not equal to '111', the instruction is a memory reference type and the bit in position 15 is taken as the addressing mode I.
- if the 3 bit opcode is equal to '111', control then inspects bit in position 15. If the bit is '0', the instruction is REGISTER type.
- if this bit is '1', the instruction is INPUT-OUTPUT type.
- only 3 bits of the instruction are used for opcode
- it may seen that the computer is restricted to a

maximum of 8 distinct operations

- However, since register reference and input-output instructions use the remaining 12 bits as part of the op-code, the total no. of instructions can exceed eight.
- Total no. of instructions chosen for basic computer = 25 → Table 5.2.
- Three letter word is utilized to represent the op-code. The Hexadecimal code is equal to the equivalent hexadecimal number of the binary code used for the instruction.
- ADVANTAGE: Reduced the 16 bits of an instruction code to 4 digits.

- Memory reference instruction: Address part of 12 bits is denoted by three x's. (3 base decimal digits corresponding to the 12-bit addresses).
- The last bit of the instruction is designated by the symbol 'I' when I=0, the last 4-bits of an instruction have a hex equivalent from '0' to 'F' (since the last bit is '0') . It is equivalent from '8' to 'E'.

INSTRUCTION SET COMPLEXITY:

- A computer should have a set of instructions so that the user can construct machine language programs to evaluate any function that is known to be computable. $f(n) = n^2 - 2 \equiv 0$

with external word

Input & output instructions. Needed for communication
Program control instructions e.g.: Branch & Load -

DIGITAL COMPUTER'S (DC's)

Decision Making capabilities are an important aspect of

between these two units.

The user must have the capability of moving information memory, but all computations in processor register. Therefore

computer (all of stored program) is stored in set. or best instruction

The bulk of binary information in a digital

4. Input & Output instructions.

They are execute.

In what order

that check status condition

3. Program control instructions together with instructions

processor register.

2. Instructions for moving information to & from memory and

1. Arithmetic, logical, shift instructions.

each of the following categories

computer includes a sufficient number of instructions in

- the set of instructions is said to be complete if the

execute any function that is known to be COMPUTABLE

the user can construct programs to

A computer should have SET OF INSTRUCTIONS so that

Instruction set (complete):

4 principle of Pass/money

ocean's razor. Make the model simple, but not any simpler.

Modeling is done for predicting things.

Page 11-4

Hardware Organization: control logic is implemented

Training and Control: control unit organization

are used in quantum computers

with gates, flip flops, decoders.

Control Logic is implemented

Hardware Organization: Hardwired Control of microprogrammed Control

There are 2 major types of control organization

microprogram operations for the accumulator

common bus, control inputs in processor register and

and provide control inputs for the multiplexers in the

The control signals are generated in the control unit

control signal

register unless the register is enabled by a

clock pulses don't change the state of a

registers in the control unit.

Registers in the system, including flip flops and

clock pulses are applied to all flip flops and

registers in the control unit.

The timing for all registers in the BASIC COMPUTER

is controlled by a master clock generator.

TIMING OF CONTROL

MINIMUM SET that provides most of capabilities

The instructions listed in Table 52 constitute a



The subscripted decimal number is equivalent to the binary value of the corresponding operation code B15 of the instruction is transferred to a buffer flip through it by the symbol I. Bits 0 designates by the symbol E.

The opcode in bits 11 through 14 are decoded with the decoder. The outputs of the decoder are designated by the symbols D_0 , D_1 , D_2 , D_3 , D_4 , D_5 through D_7 .



• 11. through ① sing

It is divided into 3 parts : the I bit, the opcode, and

RELIESTER (LR) in Fig 5-6.

An instruction read from memory is placed in the ~~instruction~~

c) a number of control logic gates.

g) Two Decade counters

QUESTION ANSWER Block diagram has 3 parts.

PROGRAM IN CULTURAL MEMORY.

model modifications can be done by updating the Micro

In the microprogramming mode control, any required changes or

modified or changed.

The various compounds if the design has to be

Hardboard world equities changes by country among.

Wavelength Selection of Microscopic Images

→ central information is stored in a central MEMORY

MICROPROGRAMME OR MICROPROGRAM:

• of operation.

Advantages: It can be optimized to produce soft mode with goals, help-flops, decoders and other DC's.

Timing signals to output of the decoder, so

Sequence counter to, which when activated by

- The first positive transition of the clock clears

- Initially, the CLR input of SC is active

clock

- The SC responds to the positive transition of the control signals.

The timing diagram (Fig. 5.7) shows the time relationship

$$D_{31} \equiv SC \rightarrow 0$$

outputs are active. This is expressed symbolically by the

at the time T_1 , SC is cleared to 0. If decoder

to provide timing signals T_0, T_1, T_2, T_3, T_4 in sequence

EXAMPLE: Consider the case where SC is incremented

the next active timing signal to be T_5 .

Once in a while, the counter is cleared to 0, causing

signals out of sync decoder.

incremented to provide the sequence of timing

synchronously. Most of the time the counter is

- The sequence counter can be incremented or cleared

timing signals through T_5 .

- The outputs of the counter are decoded into 16

from a through 15.

- The 4 bit sequence counter can count in binary

are applied to the counter loci (all 1's

8

aff

27

de la fin

24

१५९

Fayne

PULMONO

5

1270

24

The positive voltage across the diode will be equal to the voltage drop across the diode and the voltage drop across the resistor. The current through the diode will be the same as the current through the resistor.

CLR input is active. This produces the sequence of training signals to SC. SC is interconnected with every five clock transitions, unless it's connected to the binary signal, To.

To is active during one clock cycle. The positive clock transition latches information received from the binary signal To out of the decoder (N16). The first positive transition of the clock clears SC to 0, which then activates the latch (N16).

If happens if Pipeline process stuck up? jump in J17

→ In a non Pipeline system it takes 30 cycles

IF ID EV MEM WB

IF ID EV MEM WB

→ ADD R3 LOH, OAH (Register immediate) ADD R3 R1, 100
→ Every program is prefaced by base, limit

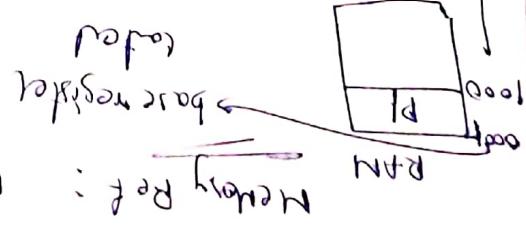
much we need to add is called offset

register offset from base to limit here



Memory Ref: We need to calculate base register + offset another register.

Ex: Reg's Ref: ALU applies operation immediately and stored.



ID: 12, 13, 14 bits passed to decoder.

1) Relationship between timing signals & memory cycle
 It is not closed & the timing signal will continue with $D_{3T_6} \rightarrow D_{3T_7} \rightarrow T_0$.
 Fig. 3 waveform in Fig. 5.7 shows how SC is delayed when D_{3T_7} is applied to it from the operation control register because output of the AND gate is timing signal. T_2
 At implementation the control function D_{3T_7} becomes active, so the output signal is
 applied to the CLK input of the SC. On the next positive clock transition it is applied to the write cycle (T) which is initiated with the rising edge of a memory read (T) write cycle counter is decremented to 0. Do which causes the timing signal T_0 to become active, instead of T_5 by the time the next clock goes through it's positive transition.
 A memory read (T) write cycle initiated by a timing signal will be completed a timing relationship is not valid in many computers, because the memory cycle is usually longer than the previous clock cycle.
 In such case it is necessary to provide ~~next~~ cycles in the processor until the next memory word is available.
 To facilitate the presentation we will assume that a WAIT period is not necessary in the basic computer.

2) To fully operate the operation of the computer it is very difficult that one consider the timing relation ship between the clock transition & the timing signals.
 In Fig. 5.7 shows how SC is delayed when D_{3T_7} is applied to it from the operation control register because output of the AND gate is timing signal. T_2
 It is usually longer than the previous clock cycle.
 In such case it is necessary to provide ~~next~~ cycles in the processor until the next memory word is available.

3) To fully operate the operation of the computer it is very difficult that one consider the timing relation ship between the clock transition & the timing signals.
 In Fig. 5.7 shows how SC is delayed when D_{3T_7} is applied to it from the operation control register because output of the AND gate is timing signal. T_2
 It is usually longer than the previous clock cycle.
 In such case it is necessary to provide ~~next~~ cycles in the processor until the next memory word is available.

Fig 5.8 shows how the first two register transfers are implemented in the bus system to provide data path for the transfer of PC to AR we must apply timing signal to achieve the following connection:

1. Place the content of PC onto the bus by making the bus selection inputs $S_2, S_0 = 010$.

2. Transfer the content of the bus to AR by enabling the LD (load) input of AR.

The next code transition initiates the transfer from PC to AR since $T_0 = 1$

$\rightarrow I \rightarrow IR(1)$

shows how the first two groups of figures are related.

$$\left. \begin{array}{l} T_0: AR \rightarrow PC \\ T_1: IR \rightarrow H[AR], PC \rightarrow PC+1 \\ T_2: D_0, D_1, \dots, D_7 \rightarrow DECODE\ IR(12-1H), AR \rightarrow [R(0)] \end{array} \right\} \text{Decoder}$$

RECEIVER TRANSFER software

MICRODEGRADATION CYCLE: The microdegradation cycle can be followed by the following steps:

- This sume the calculate translation into results the SC from 1000 to 0000
- After calculate cycle has to achieve & to indicate
- This sume the calculate translation into results the SC from 1000 to 0000
- Translating unit of a segment of instruction
- A program residing in the memory through a cycle for each instruction
- It is performed in the computer by going through a cycle for each instruction
- Each instruction cycle is divided into a sequence of subcycles
- Each instruction cycle in turn is subdivided into a sequence of subcycles
- In order to a sequence of subcycles
- In order to a sequence of subcycles
- The basic computer each instruction cycle consists of the following phases.
- (a) phases.
- (b) subcycles
- Fetch & instruction

DDD - lawline

In order to implement the second statement

$T_1 : IR \leftarrow M[AR], PC \leftarrow PC + 1$ it is necessary to use timing signal T_1 to provide the following connections in the bus system.

i) Enable the Read input of memory if necessary
ii) Place the content of memory onto the bus by making $S_2, S_1, S_0 = 111$

iii) Transfer the content of the bus to IR by enabling the LD input of IR.
iv) Increment PC by enabling the INR input of PC.
The next clock transition indicates the read and increment operations since $T_1 = 1$

→ Fig 5.8 is not entire bus & it doesn't contain all registers

→ It duplicates a portion of the bus system and shows how T_0 and T_1 are connected to the control inputs of the registers, the memory and the bus selection inputs.

→ Multiple input or gates are included in the diagram because there are other control functions that will indicate similar operations.

April-16

→ Illegal jumps → control hazards.
Memory: Storage area used to store programs, and data to perform specific task (Auxiliary, Main, Cache)
Auxiliary Memory: Not directly connected to CPU
Secondary storage devices
- Tapes, discs.

Main Memory: RAM

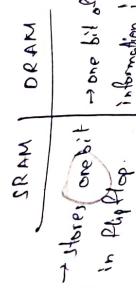
Code Memory:

Memory organization

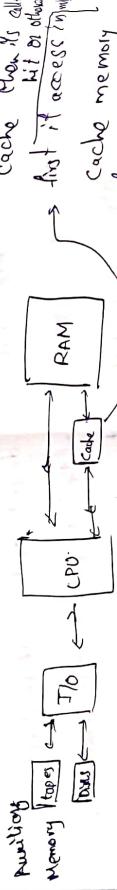
Memory organization

RAM: Random Access Memory

↳ facilitates direct access to a memory cell

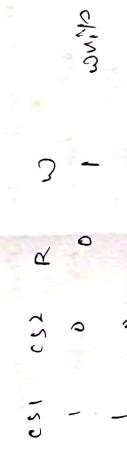
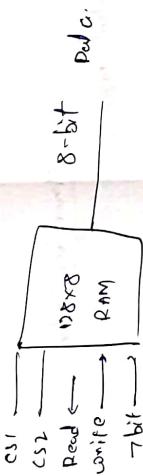


- SRAM | DRAM
- stores one bit of information in parallel.
- As long as power is supplied, data is alive → generally used to create cache-memory → faster & expensive
- very large in size → Requires refreshing
- used to create main memory → slower & cheaper.
- stores one bit of information in parallel.
- permanent are stored in RAM
- Bootstrap loader
- Cache Performance is measured in terms of hit ratio.
- Cache hits a word in Cache than its called hit or otherwise called miss.
- If first access to any data, it goes to DRAM to get data.
- CPU, if one word more control pins are provided.
- 128 words RAM & each word is of 1 byte



RAM is the most suitable memory to communicate with the CPU, if one word more control pins are provided.

- 128 words RAM & each word is of 1 byte
- Cache Performance is measured in terms of hit ratio.
- Cache hits a word in Cache than its called hit or otherwise called miss.
- If first access to any data, it goes to DRAM to get data.
- CPU, if one word more control pins are provided.
- 128 words RAM & each word is of 1 byte



ROM-512 Word ROM

- one each word - 8-bit.
- How to map memory
- CS1 → 512x8 ROM
- CS1 → 8-bit data.

Determine

→ During

that,

present

and it

after

→ Cache Performance

- M.

- Re-

hit ratio.

→ Cache

than its called

hit or otherwise

called miss.

→ Cache memory

of any data.

REFE

→ CPU

→ DRAM

→ Cache

This mapping can be done in 3 ways (Associative, Direct, Set-Associative)

Associative

Sequential Computer Parallel Computer

Parallel

$$T_0 : AR \leftarrow PC$$

$$T_1 : TR \leftarrow M[AR], PC \leftarrow PC + 1$$

$$T_2 : D_0, D_1, \dots, D_7 \leftarrow \text{DIODE}$$

$$TR [D_7-D_0], AR \leftarrow TR [0-11]$$

$$T \leftarrow T + 15$$

be made

it takes the execution leads to
instruction code.

now

determine the type of instruction:

- During time T_3 , the control unit determines the type of instruction that was just read from memory. The flowchart of Fig 5.9 presents an initial configuration for the Instruction Cycle and it shows how the control determines the instruction type after the decoding. (There are 3 possible instruction types):

- Memory Reference instructions

- Register Reference instructions

- I/O instructions

- Decoder output D_7 is equal to '1' if the operation code is equal to binary bit. But if $D_7 = 1$, the instruction must be a REGISTER REFERENCE or INPUT - OUTPUT type

- If $D_7 = 0$, the operation code must be one of other seven values 000 through 110, specifying a MEMORY REFERENCE instruction.

- CONTROL then inspects the value of first bit of the instruction which is now available in flip flop. If $D_7 = 0$ and $I = 1$, we have a MEMORY REFERENCE instruction with an INDIRECT ADDRESS $AR = M[AR]$. If $D_7 = 0$ and $I = 0$, we send the least significant

- It is then necessary to read the effective address from memory. The MICRO OPERATION for the INDIRECT ADDRESS condition can be symbolized by the REGISTER TRANSFER statement $AR \leftarrow M[AR]$. Initially AR holds the address part of the instruction mapping memory

This address is used during memory READ. The word at the address given by AR is read from memory and placed on the common bus. The LD (load) input of AR is then enabled to receive the indirect address that resides in the 12 least significant bits of the memory word.

REGISTER REFERENCE INSTRUCTION

- Recognized by the control when $D_7 = 1$ and $I = 0$. These instructions use bits '0' through '11' of the instruction code to specify one of 12 instructions. These 12 bits are available in IR [0-11]. They were also transferred to AR during time T_2 .

April 12

contrary to stack of trays where the tray itself may be taken out or inserted, the physical registers of a stack are always available for reading or writing. It is the content of the word that is inserted or deleted.

The two operations of a stack are the INSERTION and DELETION of items. The operation of insertion is called push (or pushdown) because it can be thought of the result of pushing a new item on top.

The operation of deletion is called pop (or pop up) because if can be thought of as the result of removing one item so that stack pops-up.
REMOVING one item so that stack pops-up
— CRUCIAL DIFFERENCE

However nothing is "pushed" or "popped" in a computer stack, these operations are simulated by "INCREMENTATION" or "DECREMENTATION" the STACK POINTER REGISTER.

REGISTER STACK: A stack can be placed in a portion of large memory or it can be organised as a collection of a finite number of memory words or register.

Fig 8.3. ORGANIZATION of a 64-word Register stack.
SP (stack Pointer Register): Contains the binary number whose value is equal to the address of the word that is currently on top of the stack.

Three items are placed in the stack: A, B, C in that order item C is on the top of the stack so that the content of SP is now 3. To remove the top item, the stack is popped by reading the memory word at address 3. Decrementing the content of the stack pointer. Item B is now on top of stack since SP holds address 2.

- To insert a new item, the stack is pushed by incrementing SP and writes a word in the next higher location in the stack.

Note: Item C has been read out but not physically removed. This does not matter because when stack is pushed, a new item is written in its place.

MEMORY STACK

- A stack can exist as a stand alone unit as in the case of Fig 8.3 or can be implemented in a RAM attached to CPU
- The implementation of a stack in the CPU is done by assigning a portion of memory to a stack operation and using a processor register as a stack pointer.
- Fig 8.4 shows a portion of Computer memory partitioned into 3 segments: Program, data, stack.

As shown in Fig 8.4, the initial value of SP is 4001 and stack grows with decreasing address.
 Thus the first item stored in the stack is at address 4000, the second item is stored at address 3999, and the last address that can be used for the stack is 3000.

Register Stack : FULL, EMPTY

Stack limit checks: FULL : one bit register is set to 1

EMPTY : set to 1, when stack is empty.
 we assume that the items in the stack communicate with a DATA REGISTER, (DR)

→ A new item is inserted with the push operation as follows:

$$\boxed{\begin{array}{l} SP \leftarrow SP-1 \\ M[SP] \leftarrow DR \end{array}}$$

The stack pointer is decremented so that it points at the address of the next word. A memory write operation inserts the word from DR into the top of the stack.

$$\boxed{\begin{array}{l} SP = SP \\ DR \leftarrow M[DR] \\ SP \leftarrow SP+1 \end{array}}$$

→ A new item is deleted with a pop-operation as follows

$$\boxed{\begin{array}{l} DR \leftarrow M[SP] \\ POP \\ SP \leftarrow SP+1 \end{array}}$$

= Pop item on the stack read from
 stack pointer is incremented to point at the top.

→ No hardware checks for stack overflow (full stack) or underflow (empty stack)

→ Stack limits could be checked by 2 processor registers : one to hold upper limit & the other for lower limit.

- How much (size of memory)
- How fast (speed of memory write/read)

Memory Hierarchy: CACHE MEMORY : very fast but costly

MAIN MEMORY: Medium

(Secondary) AUXILIARY MEMORY:



Instruction PIPELINING:

- It is a technique for implementing instructions.
- LEVEL PARALLELISM within a single processor. Pipelining attempts to keep every part of the processor busy with some instruction by dividing instructions into a SERIES of SEQUENTIAL steps ("pipeline") performed by DIFFERENT processor units with different parts of instructions processed in parallel.

Example Processors

Atmel AVR, PIC Microcontroller.

Intel Pentium (Pipeline size: 7, 10 or even 20 stages)
Protocol for later Will Netburst cores from Intel.

April 26

Reduced Instruction Set Computer Architecture (RISC)

Microprocessor Era:

- Computer scientists were trying to copy the same complex instructions used in Mainframe → 1950s. Some wanted to expand those instructions.

- Patterson & Hennessy's idea: not only should we not make instruction set it more complicated, we should make it even simpler: called RISC.

RISC architecture relied on a simpler collection of general functions the processor would perform shrinking the number of transistors needed to carry out a task.

Today, 99% of more than 16 billion microprocessors produced each year are based on RISC architecture, powering smartphones, tablets & IoT devices

✓ everything in the world connected to internet.

→ those early computers had perhaps 25% more instructions "complex".

under RISC design, but ran 5 times as fast.

Example : BERKELEY RISC and STANDARD MIPS processor

MIPS : Microprocessor without interlocked pipeline stages - Hennessy
RISC architecture.

which speeded up processing by using instructions loaded from the memory into a register which could be accessed faster.

- DOMAIN SPECIFIC ARCHITECTURES FOR "HL".

SUPERSCALAR PROCESSOR : (SSP)

SSP is a CPU that implements a form of parallelism called INSTRUCTIONS LEVEL PARALLELISM within a single processor.

In contrast to a scalar processor that can execute at most one instruction per clock cycle, a superscalar processor can

execute more than one instruction during a clock cycle by

SIMULTANEOUS DISPATCHING multiple instructions to DIFFERENT EXECUTION UNITS on the same processor.

~ It therefore allows for more throughput (the number of instructions that can be executed in a unit of time) than would otherwise be possible at a given clock rate. (They are on the same processor.)

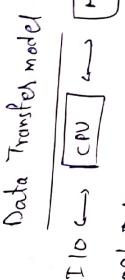
Each execution unit is not a separate processor (on a core) if the processor is a multi-core processor) but an execution resource within a single CPU such as

ARITHMETIC & LOGIC UNIT

Single core super scalar processor : 1118

30/04/2019

Tutorial



- 1) Programmed I/O
2. Interrupt - Initiate.
- 3) Direct Memory Access

Programmed I/O.

- It is the result of execution of I/O instruction written in a program
- CPU is kept in program loop and write an I/O device to transfer data.
- Constant monitoring of I/O device is required by CPU

Interrupt - Initiate I/O:

- processor with continuous task → hardware interrupt
- whenever an I/O device is ready to transfer data, an interrupt is initiated to the CPU
- Then the CPU have normal operation
- in attends to the other ISR
- (interrupt service routine)
- Shows interrupt sequence
- interrupt that occur during the program execution (RESET)

plasticine / non-manageable interrupt.
it is overridable
non-available

→ 2 Model of DMA Transfer:

- 1) Cycle Transfer 2) Burst Transfer
- ✓ Buses will go back
to the control of CPU
after the transfer of
every word

3) May 2019 1st class

CACHE MEMORY: ORGANIZATION as a STORE-CREATE MAPPING:

PROPERTY of Locality of Reference:

- Analysis of large number of typical programs has shown that the "REFERTES to Memory" at any given interval of time tend to be confined within a few localised areas in memory.

This phenomenon is known as PROPERTY OF LOCALITY of Reference.

- The reason for this property may be understood by considering that a typical computer program flows in a straight-line fashion with program loops and subroutine function calls encountered frequently.

- when a loop is executed, the CPU repeatedly refers to the set of instructions in memory that constitute

- Every time a given subroutine /function is called, its set of instructions are fetched from memory
- Thus loops and subroutines tend to locate the

for fetching instructions.

the
keeping
data in

- To lesser degree, memory references to data also tend to be localised. Table look-up procedure repeatedly refers to that portion in memory where the table is stored.
- Iterative procedures refer to common memory location and array of numbers are confined within a local portion of memory.

- The result of all these observations is the locality of reference property which states that over a short interval of time, the addresses generated by a typical program refer to a few localized areas of memory repeatedly while the remainder of memory is accessed relatively infrequently.

If the access to the memory is not to refer just at one point but

Idea : If the active portions of the program and data are placed in a fast, small memory, the average memory access time can be reduced, thus reducing the total execution time of the program.

$\frac{\text{hit}}{\text{miss}} = \frac{\text{hit frequency}}{\text{miss frequency}}$

MIPS (millions of instructions per second) such a fast, small memory is referred to as a cache memory.

The cache memory access time is less than the access time of main memory by a factor of 5-10. The cache is the fastest component in Memory Hierarchy and approaches the speed of CPU components. (eg: register)

access time will approach the access time of cache.

Although the cache is only a small fraction of the size of main memory, a large fraction of memory requests will be found in the fast cache memory because of the locality of reference property of programs.

The basic operation of cache is as follows. When the CPU needs to access memory, the cache is examined. If the word is found in the cache, it is read from the last memory. If the word addressed by the CPU is not found in the cache, the main memory is accessed to read the word. A block of words containing the one just accessed is then transferred from main memory to cache memory. The block size may vary from one word (the one just accessed) to about 8 words adjacent to the one just accessed. In this manner, some data are transferred to the cache so that future reference to memory finds the required words in fast cache memory.

31 May 2019

2nd class.

HIT Ratio: The performance of cache memory is frequently measured in terms of a quantity called HIT RATIO. When the CPU refers to memory and finds the word in cache, it is said to produce a HIT.

- When the CPU refers to memory and finds the word in cache, it is said to produce a HIT.

- If the word is not found in cache, it is in main memory and it counts as a Miss.

Total CPU reference to memory (hits plus miss) is the HIT ratio.

The hit ratio is best measured experimentally by running representative programs in the computer and measuring the number of hits and misses during a given interval of time.

After hit ratio of 0.9 and higher have been reported this high ratio verifies the validity of the

Locality of Reference Property

The average memory access time of a computer system can be improved considerably by use of a cache. If the hit ratio is high enough so that most of the time the CPU accesses the cache instead of main memory, the average access time is closer to the access time of fast cache memory.

The basic characteristic of cache memory is its fast access time. Therefore very little time or no time must be wasted when searching for words in cache.

→ THE TRANSFORMATION OF DATA FROM MAIN MEMORY TO CACHE MEMORY IS REFERRED TO AS A MAPPING PROCESS.

3 types of mapping procedures are of practical interest : when considering the organization of cache memory.

(1) All inclusive mapping

(2) Direct mapping

(3) Set associative mapping.

say

- To help in discussion of these 3 mappings procedure will use a specific example of memory organization shown in Fig 12.10
- The main memory can store 32K words of 12 bits each

(a)

- The cache is capable of storing 512 of those words at any given time.
- For every word stored in cache, there is a duplicate copy in main memory.
- The CPU communicates with both memories. It first sends a 15-bit address to cache

If there is a hit, the CPU accepts the 12-bit data from cache. If there is a miss, the CPU reads the word from main memory and the word is transferred to cache

(b)

12.2 Main Memory: $32K \times 12 \rightarrow$ of 11 bits each.

$32K \times 10 \rightarrow 2^{15} \rightarrow$ total size.

first send

15-bit address to cache.

12-bit data from cache.

- If there is a miss, the CPU reads the word from main memory and the word is then transferred to cache.
- The CPU communicates with both memory if first sends 15-bit address to cache. If there is a hit, the CPU accepts the 12-bit data from cache.

FAST

NO TIME

CACHE

- If there is a miss, the CPU reads the word from main memory and the word is then transferred to cache.

ASSOCIATIVE MAPPING: The FASTEST and MOST FLEXIBLE cache

MEMORY

ASSOCIATIVE MAPPING

ORGANIZATION

- The AM stores both the ADDRESS and CONTENT of the memory word.

radical word.

THIS PERMITS ANY LOCATION IN CACHE TO STORE ANY WORD

FROM MAIN MEMORY. The diagram shows three words

PRESENTLY STORED IN CACHE

OF 15 BITS IS SHOWN AS A FIVE

THE ADDRESS VALUE OF

Direct ~~data~~ ~~data~~ number and its corresponding 12-bit word is shown as four digital data number.

- A CPU address of 15 bits is placed in the "argument" register and the associative memory is searched for a matching address. If the address is found, the corresponding 12-bit data is read and sent to CPU.

- If no match occurs the main memory is accessed for the word.

- The address-data pair is then transferred to the ASSOCIATIVE CACHE MEMORY. If the cache is full, the address-data pair must be displaced to make room for a pair that is needed and not presently in the cache.

- The decision as to which pair is replaced is determined from the replacement algorithm that the designer chooses from Cache

- A simple procedure is to replace cells of the cache in a round robin order whenever a new word is requested from the main memory. This constitutes a FIRST-IN-FIRST-OUT replacement policy.

→ How do we divide work to the processor?

RAM
↓
Jewel address

DIRECT MAPPING: normally RAM's are expensive compared to RAM because of the added logic associated with each cell. The possibility of using RAM for the cache is very low.

- The CPU address of 15 bits is divided into 2 fields: the nine least significant bits constitute the INDEX FIELD and the remaining 6 bits form the TAG field.

The figure shows that the memory needs an address that includes both the tag and the index field. The number of bits in the index field is equal to the number of address bits required to access cache memory.

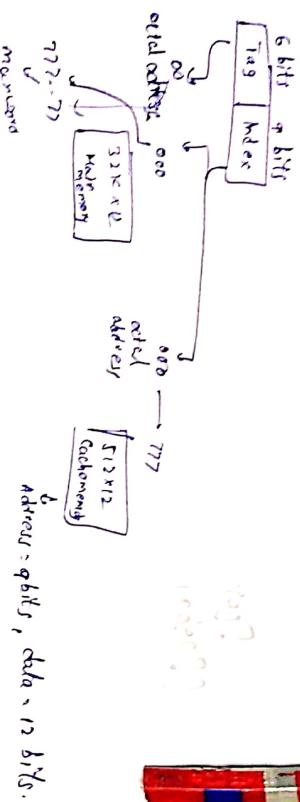
In the general case, there are 2^k words in cache memory and 2^n words in main memory ($k < n$).

The n-bit address is divided into 2 fields:

n-k bits for index field and k-bit for tag field.

The direct mapping cache organization uses the n-bit address to access the main memory and k-bit index to access the cache.

The internal organization of the word in cache memory is as shown in fig-12.13(c)



- Each word in cache consists of the data word and its associated tag.

- When a new word is first brought into cache, the tag bits are stored alongside the data bits.

- When CPU generates a memory request, the index field is used for the address to access the cache.

- The tag field of the CPU address is compared to with the tag in the word read from cache. If the two tags match, there is a hit and the desired data word is in cache.

→ If there is no match there is a MISS and the required word is read from main memory. It is then stored in cache together with new tag replacing the previous value.