

Lec 12 / 13

Data Path

- Collection of fundamental units such as ALU or multipliers.
- Perform data processing operations, a set path is used, this is called Data Path.

Definition :

- * So, a data path is a collection of functional units such as ALU, registers and buses.
- * Data path is a set of functional units that carry out data processing operations.

The Control Unit + Data Path are the two major components of the CPU.

Control Unit

It generates control signals and components work according to these components.

- Control variable \Rightarrow Name of ~~a~~ or a control signal is known as control variable.
- Control Word \Rightarrow Collection of all the control signals is called control word.

Control word is the signals generated by the CU to do a single microoperation.

Control Unit Organization

- Hardwired CU
- Microprogrammed CU

Hardwired CU

Control logic is implemented with gates, flip-flops, decoders and other digital circuits.

- Can be optimized for faster operation
- Hard to rearrange components.

A hardwired CU takes instruction as an input and a timing generator.

The timing generator will generate signals so that CU can identify which cycle it is in for a given instruction and will generate corresponding control word.

* Note : A single instruction can take multiple cycles / control words.

Microprogrammed Control Unit

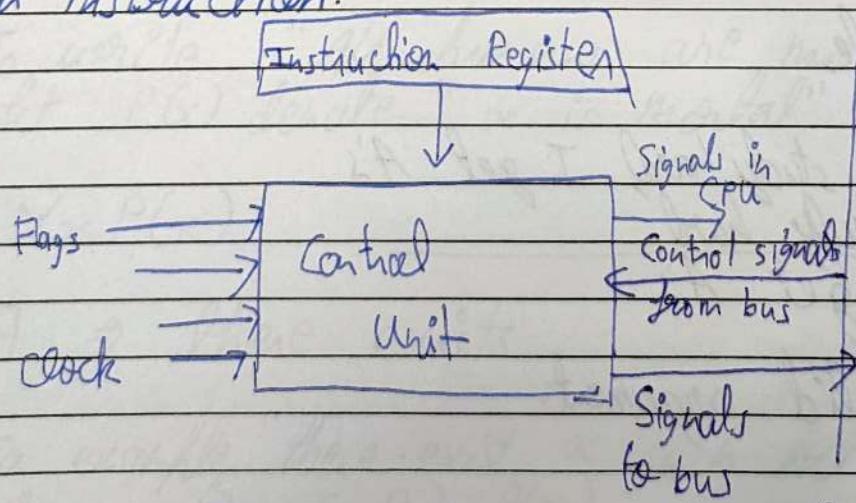
Rather than generating control word using hardware, keep all possible control words in some memory in form of micro-programs and send the corresponding control words when needed based on ALU input.

- Updating logic is easy.
- Slower than hard-wired CU

The memory to store memory where control words are stored is called control memory. Control memory is in CPU.

All the control words for a instruction are stored in a sequence in control memory.

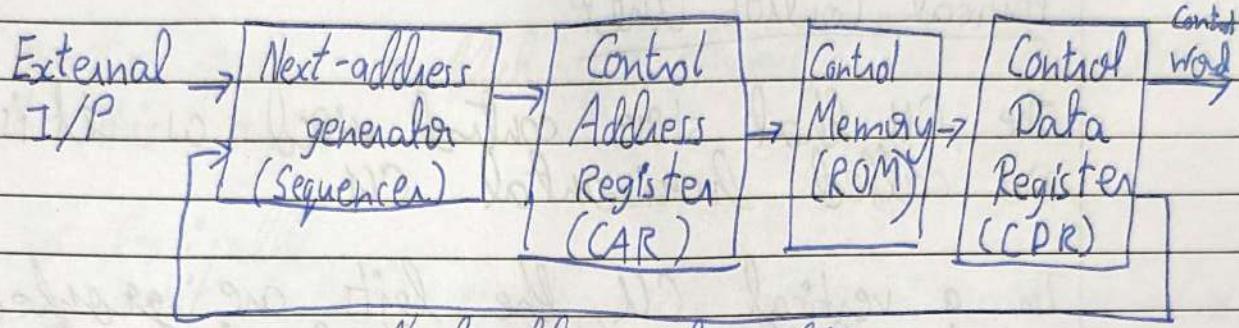
This makes it easier to generate words for an instruction.



A generic diagram for CU control bus

Lec 14

Control Word Sequencing



Next address information

On each location in control memory :

Control Word	Next address information
--------------	--------------------------

is called Micro-instructions.

The collection of micro-instructions is called micro-program.

The control unit will send control words in sequence until the next address information tells the sequencer that instruction has ended in which case external input is taken.

This following of sequence is called control word sequencing.

The sequencer uses MUX to decide what to use as next address is, so a micro instruction uses the following structure.

control signals (control word)	MUX Select	Next Address	Collection of control signals is called control word
-----------------------------------	---------------	--------------	--

The mux select will tell the MUX in sequence what to forward.

Vertical Control Unit

The CU that store control word as it is are called horizontal CU.

In a vertical CU, the bits are grouped such that we know that at a time there is only one bit high and others are zero.

We can group also in a way that almost one is zero.

Eg, if a bus is connected to a graph cell registers we know only one has IN^{right} signal 1 and rest have zero.

We use decoders to store these grouped signals.

We use n-bits to re

We can compress an n-bits signal group to $\text{ceil}(\log_2(n))$ bits.

And latter later has these compressed signals through a appropriate decoder.

If all signals in group can be 0.
we use $\lceil \log_2(n) + 1 \rceil$ bits in compressed format to account for that possibility.

The extra outputs of decoder are discarded in this case.

Q8 So basically,
a n-bit control group will have n combinations if one & only one bit is 1
and will have $n+1$ combinations if atleast 1 bit can be 1.

Ans And no. of bits in compressed signal are
 $\lceil \log_2(\text{no. of signal control words in group}) \rceil$
 $\lceil \log_2(\text{no. of combinations}) \rceil$

Lec - 15

In vertical CU if the ~~program~~ few signals cannot be grouped together, they are stored in a horizontal manner.

In terms of speed, from fastest to slowest

Hard wired CU, horizontal CU, vertical CU

~~Micro-Programmed CU~~

Nano Programmed Control Unit

In a micro-programmed CU we store the microinstructions in a sequence for easier sequencing.

But trying to maintain a sequence may ~~leads to~~ leads to a lot of repetition of microinst control signals.

And control signals are long so keeping them in sequence may be takes a lot of memory.

To reduce memory usage we use nano nano programmed CU.

All unique control signals are kept in a separate memory and in our micro instruction, we replace the

signal with the address for that signal
in our unique signals memory.

The main benefit is usage of less storage.

- * This is slower than microprogrammed CU.

RISC

(Reduced Instruction Set Computer)

- 1) Less number of instructions
- 2) Fixed length instructions
- 3) Simple Instructions
- 4) Limited addressing modes
- 5) Easy to implement hardwired
- 6) One cycle per instruction
- 7) More registers
- 8) Register to Register ALU only

CISC

(Complex Instruction set computer)

- More no of instructions
- Variable length instructions
- Complex instructions
- More & Complex Addressing modes
- Hard to implement hardwired
- Multiple cycles per instruction
- Less registers
- Any ALU input feasible.

Dec - 16

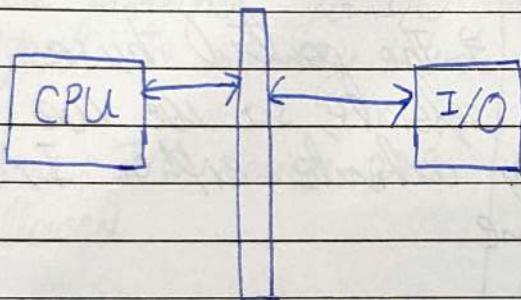
I/O Organization

Peripherals :

Other than main memory, all devices connected externally are known as peripheral devices.

Peripherals

- Input device
- Output device
- Storage devices



I/O interface
(Hardware / soft software)

Device drivers

I/O interface will help in signal conversion.
It also helps to synchronize the data signals sent between CPU and I/O.

I/O interface also helps in data format conversion between peripherals & CPU.
I/O interface also manages peripherals such

they don't interrupt CPU or other peripherals.

I/O vs Memory Buses

- Separate I/O and memory buses
- common address & data bus but different control buses
- Common bus for both I/O & memory.

We don't use completely separate buses because it is not cost effective.

Common data / address \Rightarrow The control buses are separate, so we use control bus to only activate either I/O or memory at a time.

Here we have separate address for I/O and memory.

This config is called I/O mapped I/O or port mapped I/O.

Common data, address & control buses \Rightarrow In this

config we take ~~any~~ address from memory and assign that to an I/O device.

~~Now when control bus gives signal,~~

Now when address bus ~~1~~ refers to a address related to I/O, the I/O device is accessed, otherwise memory is accessed

This config is called Memory mapped I/O

Memory Mapped I/O

- i) I/O do not have own address space
- ii) Some memory waste
- iii) Memory access instruction can be used to access I/O also
- iv) More instruction for I/O (equal to memory instructions)

I/O Mapped I/O

- i) I/O have their own address space
- ii) No memory waste
- iii) I/O and memory have different instructions.

Less instructions for I/O

- v) Only a limited number of I/O devices are allowed (depends on no. of ports).

Dec 17

Asynchronous Data Transfer

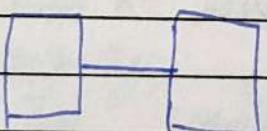
If two devices work on same clock, they are in synchronous operation. But in synchronous device clock has to be set on basis of the slowest component and hence causing bottleneck.

- ⇒ So we use asynchronous system where each component has its own separate clock.

So in async system, we need external sync.

The I/O interface will provide sync.

Serial Data Transfer

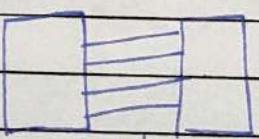


1-bit at a time



Cheaper

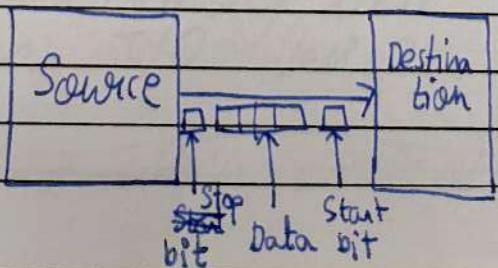
Parallel Data Transfer



n-bit at a time

where, $n = \text{no. of lines}$
Faster

Serial ~~but~~ asynchronous



To identify the start and end of data transfer, we use start bit and stop bit

- * First a start bit is sent,
- * Then ~~a~~ the data is sent based on the memory ^{word} size (byte addressable/ word addressable)
- * Then stop bit is sent.

The start & stop bits are used to keep source & destination in sync.

Note : the memory is typically byte addressable so 8 bits of data are sent at a time.

Consider memory to be byte addressable unless specified.

Mode of transfer

- Programmed I/O
- Interrupt Initiated I/O
- Direct Memory Access (DMA)

Programmed I/O (Polling)

- There is no free busision through which I/O can inform CPU of data transfer.
- I/O sets status and waits
- CPU runs program periodically to check

status of each device

- If any device has status set, the CPU performs data transfer on it.

To set the status, a flag register is used in I/O interface

Interrupt Initiated I/O

- I/O has provision (interrupt signal) to inform CPU about data transfer
- When CPU receives interrupt,
 - * It completes current instruction
 - * Saves status (PC, PSW etc) of current process on stack
 - * Branches to service interrupt
 - * Resumes process by taking value from stack.

Dec 18

Interrupt Service Routine

The program which is run, in response to a interrupt is called ISR.

Different interrupts will have different ISR.

Vectorized Interrupt

The device which sends interrupt signal will also send a vector address which is the address address to the proper ISR to that interrupt.

Non-vectorized Interrupt (Scalar interrupt)

The device only send interrupt signal.

The CPU will use Default Service Routine.

The Default service routine will tell CPU the ISR address.

ISR's are a part of the device drivers.

Maskable interrupt

CPU can accept or reject the interrupt. CPU can also keep it in pending (based on CPU design).

Nonmaskable Interrupt

CPU cannot reject these interrupts.

Internal interrupts

When interrupt occurs due to unexpected error during instruction execution.

Eg, page fault, division by zero.

External interrupts

Interrupts received by system external reasons. (i.e. not by the currently processing instruction error) by \neq I/O devices mostly

★ Internal / External interrupt are also called software / hardware interrupts, respectively

Simultaneous Interrupts

Interrupt from higher priority device will be serviced first.

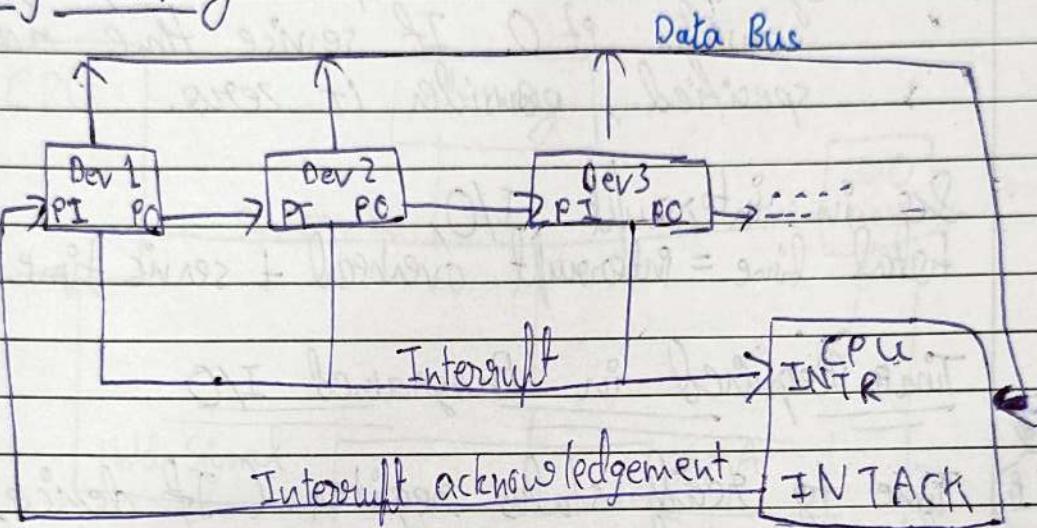
- Software solution
- Hardware solution

Software solution : When receiving a simultaneous interrupt, a software is executed to decide device with higher priority. (It uses polling)

Hardware Solution

- Serial (Daisy chaining)
- Parallel (Not in Syllabus 😊)

Daisy chaining



- * If interrupt is on from a device, signal is received in INTR
- * CPU generates acknowledgement and send it to serial chain
- * In priority Dev 1 > Dev 2 > Dev 3
- * If the dev was causing interrupt receives signal acknowledgement, it can return its vector address
- * Else device forwards acknowledgement
- * The INTR pin in CPU will remain enabled until all interrupts are resolved.

Usually, devices which is nearest to CPU is given highest priority.

Time Required in Interrupt I/O

- i) Interrupt overhead (completing current inst, storing in stack etc)
- ii) Service time (remaining all).

Note: if overhead not specified in ques, consider it 0. If service time not specified, consider it zero.

So in interrupt I/O,
total time = interrupt overhead + servc time.

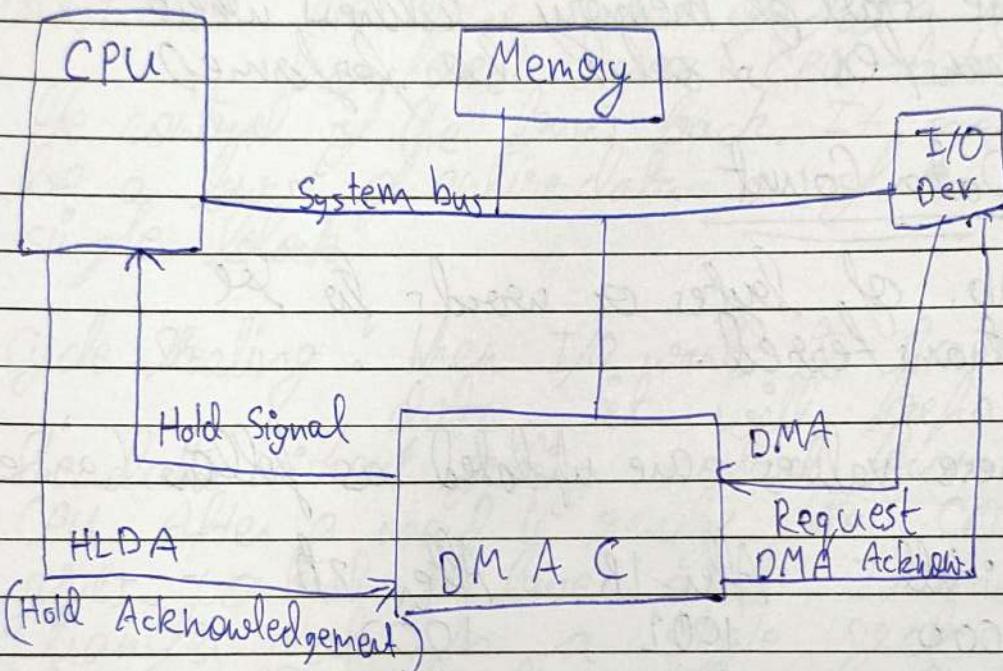
Time Required in Programmed I/O

- i) Time to read status register. (If device transfer rate is given, take size of register in consideration. If size not specified consider what size-addressable / transfer size of machine is using, else 1 byte)
- ii) Transfer time between components

Lec - 19

DMA I/O

A hardware named DMA Controller is used to transfer data without CPU.



CPU sends two data to DMAC before starting transfer,

- ① Starting address
- ② Data count

The complete process is

- 1) I/O device sends DMA Request to DMAC
- 2) DMAC sends Hold signal to CPU
- 3) CPU sends data to DMAC using system bus
- 4) CPU sends HLDA signal to DMAC
- 5) DMAC sends DMA Acknowledgment to I/O

⑤ Data transfer between device and memory using system bus.

I Starting Address

The start of memory address where transfer should be performed

II Data Count

No. of bytes or words to be transferred.

These values are updated as follows (an example)

	Initial	After 1B	After 2B	---
S. A.	1000	1001	1002	-
D. C	500	499	498	0

During DMA Transfer, CPU can perform only those operation which do not require system bus. so mostly CPU is blocked

If CPU is block for too long, performance will suffer.

So we use DMA Transfer modes to solve this issue

Modes of DMA Transfer

- Burst Mode
- Cycle Stealing
- Interleaving DMA

Burst Mode : The entire burst of data is transferred before CPU takes the control of the bus back. It can be a burst of entire data or of a single block.

Cycle Stealing : When I/O wants to transfer data, it will prepare the data without getting control from CPU. After a word is ready, the CPU gives control of buses and I/O will transfer data in a single cycle (I/O cycle) and then return control.

If time required to prepare data = t_x
 " transfer data = t_y

% of time CPU blocked = $\frac{t_y}{t_x + t_y} \times 100$ (Burst mode)

If transfer time is overlapped with fetch/decoder time

% of time CPU is blocked = $\frac{t_y}{t_x} \times 100$ (Cycle stealing)

Lec - 20

Interleaving DMA

Only when CPU is not using system bus, DMAC will use the system bus.

CPU will not be blocked due to DMA.

- * DMAC is a special purpose processor in this, as it can control data transfer between memory & I/O (can generate address & control signals for memory). (DMAC is often referred to as special purpose processor).
- * DMA service can be done independent of, CPU instruction state, so unlike interrupt service, it can be done in the middle of an instruction.
 - ⇒ If data size is very small, we use interrupt I/O
 - ⇒ If data size is very large, we use DMA.

⇒ DMAC has 3 registers

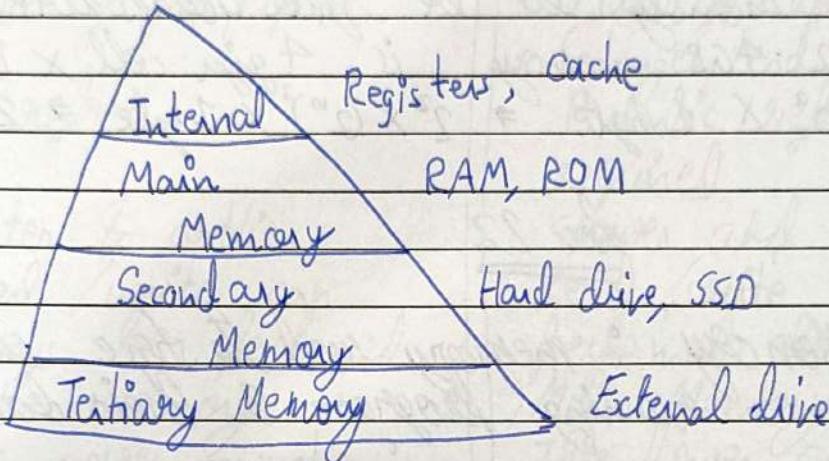
- * Address Register - Address to destination in memory
- * Word Count Register - The no. of words to be transferred
- * Control Register - Specifies the transfer mode (Interleaving, cycle stealing, Burst)

Lec - 21

Memory \Rightarrow Physical device, used for storage.

Goal of Memory hierarchy

1. Reduce access time
2. Reduce cost per bit



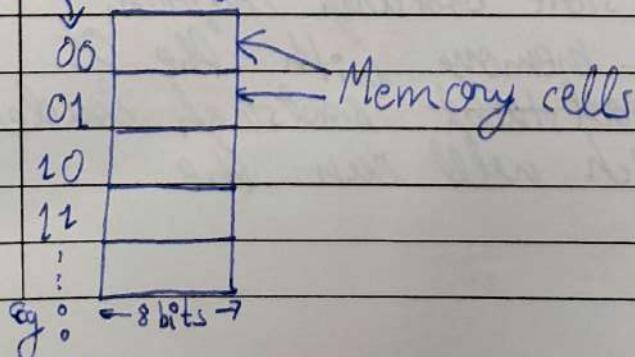
The goal is to have system access memory more on the higher level (top level).

Average cost of memory = $\frac{\text{Cost of total memory in system}}{\text{size of total memory in system}}$

Memory Representation

Memory is represented by

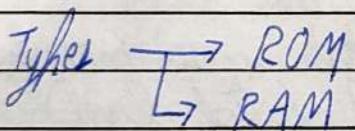
$$\text{cell no./ addresses} \xrightarrow{\text{capacity}} \text{No. of cells} \times 1 \text{ cell capacity}$$



- ⇒ So a 8×16 bits memory has 8 cells each of size 16 bits.
- ⇒ Each cell has a unique number we call memory address.
- ⇒ 128×8 bit can be written as 128×1 byte or 128 bytes. If not specified, we can consider memory to be byte addressable.
- ⇒ So a 4GB memory is $4 \text{ giga cell} \times 1 \text{ byte} \Rightarrow 2^2 \times 10^9 \times 1 \text{ byte} \Rightarrow 2^2 \times (2^10)^3 \times 1 \text{ byte} \Rightarrow 2^{32} \times 1 \text{ byte}$

Lec 22

Main Memory : memory used to store current running programs & their data.



ROM ⇒ It stores initial ~~operation~~ instructions when CPU is started. The first task function of these instructions is

- i) Hardware check (Power On Self Test) (POST)
- ii) Booting (Bootstrap Program performs booting)

* ⇒ We can also store Bootstrap Program (BSP) in secondary memory with the OS and store a ~~Bootstrap~~ Bootstrap Loader in our ROM which will run the BSP.

RAM \Rightarrow used to store current running programs.

- Static RAM
- Dynamic RAM

Static RAM

- i) Made of flip flops
- ii) No refresh
- iii) Faster & costlier
- iv) Read / write can always be performed
- v) Cache memory
- vi) Idle power consumption is less compared to DRAM
- vii) Operational power is more compared to DRAM

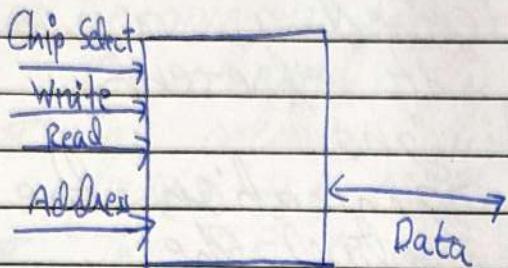
Dynamic RAM

Made of capacitors
(charges tend to discharge)
Periodic refresh is required

slower and cheaper
Read / write cannot be performed during refresh
Main Memory
Idle power consumption is more compared to SRAM
Operational power is less compared to SRAM

Lec 23

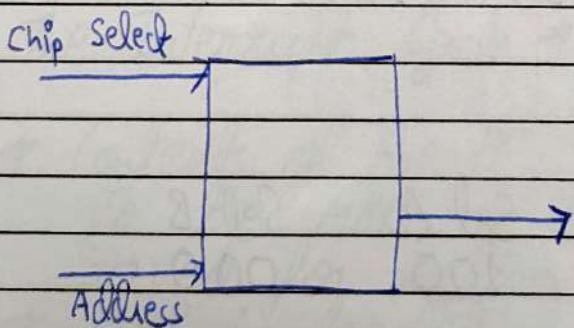
RAM Chip



chip select \Rightarrow Use by CPU to select either RAM chip or ROM chip.

CS	RD	WR	Operation
0	X	X	no op.
1	0	0	no op.
1	0	1	Write
1	1	X	Read

ROM Chip



CS	Operation
0	no op.
1	Read

Multiple Chip Support

Assume we have two chips in RAM

RAM 0	\Rightarrow	128 X 8 bits	address	7 bits
RAM 1	\Rightarrow	128 X 8 bits	address	7 bits
Total	\Rightarrow	256 X 8 bits	address	8 bits

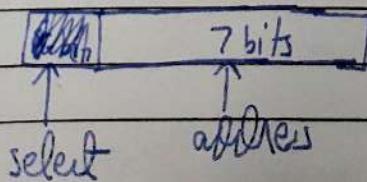
So we need 7+7 bits for both,
but CPU only generates address
for the total memory. Hence only
8 bits.

So we use the MSB of 8 bit address
to split them to two 7 bit signals.

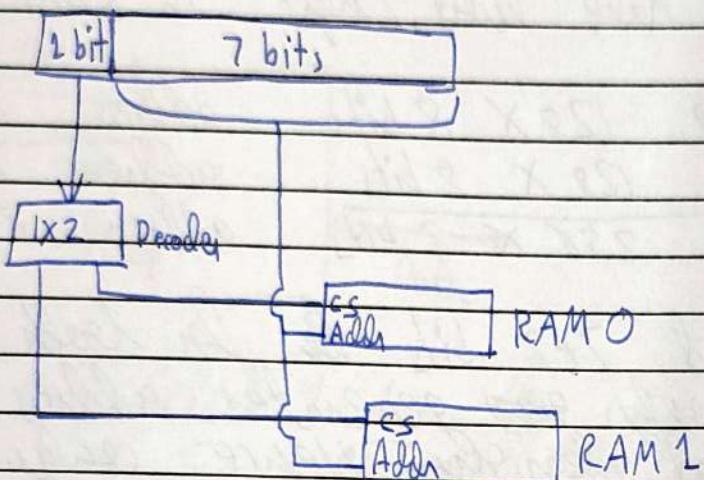
They are divided as follows

Decimal	Binary
0	00000000
1	: :
127	01111111
128	10000000
1	:
255	1111111

So for CPU generated address in our examp



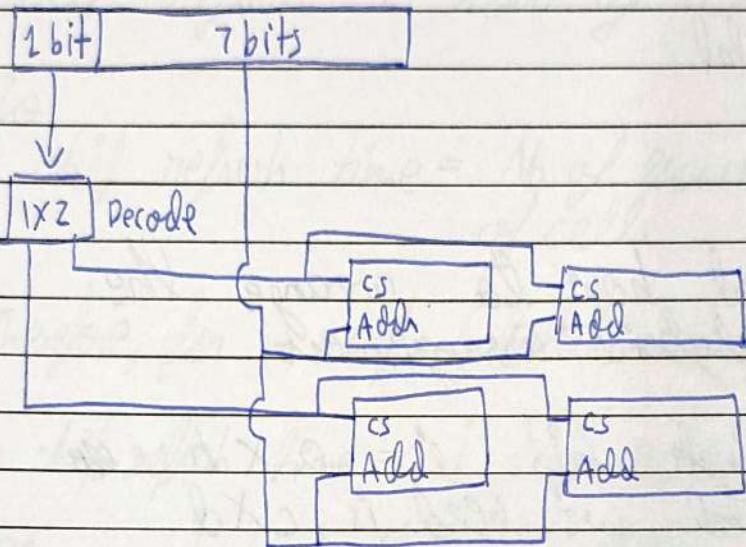
These can be drawn in hardware as CPU address



Similarly we can expand this idea for more chips.

- ⇒ So when we keep the amount of total addressable data same and increase no of cells, it is called vertical arrangement.
- ⇒ We could also rather keep number of cells constant and increase addressable size, this is called horizontal arrangement
- ⇒ Using both together is called hybrid arrangement.

Eg, ~~256X 256X 16 bits~~ using ~~128X 8 bits~~



Lec 24

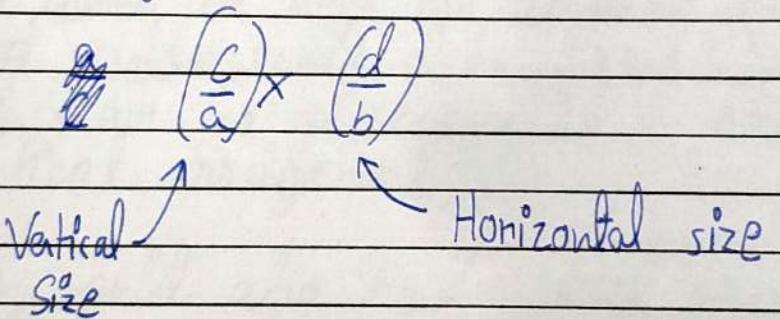
Multiple chip Support

- Horizontal
- Vertical
- Hybrid

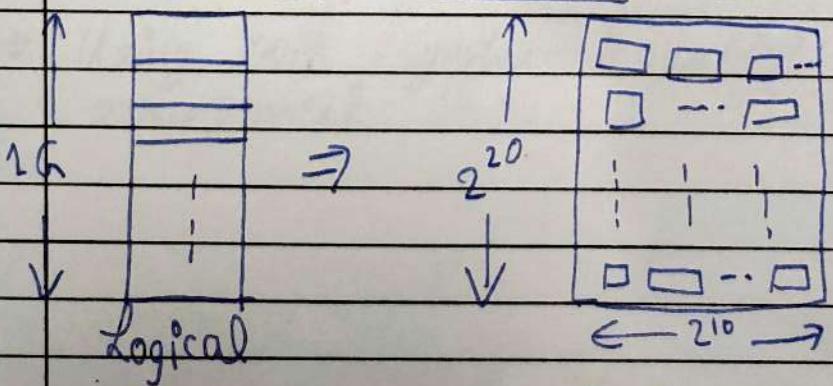
If we want how to arrange the chips in hybrid arrangement

If a single chip is $a \times b$ and total need is $c \times d$

The arrangement is



DRAM Chip Cell Arrangement



Practical arrangement

DRAM Chip Refresh

In a DRAM chip, one refresh operation will refresh 1 row of cells.

So

$$1 \text{ chip refresh time} = \text{No of rows} \times 1 \text{ refresh operation}$$

of cells

Therefore, for n chips (Total memory)

$$n \text{ chip refresh time} = 1 \text{ chip refresh time}$$

Since, refresh happens parallelly

Note : We are referring to rows as in a row in the practical design of a chip as shown in previous page.