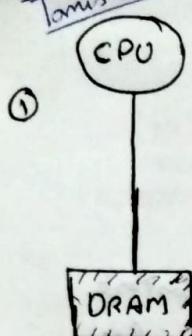
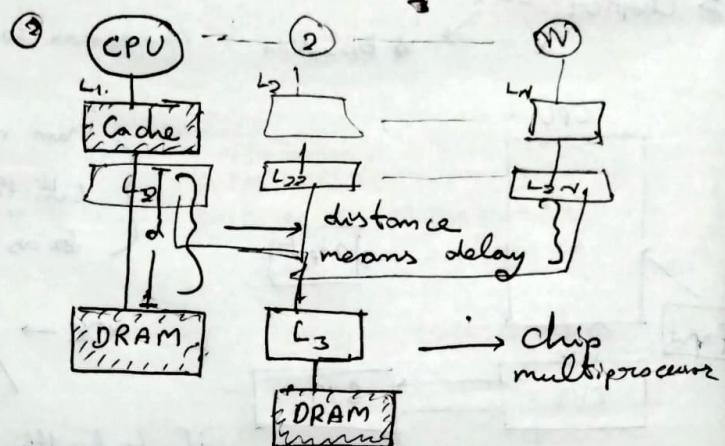


Tariq Ahmed

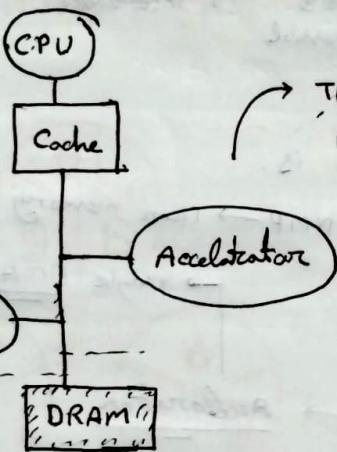


DRAM

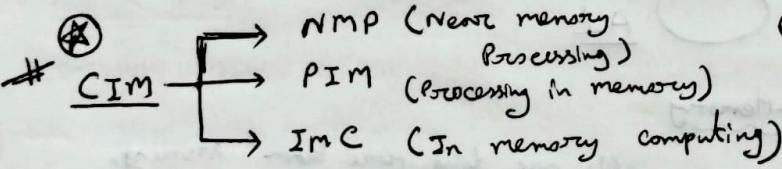
↳ When power is on, data can still decay. ③



④



These accelerators will produce high bandwidth

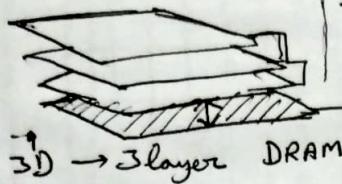


CIM is divided into three classes.

Requirements for CIM:

- ① design & compute capable memory & controller.
- ② design a processor chip with in-memory unit.
- ③ design of system & hardware interfaces.
- ④ Fixing of system softwares & languages.
- ⑤ Introduce memory within the CPU.
- ⑥ developing suitable algorithm for CIM
- ⑦ Develop a new memory architecture.

## DNNMP :

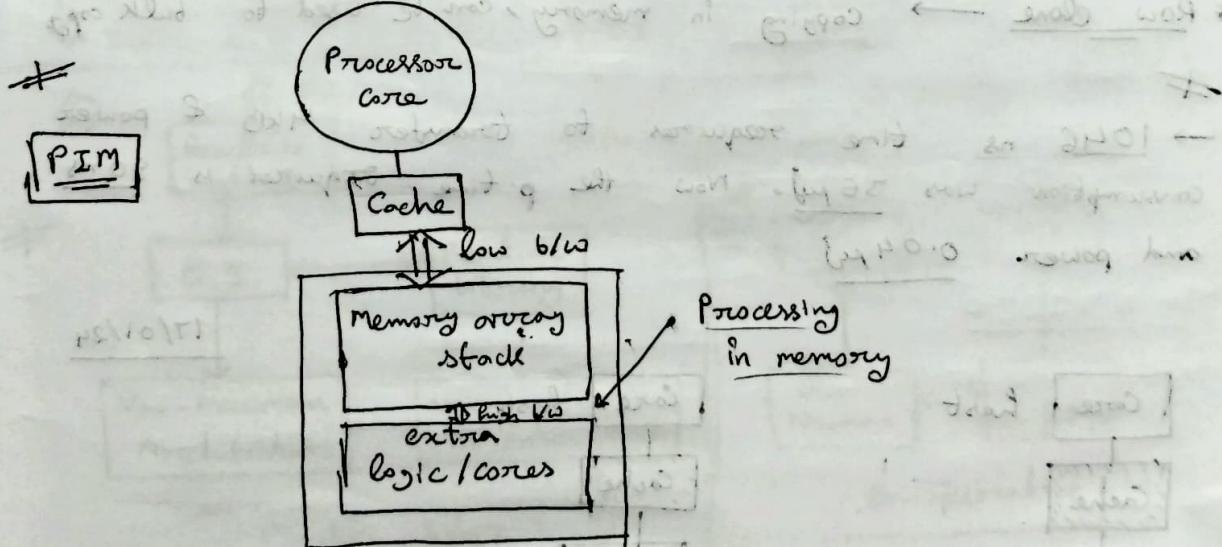
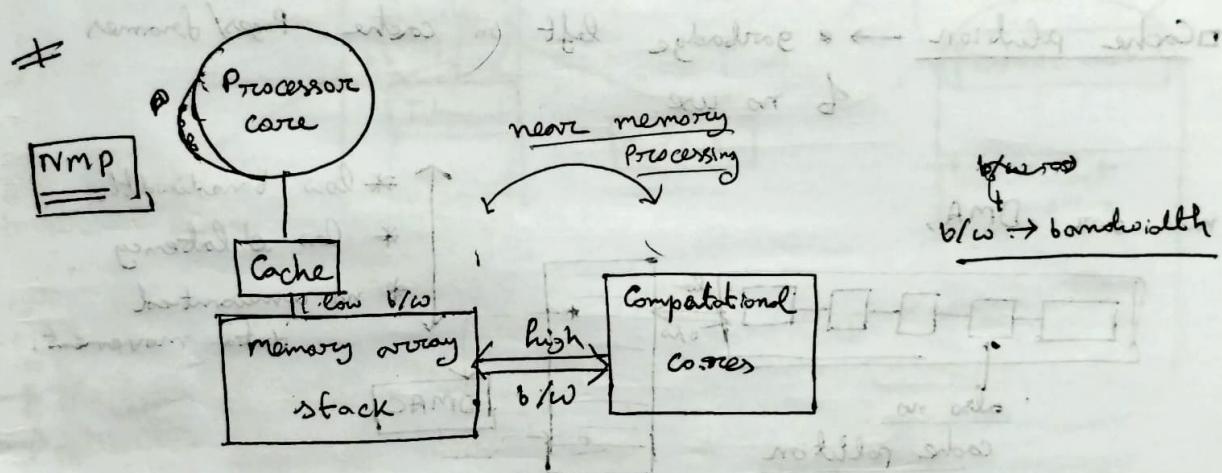


These layers are connected via TVS  
→ another layer which have a processing core  
↓  
This is part of (PIM)

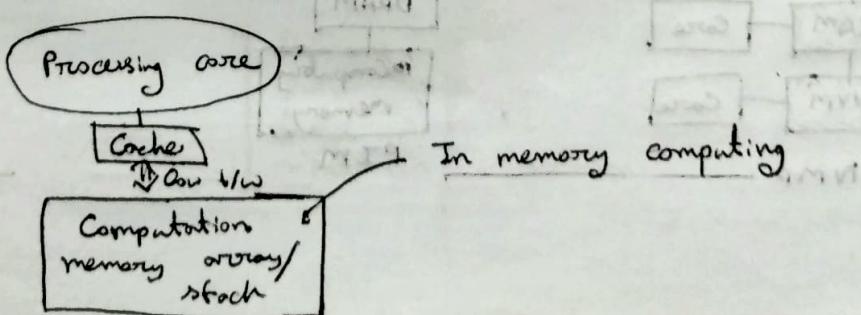
2.5D DRAM  
3D DRAM

→ NVM: The computation will be lost as memory will be gone if it is a volatile memory.

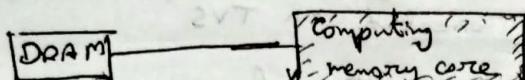
so, IMC has two sides      → Using DRAM  
                                  → using storage



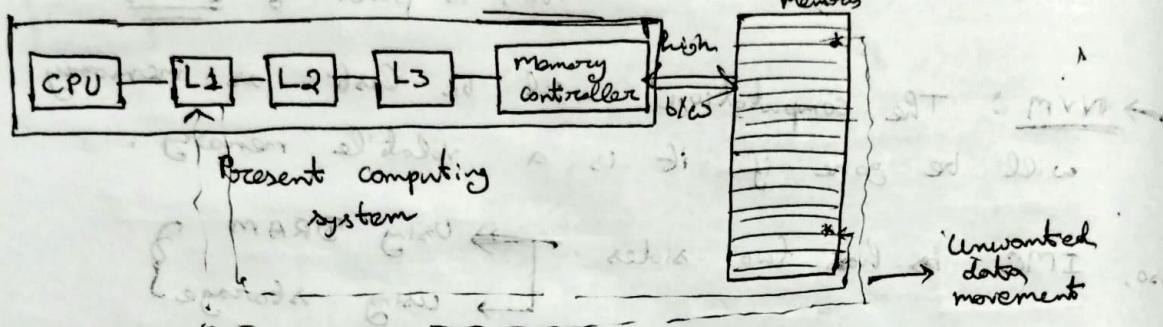
#  
IMC



CIM

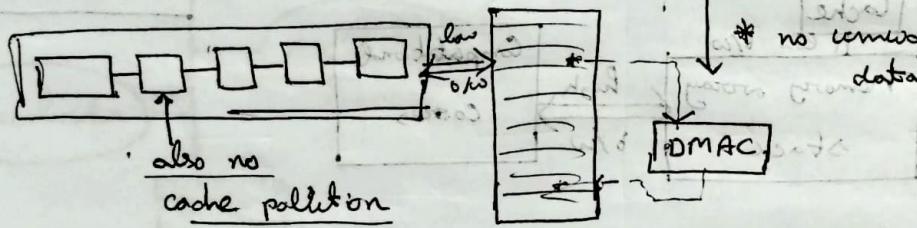


▷ DMA can be considered as near memory processing.



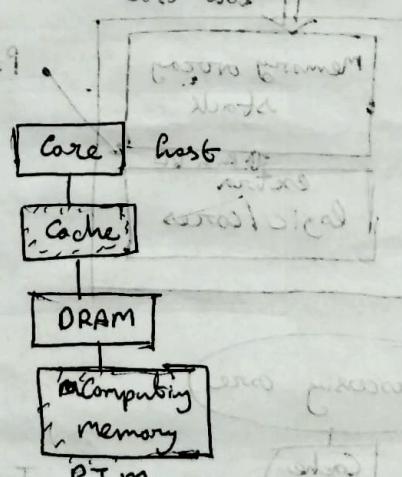
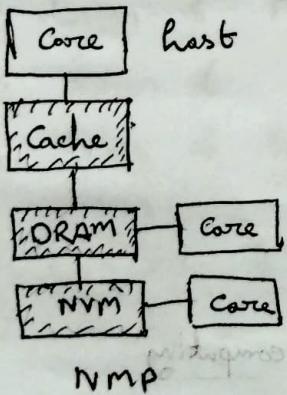
□ Cache pollution → • garbage left on cache. Pages becomes & no use.

now for DMA,



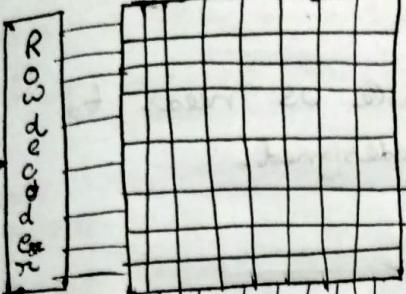
\* Row clone → copying in memory, can be used to bulk copy.

→ 1046 ns time requires to transfer 4KB & power consumption was 36 μJ. Now the p-time required is 90 ns and power 0.04 μJ



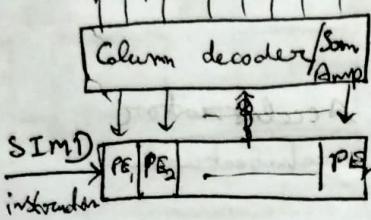
17/01/24

PIM

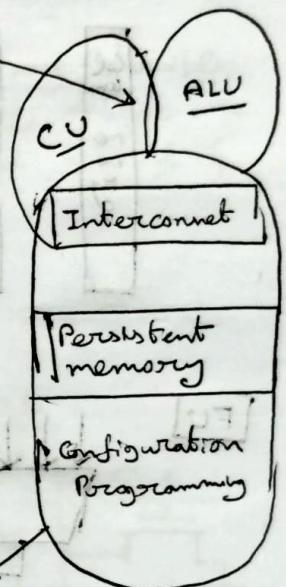
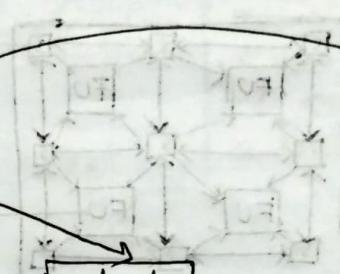
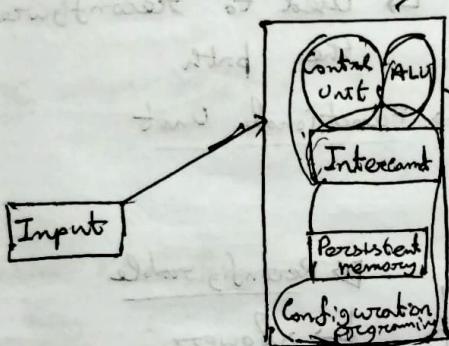


{ Error correction scheme

As this is within a chip there's less possibility of error during data transfer.



Cell is not computing just this SIMD processor.



→ The types of realisation with CIM & V-N

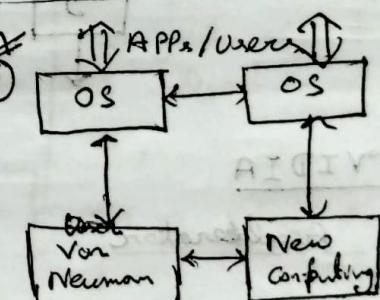
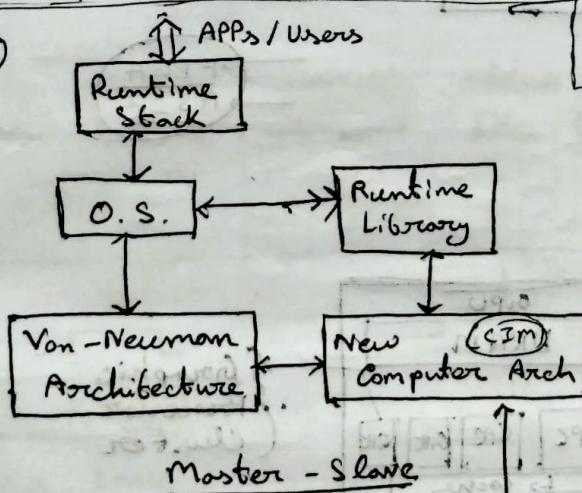
→ Master-Slave (i)

→ Co-operative (ii)

→ Independent (The picture above) (iii)

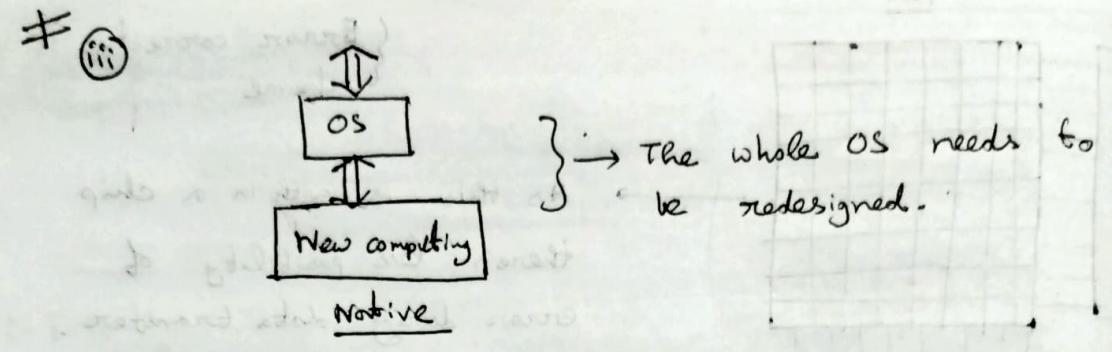
CIM without interaction with Von-Neumann Machine:

↳ not really feasible.



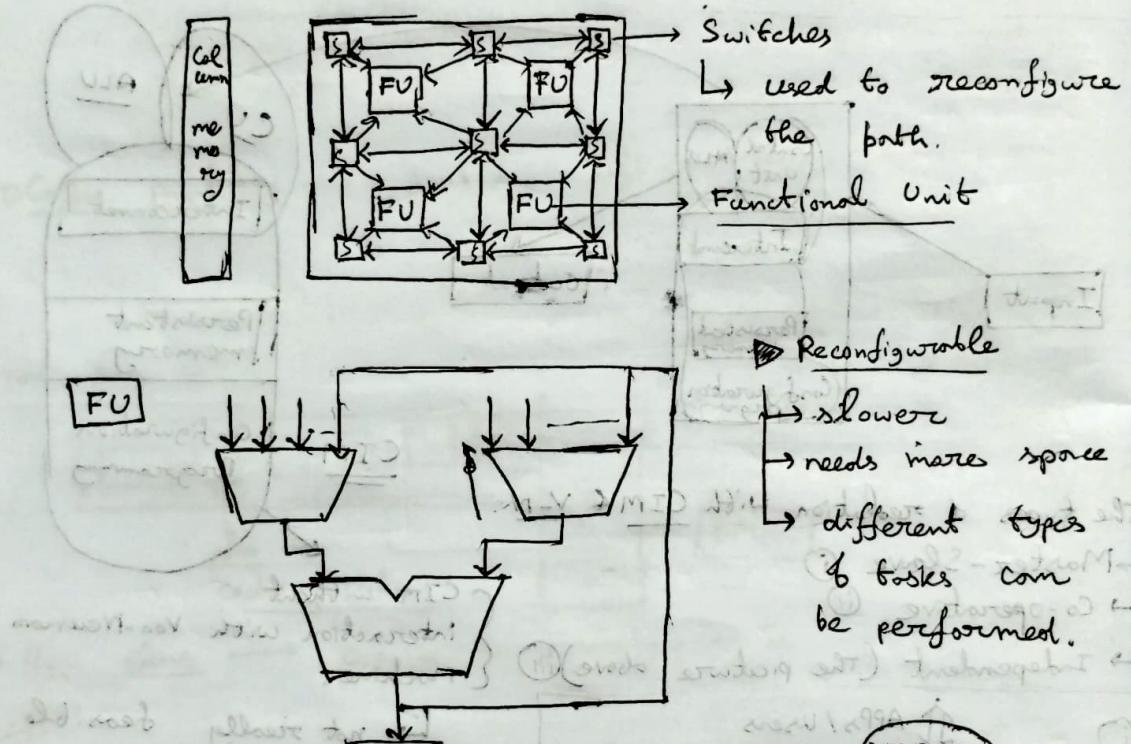
As a part of Von-Neumann

↳ need reorganisation



→ NMP:

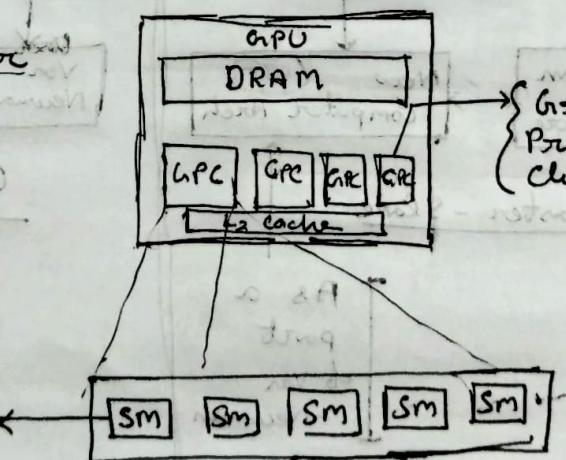
CGRA → Coarse Grain Reconfigurable Accelerator Architecture



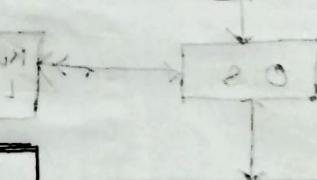
→ NVIDIA

GPU Accelerator

{ series of multiprocessor

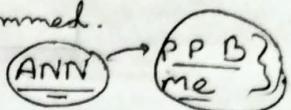


FPGA  
ASIC



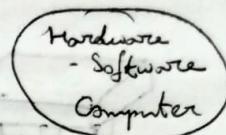
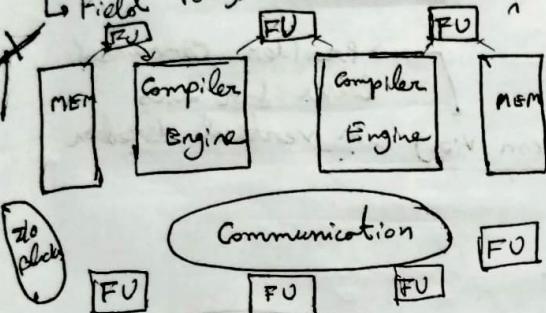
FPGA → Reprogrammable designs

↳ Simulations are run and re-programmed.



~~DETAILED~~ ▷ FPGA-based accelerator: Array.

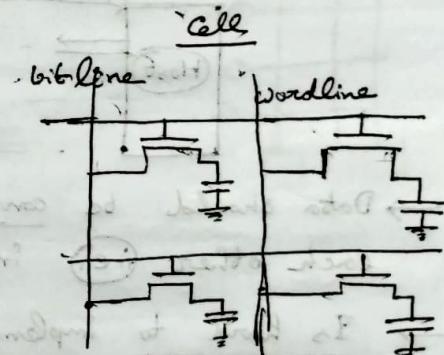
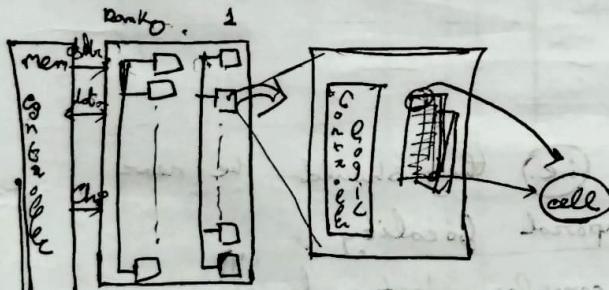
Field Programmable Gate Accelerators need processors.



[FPGA] designs are flexible and many designs are centered around it.

HMC → Hybrid memory cube  
(Designed from DRAM)

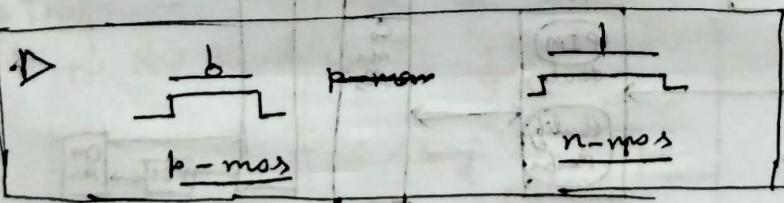
▷ Basic structure of DRAM



The wordline is basically address line

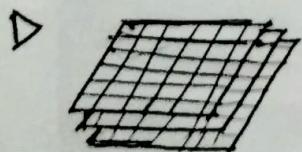
The bitline is data line

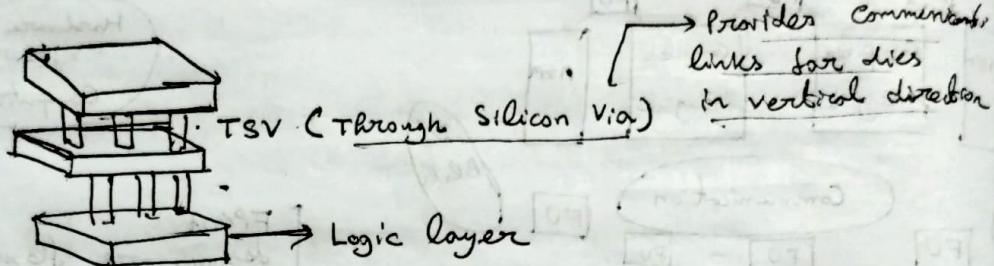
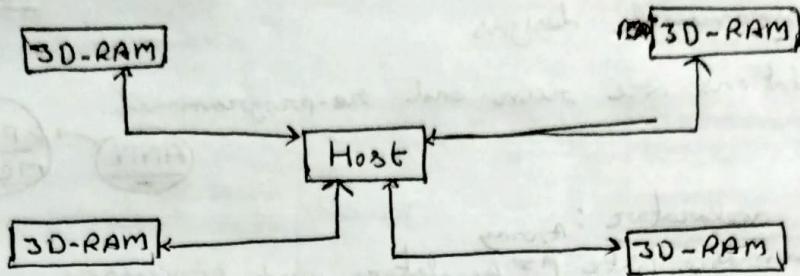
if capacitors are charged the it is '1'. Discharged means '0'



SA → Sense Amplifier  
↓  
considers difference b/w bitline & b-line

This basic structure of DRAM will not be used in our subject, we will use HMC



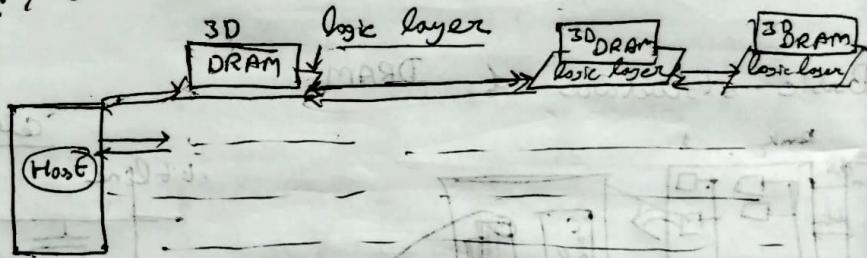


► TOP - PIM:

Map Reduce Workload

→ Daisy-chained:

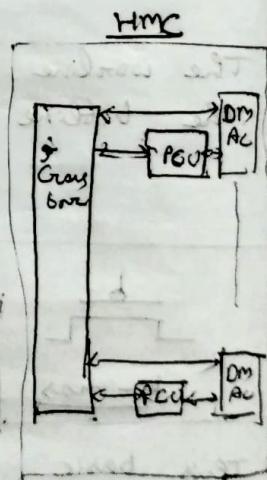
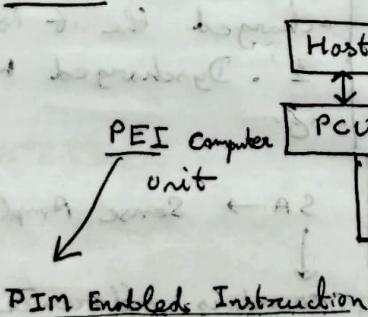
(Map) Daisy chained (Series)



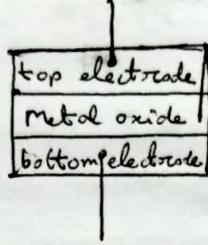
→ Data should be congruous i.e. they should be near each other i.e. in temporal locality.

→ Is hard to implement complex instructions.

→ ISA?



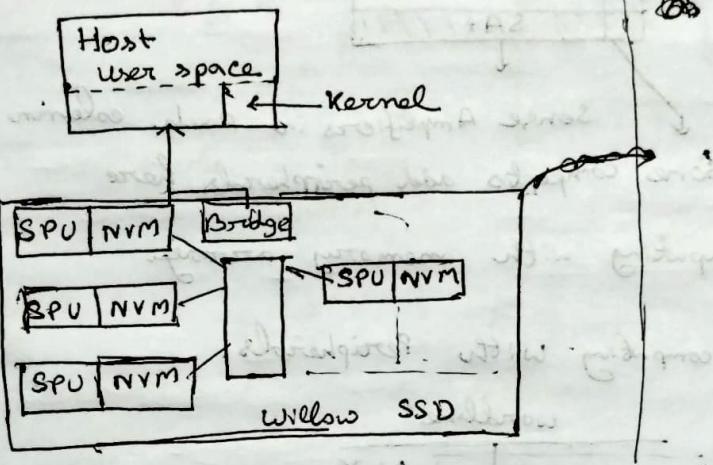
## Re-RAM → Resistive RAM



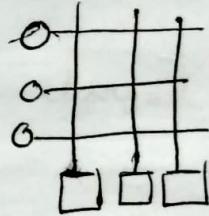
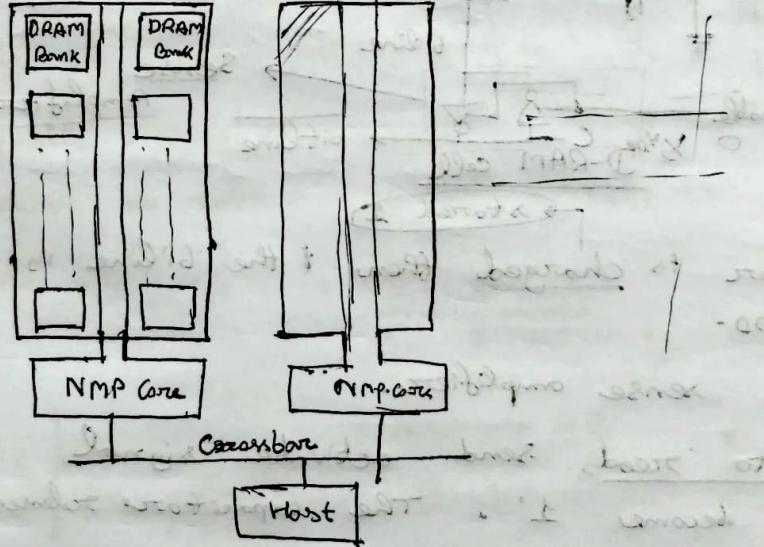
Memristor

key component  
in IMC

implements

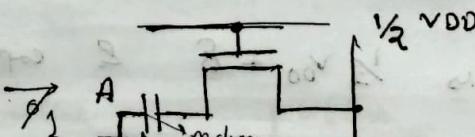


## NMP Vault



## Ambit

TRA  
(Triple Row Activation)



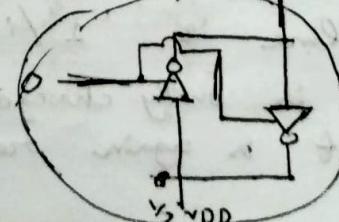
$$BC \rightarrow \frac{1}{2} VDD$$

$$A \rightarrow \frac{1}{2} VDD$$

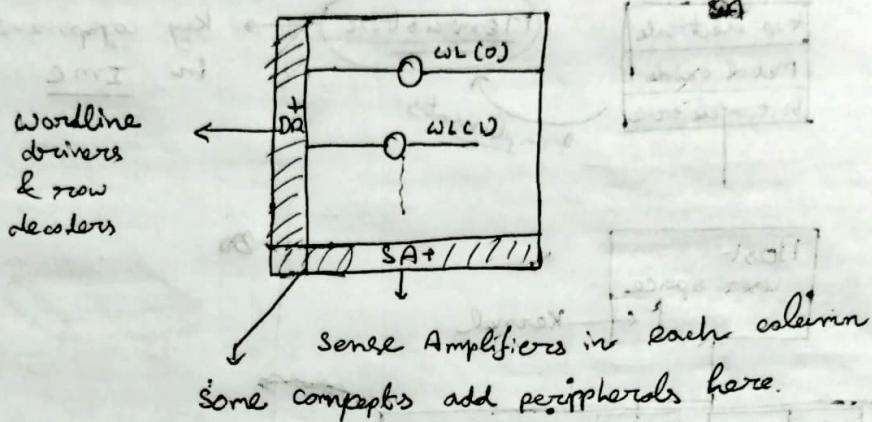
charge sharing

If we control ABC & the below input  
& we can use different  
logics.

$$\text{e.g. } AB + BC + CA$$



D-RAM is a series of transistors & capacitors.

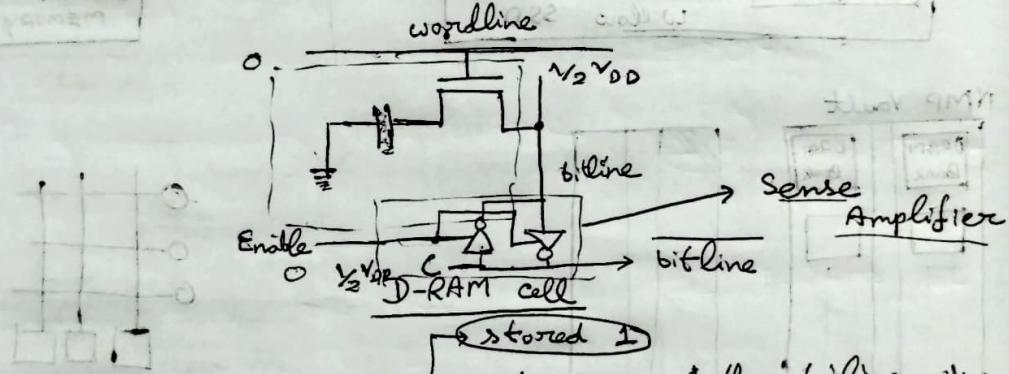


IM-A

↳ In memory computing with memory average

IM-P

↳ In memory computing with Peripherals.



- If the capacitor is charged then the bitline is positive  $\frac{V}{2} DD$ .
- Enable controls sense amplifier.
- If we want to read, send activate signal, wordline will become '1'. The capacitor releases charge.

bitline charge is  $\frac{V}{2} DD + 8$  & capacitor is  $\frac{V}{2} DD + 8$ .

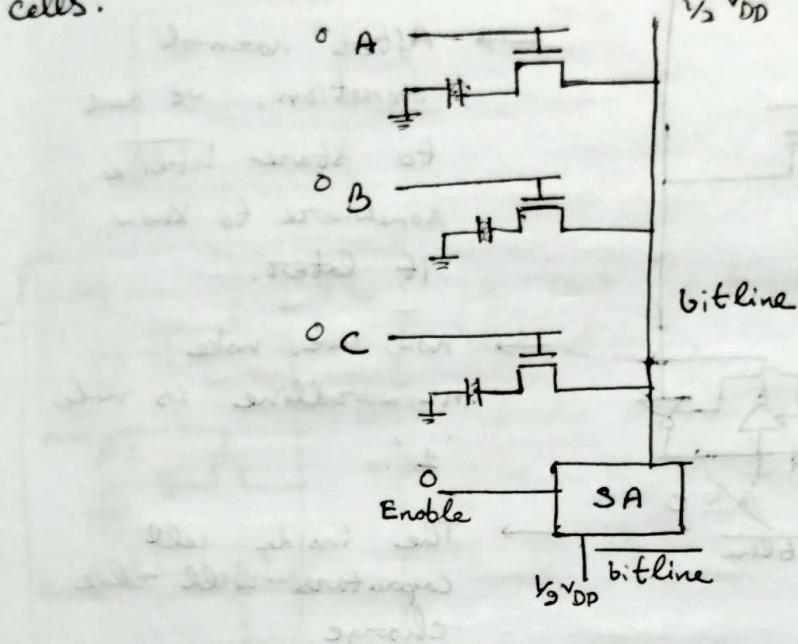
- Now, the enable is '1', The twin inverse will be in a feedback loop. Now, the ~~bit~~ bitline goes to ' $V_{DD}$ ' & C point will go to '0'.

→ Hence, the ~~bit~~ bitline is '1' i.e.  $V_{DD}$ .

Also the capacitor is fully charged again. Hence the data we lost is again recovered.

## ► TRA :

A simple AND/OR logic was implemented using three cells.



~~SA~~  
Lo Same as  
last picture

$$\delta V = \frac{(2K-3)C_c}{6C_c + 2C_b} V_{DD}$$

$C_c \rightarrow$  cell capacitance

$C_b \rightarrow$  bitline capacitance

$K \rightarrow$  no. of cells in a fully charged state.

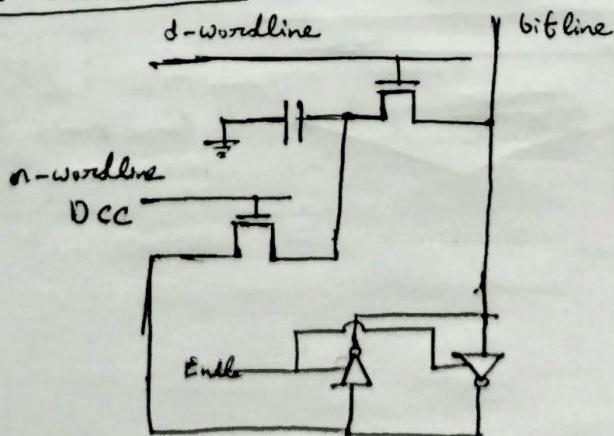
→ So, if  $K=2$  or  $3$ , then the ~~if  $B \rightarrow +ve$  from expression~~ bitline will make '1'

→ Enable = 0,  $A=B=C=0$ , Capacitors are charged this ~~is~~ data stored state.

→ So, bitline =  $\frac{AB+CA+BC}{3}$   
So, if majority is one.

$$\begin{cases} \text{if, } C=0, & \begin{array}{l} \text{AND} \rightarrow AB \\ \text{OR} \rightarrow A+B \end{array} \\ \text{if, } C=1 & \end{cases} \quad \left. \right\}$$

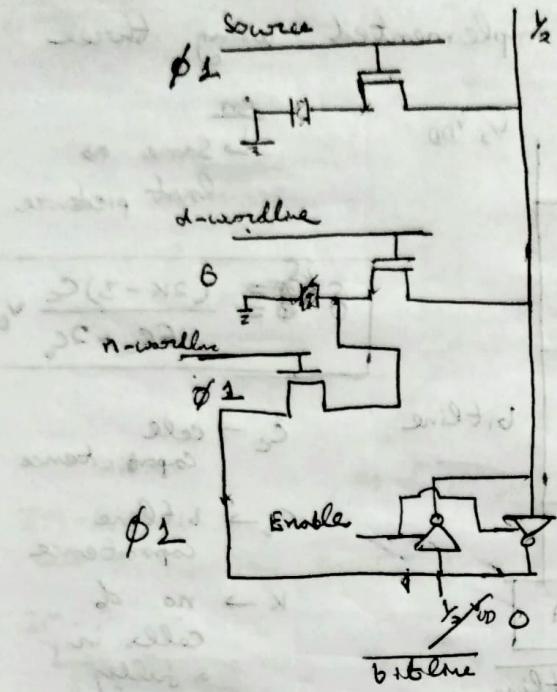
## ► Dual-Connect Cell :



d-wordline  
n-wordline

d → data  
n → negation

now,



bitline

$\frac{1}{2} V_{DD}$  +  $\frac{V}{2}$

1

signal voltage depends on

1. After normal operation, we need to store bitline somewhere to know if later.

now; we make  
n-wordline is node  
'1'.

The inside cell capacitors will store charge.

so, '0' represents no charge }  
'1' charged }      **I/O** in bitline

Through this, we can make the DRAM cell compute simple logic operation.

$$DA + DA + DA = \text{and} \text{ed}$$

and of course it is

$$\left\{ \begin{array}{l} DA \leftarrow D(A) \\ D(A) \leftarrow A \end{array} \right. \quad \left| \begin{array}{l} 0=0 \rightarrow 1 \\ 1=1 \end{array} \right.$$

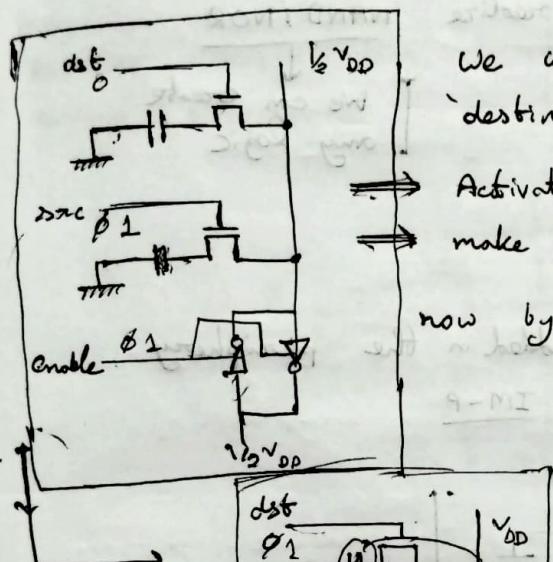
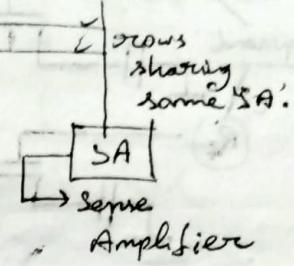
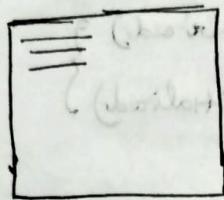
subword address  
interior  
content  
contents

and

1. terminal - load  
without



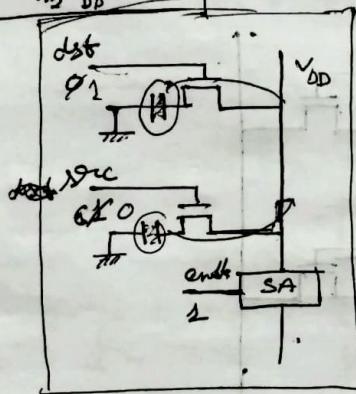
## Bulkcopy / Initialization:



We want to transfer 'source' to 'destination'.

Activate source  $\Rightarrow (0 \rightarrow 1)$   
make enable  $\Rightarrow 1$

now by ' $\frac{1}{2} V_{DD} + 8$ ' & 'SA'



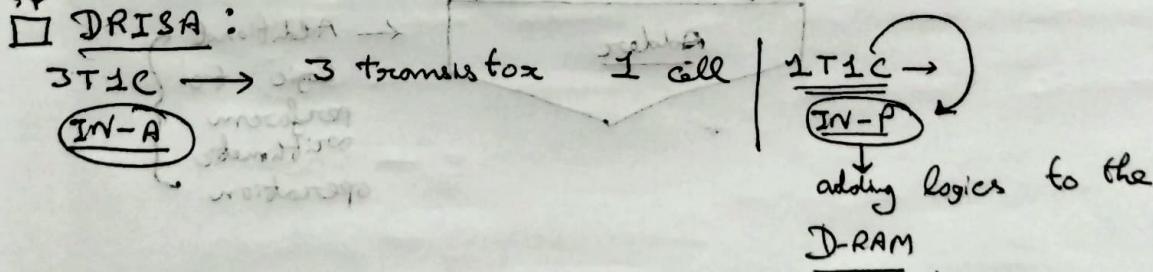
$\Rightarrow$  deactivating 'source' ( $1 \rightarrow 0$ )

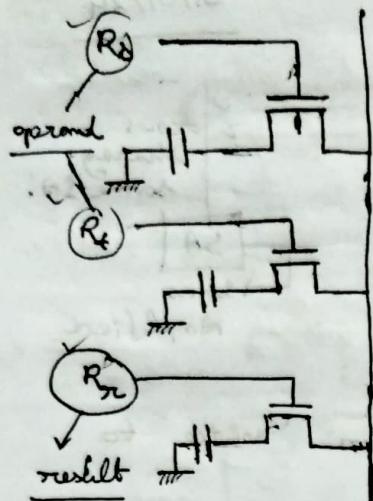
$\Rightarrow$  activating 'destination' ( $0 \rightarrow 1$ )

now, charge is shared to destination.

Hence, source is copied to destination.

- $\Rightarrow$  In reality, charge sharing is not 100% error proof.
- $\Rightarrow$  If we activate one row, the whole subarray will be copied. We cannot copy a part of subarray.





$$\rightarrow R_3 R_4 + R_4 R_n + R_n R_3$$

Adder

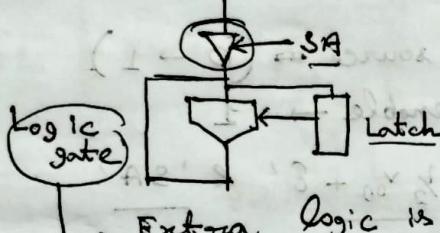
$$\begin{aligned} &\text{AND} \\ &R_3 = 0 \text{ (initialized)} \\ &\text{OR} \\ &R_n = 1 \text{ (initialized)} \end{aligned}$$

using extra logic,

we can realize

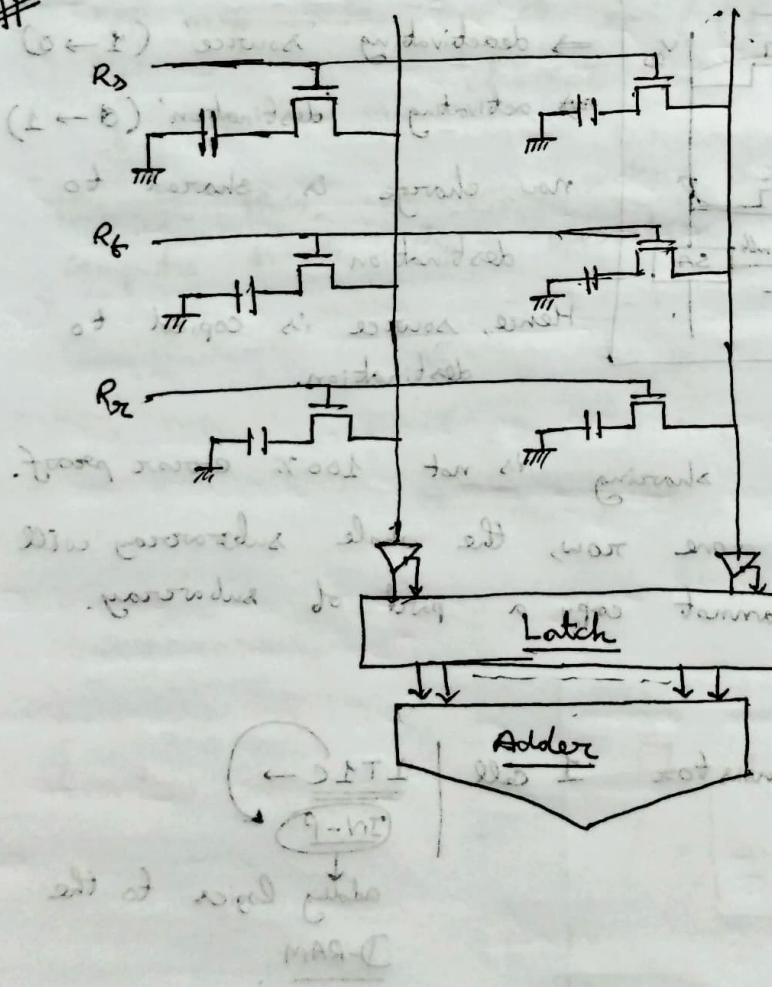
NAND / NOR

We can create  
any logic



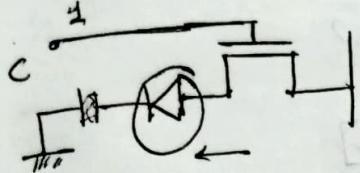
Extra logic is added in the periphery  
hence, it is called IM-P

#



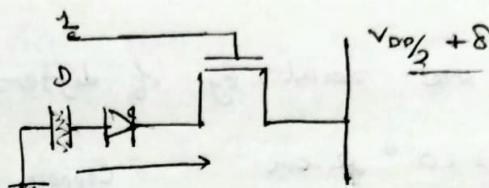
Adder  
Logic to  
perform  
arithmetic  
operation

► ROC :

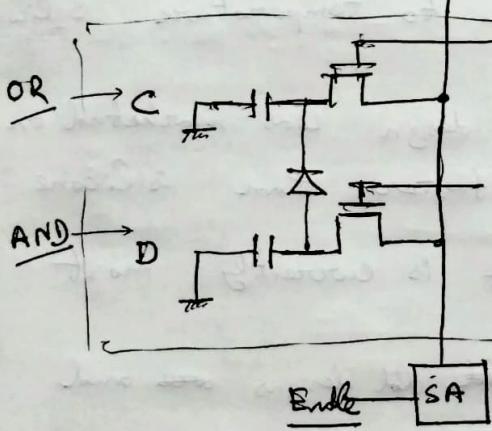
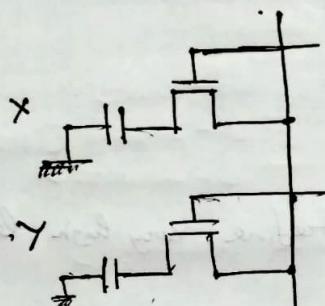


$$V_{DD/2} - 8$$

if this is charged  
nothing happens.  
only when discharged  
~~charge~~ moves to the  
showed direction



If this is discharged  
nothing happens.  
only when discharged, charge  
moves to the showed direction.



$$X + Y$$

- 1) Copy 'X' to 'C'
- 2) Copy 'Y' to 'D'

if D @ has charge, then  
charge can flow to  
 $D \rightarrow C$  but not  
vice versa, and  $\leftarrow$   
by this we can make  
sure if only 'D' is  
charged, charge is  
moved to 'C'.  
by this AND/OR can  
be realized.

Q. What are resistivity of different materials?

$$\rightarrow \rho_{Cu} = 10^{-6} \Omega \cdot \text{cm} \quad \text{Copper [conductor]}$$

$$\rightarrow \rho_{Ge} = 50 \Omega \cdot \text{cm} \quad \text{Germanium}$$

$$\rho_{Si} = 10^3 \Omega \cdot \text{cm} \quad \text{Silicon}$$

$$\rightarrow \rho_{Mica} = 10^{12} \Omega \cdot \text{cm} \quad \text{[Insulator]}$$

Compounds  $\rightarrow$



Properties:

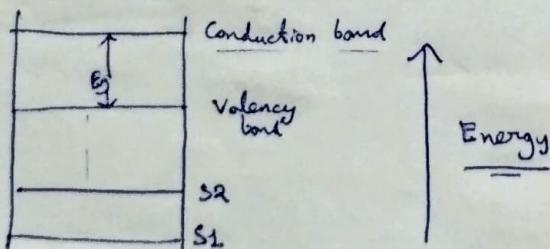
- i) Germanium is easily extractable, easy to refine, very high level of purity.
- ii) Fabrication process.
- iii) Silicon is more sensitive to temperature. Silicon is more abundant.
- iv) Gallium arsenide based design was marketed in 1970 and it was 5 times faster than silicon.
- v) Currently gallium arsenide is currently most used.
- vi) Silicon is favoured by old designs and research is being done.

Silicon has  $4e^-$  in valence band ( $14e^-$ )

Germanium has  $4e^-$  " " " ( $32e^-$ )

Gallium has  $3e^-$  in valence band ( $31e^-$ )

Arsenic has  $5e^-$  in valence band ( $32e^-$ )



$E_g \rightarrow > 5\text{eV}$  (Insulator)

$E_g \rightarrow > 0.67\text{eV}$  (semi)

$E_g \rightarrow \sim 0$  (Conductor)

\* they overlap.

► Let's assume

Si,  $1e^-$  in conduction band ~~are~~ among  $10^{12}$  atom.

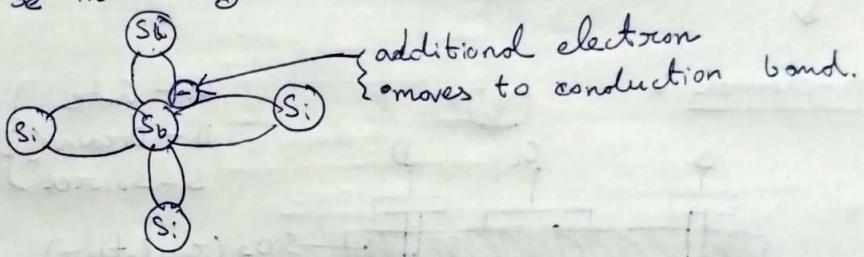
now, 1 impurity in  $10^7$  atoms.



$10^5$  times increase in conduction.

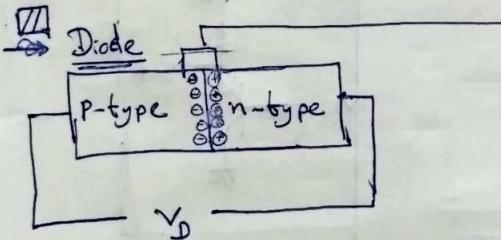
### ► n-type \ P-type :

$Sb \xrightarrow{\text{Antimony}}$  ~~Si~~  $\xrightarrow{\text{Sb}}$   $Se^-$  in valency bond.



n-type  $\rightarrow$  majority carriers  $e^-$

p-type  $\rightarrow$  " holes  $(h^+)$   
[I'll be working like this]



### Depletion region

→ can be biased using  $V_D$ .

no-bias  $\longrightarrow (N_0 = 0)$

forward-bias

reverse-bias

→ In reverse bias  $\rightarrow$  The depletion region is extended.

### VII Field effect Transistor :

- JFET
- MOSFET  $\longrightarrow$  metal on ~~on~~ metal oxide Field
- MESFET

### Comparison b/w FET & BJT;

i) FET is unipolar.

ii) In BJT, ~~is~~ current control.

FET, Voltage control.

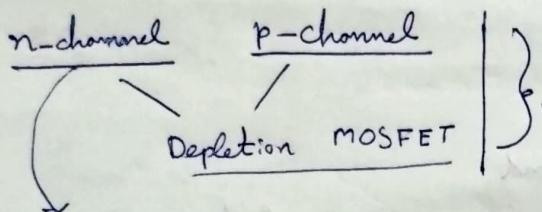
↓ Output current,  $I_D$  is a function of gate voltage.

iii) FET are more temporally stable.

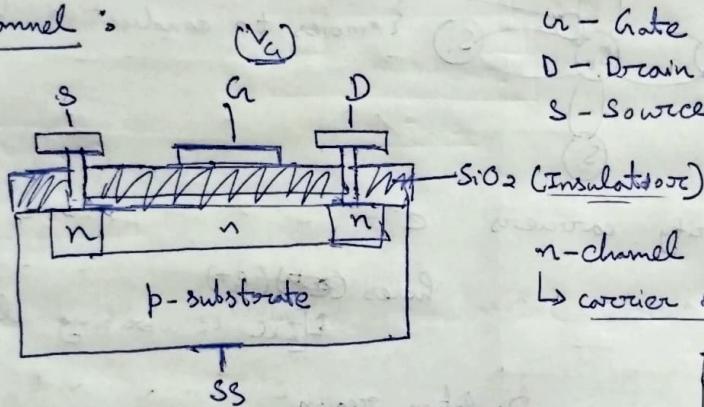
iv) FET are smaller.

### ► MOSFET:

- Depletion MOSFET  
→ Enhancement MOSFET }



### ► n-channel:



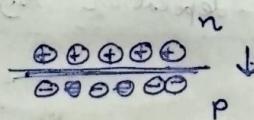
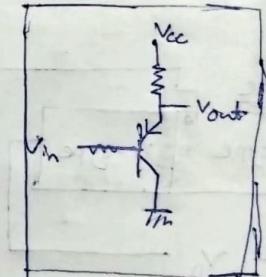
G - Gate  
D - Drain (Destination)  
S - Source

n-channel  
↳ carrier is electron

$V_{GS}$  → Gate to Source voltage  
↳ controls the current flow.

When,  $V_G = 0$  (high impedance)

The, n-p junction will have  $\text{∅}$



$V_G$  → Input,  $I_{DS}$  → Output (Drain to source current)

⇒  $V_G = 0$

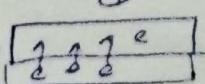
$V_{DS} = +V_{DD}$  →  $I_D$  current

⇒  $V_G = \text{'-ve'}$

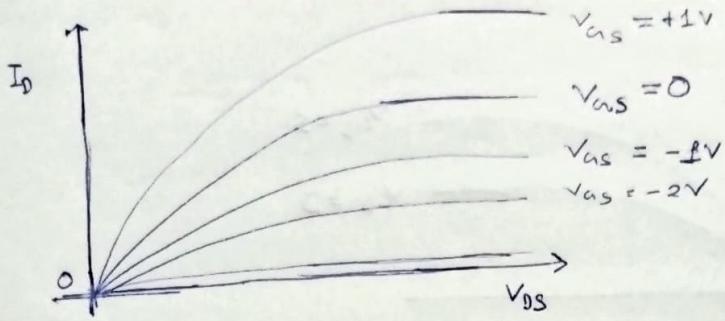


→ no. of free electrons is reduced.  
hence,  $I_D$  is less.

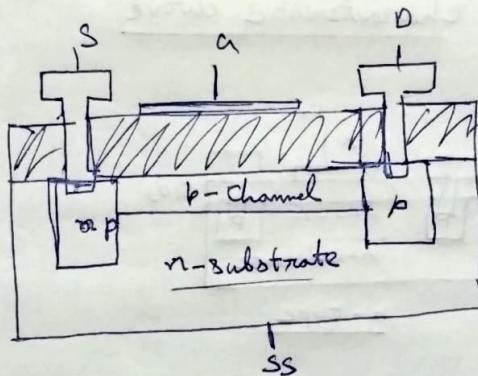
⇒  $V_G = \text{'+ve'}$



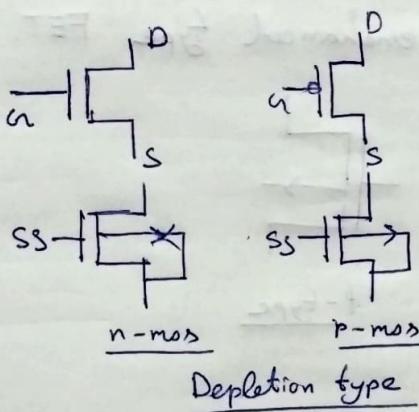
→ no. of free electrons is increased.  
hence,  $I_D$  is more.



Characteristic curve



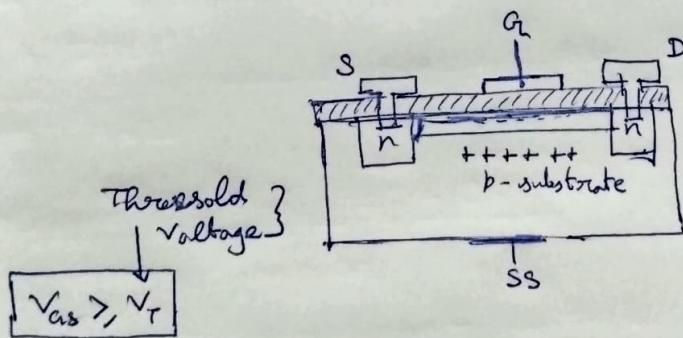
The logic is same as, n-channel.



07/02/24

#

Enhanced Type :

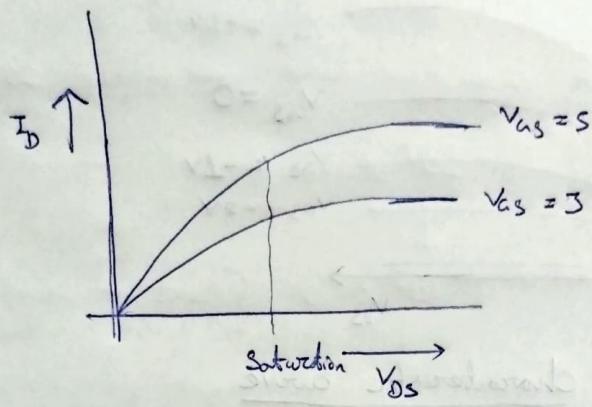


p-type

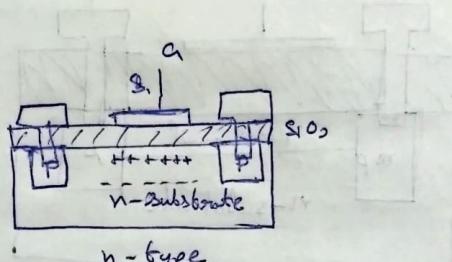
$$\frac{V_{GS}}{V_{DS}} > 0$$

$$\frac{V_{GS}}{V_{DS}} > 0$$

If we keep increasing  $V_{DS}$ , at some point current will be saturated as no more minority carrier will be available.

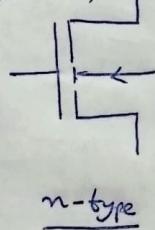


characteristic curve

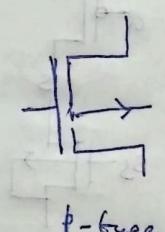


$$\begin{aligned} V_G < 0 \\ V_{DS} < 0 \end{aligned}$$

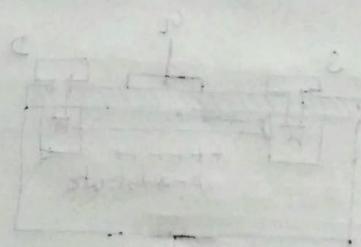
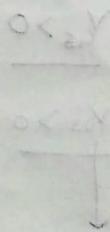
- ④ There is no channel in enhanced type FET.  
hence, the symbols,



n-type



p-type



↓ Blue light  
Emissivity