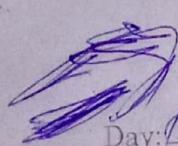


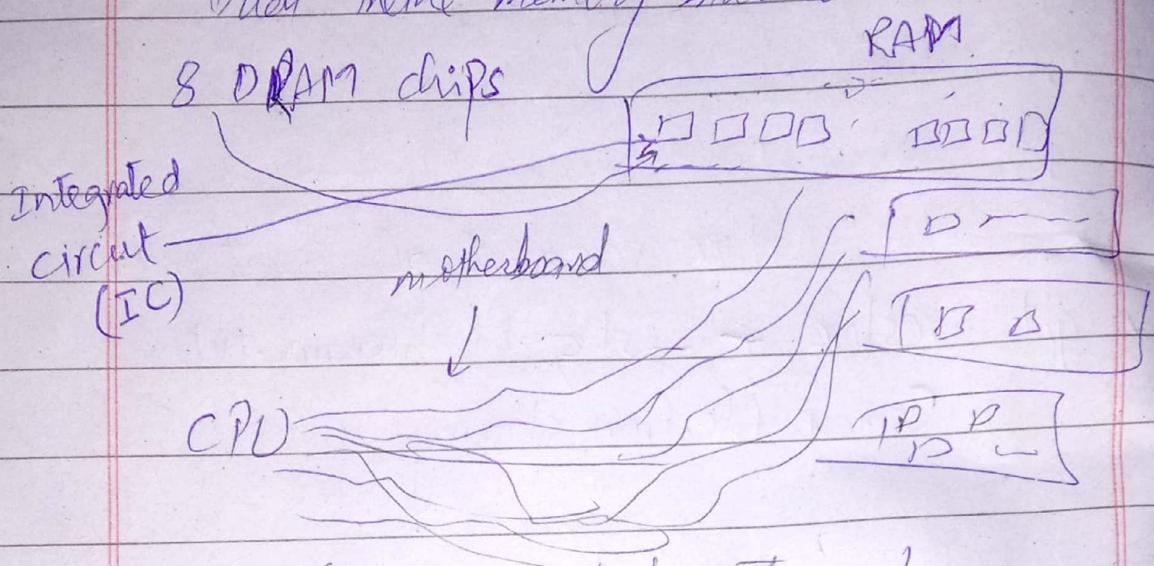
How Computer memory works



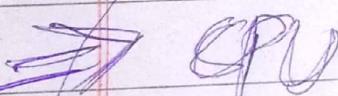
Day: How does computer memory work?

- DRAM is also called (DIMM)

Dual inline memory module



~~if 4 DRAM then two has different channel and other two diff.~~



- In CPU memory controller manage the communication with DRAM
- There is a one other section that control the communication with M.2/SATA HDD.

⇒ Memory channel DDRS

memory is divided into two channels

A S B

- A & B transfer 32 bits at a time using 32 data wires.
- 21 additional wires to memory that carries address where to read & write data
- 7 control signal wires, Command are relayed. These 7 signal are used to send and receive all the bits

Power (electricity) supplied by the motherboard and handled by the chips that are placed in center.

- ⇒ RAM microchips (8 bit)
 - there are four layers between two at center interconnected the lower (ball grid array) to upper (Die).

- ⇒ Die
 - It has 8 Bank groups
 - one Bank group consists of 4 Bank.
 - so total Bank $8 \times 4 = 32$ Banks

~~Each~~ ⇒ Bank
Each Bank has a massive

array (65536 ^{rows} memory cells) tall / length
 (8192 ^{columns} cells) across / horizontally

There is a circuitry outside the Banks which consists of thousand of wires for communication.

⇒ How to access data in banks
 31 Bit address

⇒ first 3 bits to identify Bank group
~~4~~ $2^3 = 8$
 1 1 1 8 combination as group also

⇒ next two Bits to select Bank
 $2^2 = 4$ as per group
 There are ~~two~~ 4 Banks.

⇒ Next 16 Bits to identify the row of data for accessing like
 $2^{16} = 65536$

⇒ 8192 cells group of 8 in one Bank - why in 8 columns Because it can only read and write 8 bits

at a time.

first 5 bits Bank address

Next 16 bits Row address

Next 10 bits Column address

⇒ Cell / I^TC DRAM memory cell.
One memory cell

Two Parts

→ Capacitor: To store 1 bit of data either 1 or 0 in the form of electrical charges or electrons. Its shape is like deep trench access.

→ Transistor: It allows for reading and writing the data. It controls the flow of electrons. Like amplify which represent (1) or to zero.

If capacitor charged up with electron it means 1. and if no charge present and it is zero ~~large~~ volt then

it is a binary zero.

Types of capacitors / evaluation.

- ① MOS capacitor ② stacked
- ③ trench ④ substrate plate trench capa-
- ⑤ 3-D stacked capacitor

~~wordline~~

~~Bitline~~

~~There is a gate which bias on which if wordline act then the gate will open and through channel the capacitor charged up.~~

⇒ How to read stored value in capacitor.

~~we can read the value in the capacitor by measuring the volt amount of charge in capacitor.~~

⇒ When word horizontal

Wordline

Bitline ↑

↓ vertical

Applying voltage on the wordline turn on the transistors which leads to the full charge of capacitor.

Now electron can flow \downarrow and connect to bitline.

- when the wordline is off transistor also turned off ^{as well} and capacitor isolated from bitline saving the data/charge that were previously written.

- Electron Leakage \rightarrow Refresh memory

The transistor is incredibly small so the electrons leakage happened slowly, ^{through channel} so we need to capacitor needs to be refreshed to recharge ~~over~~ the leaked electrons.

To resolve this issue we use refreshing memory cells technique.

~~cover later.~~ \rightarrow Solution?

Each column is connected to a sense amplifier. When a 32 bit address come then first 8 bits select the specific Bank and next 16 bit select the row.

in the bank, all the wordline
except the selected one turned
off. now ~~other~~ in that one
active row which contain
892 columns is active. Some
capacitors are charged (\oplus) some
not (\ominus). So in the charged
capacitors the electron started
to leak in the bitline. however
the bitline is connected to the
sense amplifier so it sense
the leakage of electrons so
it amplify the electrons and
make it full. Also in the
capacitor that store (\ominus) zero
electron the electron started
to charge in the capacitor
through bitline channel. So
sense amplifier then reduce
the electron from that bitline
to zero.

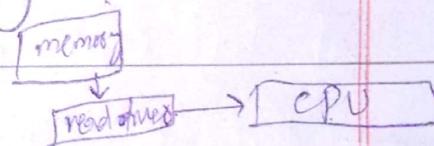
Now the row is selected from 65536 lines how but what about columns we only want to access 8 bit but one row contain 8192 bits so how?

~~Read Memory~~

The last 10 bits from the 31 bit address are used to identify which columns are need to be accessed.

These 10 bits select the 8 bit from columns and then send to

Read driver through wires and then send to CPU



~~Write Memory~~ (31 bit address)

First Five bits select the Bank

Next 16 Bit select the row.

Next 10 bit select the specific columns and connect to the

Write Driver. At this time

the bits that should be

written in memory was already sent to the write driver by ^{CPU} which is in the DRAM. ~~After~~ the write drivers are much more stronger than sense amplifiers. So it charges ^{overvoltage} the capacitor if (a) needed otherwise if not charge if (c) needed. ~~when~~

At the time of row selection all the capacitor charged to 0.5 Volts which is the half of full signal voltage. It is called Precharge.

\Rightarrow Read & write happen concurrently to the all four IC's with shared address & command wires but separate data wires

voltages for 1 in different DDR

DDR 2.5 V

DDR 4 1.2 V

DDR 2 1.8 V

DDR 5 1.1 V

DDR 3 1.5 V

precharge is half for each.

DDR5 precharged voltages 0.55V
it store 1.4V because of
leaked electron/voltages.

2^3 8 Bank group

2^{13} 8192 bitline

2^5 32 Bank group

2^{16} 65536 wordline

⇒ Refreshing the rows because
of volt leakage.

Sense amplifiers open the first row
and then refresh the whole row like
1 to full charge 0 to zero charge.
This happen row by row.

⇒ Refreshing one row takes 15 nano sec time.
⇒ 3 milliseconds to refresh 65536 rows
1 whole Bank.

Refresh action occurs after every
64.00 milliseconds for each Bank.

Each IC ^{on} ~~off~~ Memory handles

4800 million Requests/second

In one second the whole IC⁰ Banks refreshes 16 times.

⇒ Why need this much faster memory
e.g.: for video games
each moment and changing environment shadow and many more features need fast and on-time calculation in the memory.

⇒ Row HIT / Page Hit

If there is data that is already loaded on the memory and we are also reading and writing that data (same address). It means the row is open and then we can directly read or write the data skipping some steps like. For opening a row we need to

Times ~~the~~ - Perform
actions in clock
↓ cycle

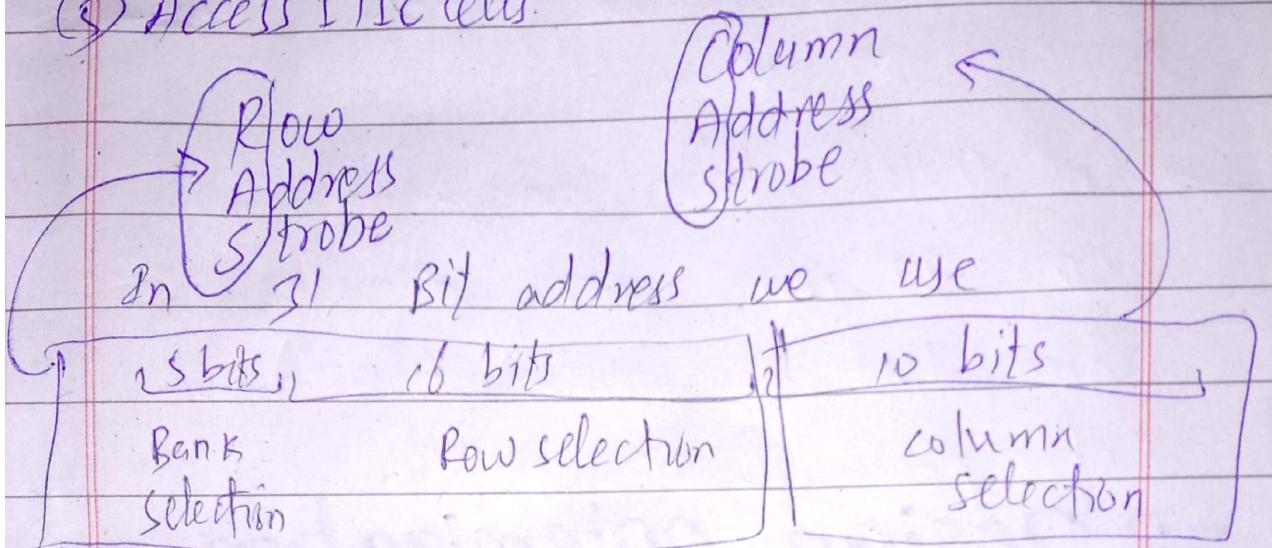
(2) Row close

(2) Precharge Bitlines → 39 Precharge + RP

(3) Row open → 39 RAS to CAS Delay + RCD

(4) Column address → 40 CAS latency.

(5) Access ITIC cells.



If a row is already selected we can skip this step and directly use 10 bits to select columns.

It reduce the amount of time.

~~If~~ If the same row is requested that was already open for read/~~or~~ write we called page hit/row hit.

~~If~~ If another row requested for read and write except the open one. then it is called Row miss.

⇒ Row Thashing
when a program ^{request} jumps from one row to another over and over. It becomes inefficient in terms of both energy and time.

⇒ In DDR5 there are 32 banks so 32 rows open at the same time to increase the likelihood of Row hit in each IC.

⇒ Design Optimization

⇒ Burst Buffer
Read and write driver each has 8 wires that connects to memory.

In Burst Buffer now we place Burst Buffer between Read ~~and~~ write driver and memory. Now memory connected to each Burst Buffer is with 128 wires.

For column selection 10 bits divided into, 6 bits, for column multiplexer and 4 bits for Burst Buffer.

⇒ Burst Buffer for reading.

It consists of 128 bits

16 columns

$$2^4 = 16$$

8 Rows

Burst Buffer one column for which has 8 bits connected to Read driver through 4 bits

★ Same for writing.

3 Benefits

16 sets of 8 bits per microchip $= 128$ bits

Total 8 microchip SO $= 128 \times 8 = 1024$ bits

⇒ Another optimization

each IC has 65536 rows and 892 columns

So the ~~so~~ the bitlines are very long

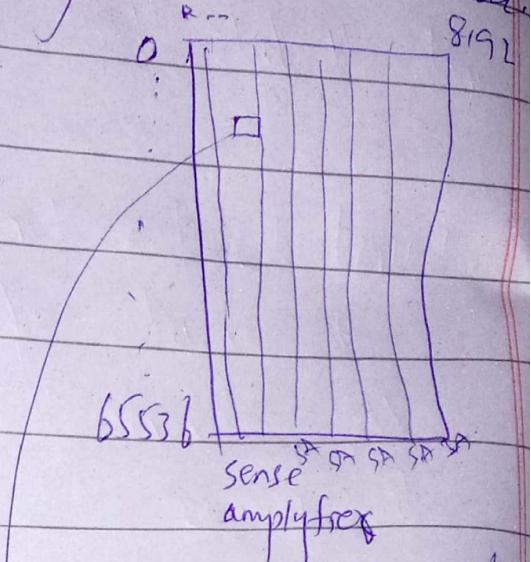
SA = sense amplifier

Day:

Date:

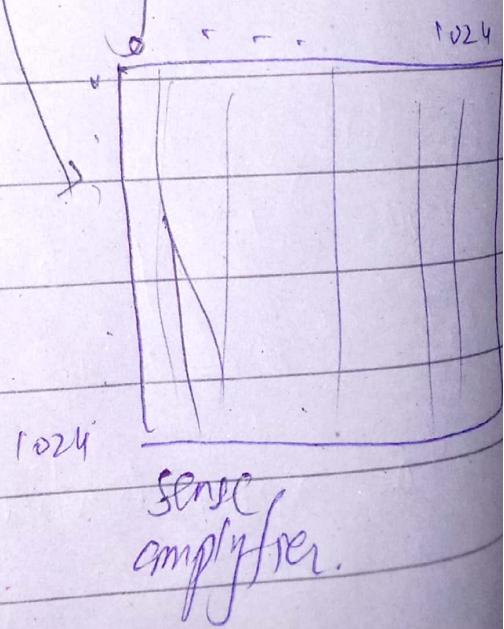
and also charging the capacitor
for sense amplifiers is hard.

It is a massive.
structure.



So, we divide and grouped the
block of 1024×1024 bits and each
block has little ^{sub}sense amplifiers

So, for charging the capacitor becomes
fast because of having SA near
to capacitors.



Another

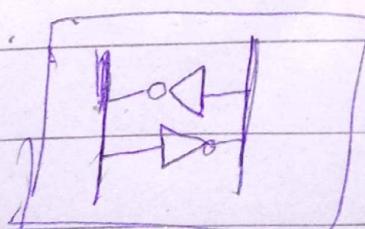
-Sense Amplifiers connected to each bit line or column.

Now we connected two bitlines to single column but alternating the rows of memory cells

Now if one row is open we can then half of the bitline active. Other half are passive(off) also vice versa.

→ How above works?

Inside the sense amplifier there is a cross coupled inverter. It turns on one bit line and off the other bitline.



Sense Amplifier.

→ Benefits

- ①. Precharge
- ②. Noise Immunity
- ③. Parasitic Capacitance.