

Subject: Computer architecture. Lec 1 - Memory!

Main Memory:

- RAM \Rightarrow Random access memory
- ROM \Rightarrow Read only memory.

RAM:

- DRAM \Rightarrow Dynamic RAM.
- SRAM \Rightarrow Static RAM.

DRAM consists of capacitors which slowly leak their charge. Thus, they must be refreshed periodically to prevent data loss.

- Simple in design and cheap.

SRAM consists of circuits similar to D flipflops

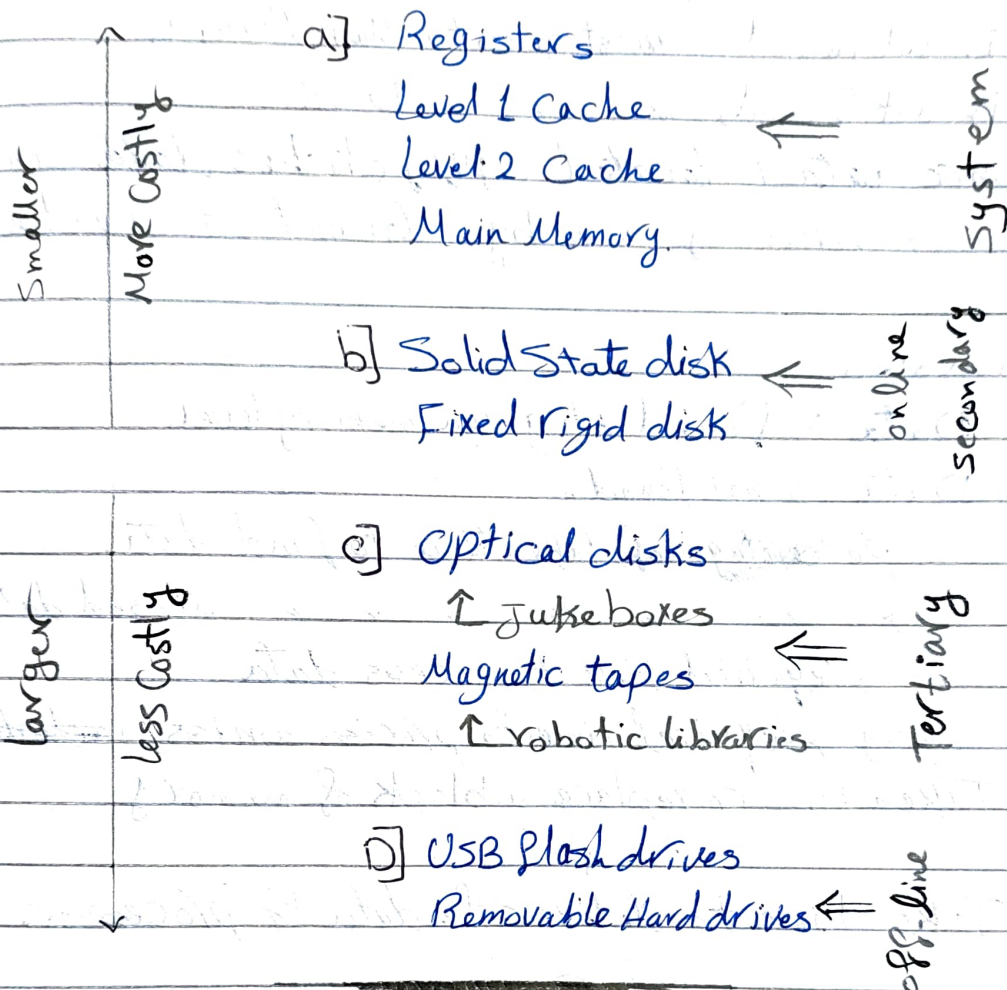
- Doesn't need to be refreshed
- Very fast memory.
- Used to build ~~the cache~~ cache memory.

ROM

- Doesn't need to be refreshed. "needs a little charge".
- Used to store permanent or semi-permanent data.
- Data ~~persists~~ persists even when system is turned off.
- Cheaper than RAM and slower than RAM.

- * Generally faster memory is more expensive than slower memory.
- * Hierarchical fashion used to provide best performance at the lowest cost
- * Small, fast storage elements are kept in the CPU
- * Larger, permanent storage are kept further from CPU.

Storage organization:-



- **Registers:** Storage locations available on the processor itself.

- **Virtual Memory:** implemented using a hard disk.

- Used to extend address space from RAM to hard drive.
- Provides more space.

Cache Memory • provides speed.

- Data Fetching.

1. CPU send request to cache. "Nearest memory"
2. if not found, request goes to Main Memory.
3. if not found, " " " " Disk.
4. Once data is located, CPU fetches data and number of nearby elements to Cache.

- Definitions:-

- hit \rightarrow Data is found at given memory level.
- miss \rightarrow Data is not found.
- hit rate \rightarrow percentage of taken time to find data.
- miss rate \rightarrow " " " " " " missed data.
- hit time \rightarrow required time to access data
- miss penalty \rightarrow required time to process a miss, ~~includes~~
 - \uparrow = Taken time to replace a block of memory

+
Taken time to deliver the data to processor.

$$\text{Miss rate} = 1 - \text{hit rate}$$

- Locality:-

principle of locality:- Once a byte is accessed, The nearby data will be needed soon.

so, while hitting, ~~a block~~ an entire block is copied.

Locality Forms:-

- Temporal locality → Recently accessed data tends to be accessed again.
- Spatial locality → Accessed data tends to ^{be} cluster.
↑ جميع البيانات التي تم الوصول إليها
- Sequential locality → Accessed instructions tends to be accessed sequentially.

Cache Memory:-

- purpose → Speed up data accessing by storing recently used data closer to CPU, instead of Main memory.
- Content addressable Memory → accessed by Content.
- Smaller than Main Memory but faster than it.
- Single large Cache Memory not desirable because it takes longer time to search and that is because it searched by Content not address.

* Main Memory's data accessed by address while ← Fast
Cache Memory's data accessed by Content ← Slow

Cache Mapping Schemes:-

1. Direct Mapped Cache.
2. Fully associative mapped Cache.
3. N-Way set associative Mapped Cache.

1. Direct cache mapping:-

• The simplest cache mapping scheme.

• $Y = X \text{ modulo } N$

$Y \rightarrow$ Number of block in Cache memory.

$X \rightarrow$ Number of block in Main memory.

$N \rightarrow$ ~~How~~ Total blocks number in Cache.

• Main memory address:-

Tag	Block	offset
-----	-------	--------

- ~~Tag~~ offset.

- offset: identifies address within block.

يعرف مكان الـ data جوه الـ block

- Block: identifies which block contains data.

يعرف الـ block الى جوه الـ data

- Tag: Marks data in block.

- Mapping Mechanism:

1. Determine address format for mapping:-

• number of bits for offset $\Rightarrow 2^n = x$

$x \rightarrow$ Bytes per block. "Size".

• n-bits for Block $\Rightarrow 2^n = x$

$x \rightarrow$ Number of Blocks in Cache

• n-bits for Tag \Rightarrow remaining bits.

• Use modulo algorithm to get mapped block in Cache.

• n-bits for address \Rightarrow Memory Size = 2^{main}

ex:- Byte addressable main memory has 4 blocks, and a cache with 2 blocks where each block is 4-byte.

⇒ N -blocks in main memory = 4

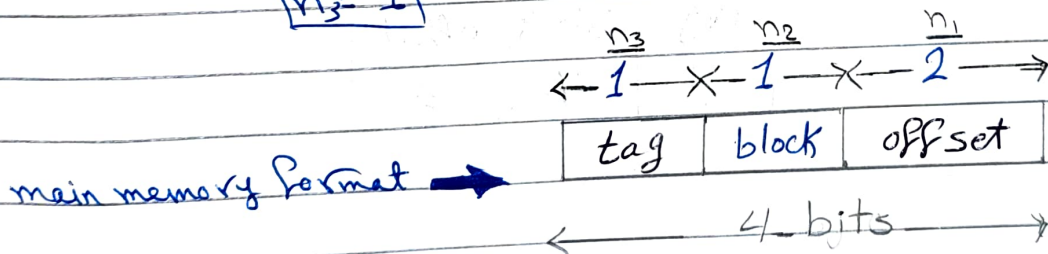
N -blocks in Cache = 2

Address = ~~4 byte~~ 16-bit = 2^4 ⇒ address bits = 4

(a) n -bits for offset ⇒ $2^n = 4$ ← each block is 4-byte "size"
 $\boxed{n_1 = 2}$

(b) n -bits for block ⇒ $2^n = 2$ ← 2 blocks in Cache
 $\boxed{n_2 = 1}$

(c) n -bits for tags ⇒ Size - ($\boxed{n_1}$ bits for offset) - ($\boxed{n_2}$ bits for block)
 $n = 4 - (2) - (1)$
 $\boxed{n_3 = 1}$



ex 2: Assume a byte addressable memory consist of 2^{14} byte,
Cache has 16 block, each block has 8 byte.

→ word = 1 byte
memory size = 2^{14} byte
∴ Address bits = 14 bit.

n-bits for offset \Rightarrow block size = 8 = 2^3 byte
~~n-bits~~ n = 3 bits.

n-bits for block \Rightarrow Number of blocks = 16 = 2^4
n = 4 bits.

n-bits for tag \Rightarrow remainder bits
n = $14 - (3 + 4)$
n = 7 bits.

tag	block	offset
7	4	3

ex: 16-bit memory address and 64 block of Cache where each block contains 8 byte.

⇒ Address bits = 16-bit.

offset ⇒ block size = 8 = 2^3 ⇒ offset bits = 3 bits

Block ⇒ Number of blocks = 64 = 2^6 ⇒ block = 6 bits

Tag ⇒ Address bits [offset bits - block bits]
Tag = 7 bits.

- Direct cache mapping depends on:

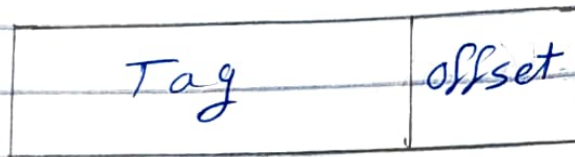
- Number of bits in the main memory address. ⇒ address bits
- Number of blocks in cache memory. ⇒ size of block field
- Number of addresses in block. ⇒ size of offset field.

- Direct cache mapping issue:

Memory block is replacing cache block depends on memory address without checking for free space, which leads to Data loss in cache.

2. Fully associative cache mapping:-

• Main memory address:-



- When Searching in Cache, all tags are being searched in parallel to retrieve data quickly.
- Requires special and costly hardware.
- Mechanism:
 - Searching for ~~empty~~ free space in Cache, then memory block mapped to this place.
 - if there is no free space, There is an algorithm used to find victim block which is evicted.

3. Set associative Cache: "N-Way set associative".

- It Combines the ideas of direct mapping and Fully associative cache.
- Memory reference is being mapped to free space in subset of Cache slot.
- Number of cache blocks in a set varies according to System design, ex:
 - 2 way set associative cache
 ↗ Set = 2 blocks
- Main memory address:

Tag	Set	offset
-----	-----	--------

- offset → Chooses the byte within cache block.
~~to be mapped.~~

- Set → Determines which set has the block to be mapped.

ex: Suppose a byte addressable main memory contain 1 MB.
Cache consist of 32 line, each block contains 16 byte, using:

A) Direct: Memory size = 2^{20}
Address bits = 20 bit.

offset \rightarrow block size = $16 = 2^4$
offset bits = 4 bit.

Block \rightarrow Number of block = $32 = 2^5$
Block bits = 5 bit.

Tag $\rightarrow 20 - [4 + 5] =$
Tag = 11 bit.

Tag	Block	offset
11	5	4

B) Fully associative:
Address bits = 20 bit.

offset bits = 4 bits

Tag bits = 16 bits

Tag	offset
16	4

c) (4) way set associative:

Address bits = 20 bit

~~n. set = 4 set~~

offset bits = 4 bit.

number of sets \rightarrow Number of blocks in a set = 4

$$n. set = \frac{\text{Number of blocks in cache}}{\text{Number of blocks in a set}} = \frac{32}{4} = 8$$

$$n. set = 8 = 2^3$$

set bits = 3 bit

$$\text{Tag} = 20 - [3 + 4] = 13 \text{ bit.}$$

Tag	set	offset
13	3	4