

★ Cache Memory

→ Locality of Reference:

During execution of a given program memory references tend to cluster.

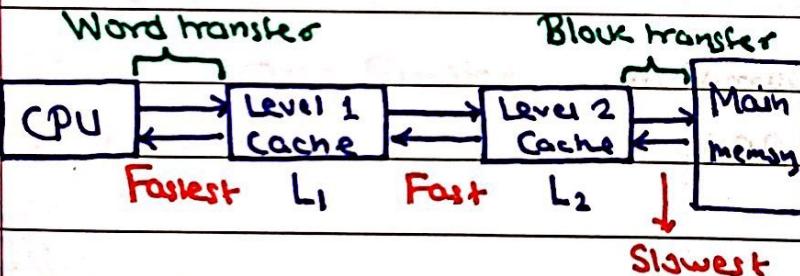
e.g.: sum of elements of an array. (# contiguous memory locations)

→ Cache

① Small amount of fast memory

② Between CPU & Main memory

③ Maybe located on CPU chip



★ Memory Hierarchy

→ Registers

CACHE (L₁)
L₂)

Main Memory

Disk Cache

Disk

Optical

Tape

↓ cost per bit

↑ capacity

Inbound ↑ access time

↓ freq of access time by the processor

Outboard

Offline

* Cache Operation:

- ① CPU req contents of memory location.
- ② Check Cache for this data
(IF present get from Cache) [Cache Hit]
- ③ IF not present [Cache Miss]
Read required Block from main memory to cache
then deliver from Cache to CPU

* Cache includes tags of blocks of main memory in each cache slot.

* Cache Read Operation:

→ Read address (RA)

Start

Receive RA from CPU

Is block contains RA in

No

Access MM for block w/ RA

Yes

Fetch RA word & deliver to CPU

Allocate Cache line for MM block

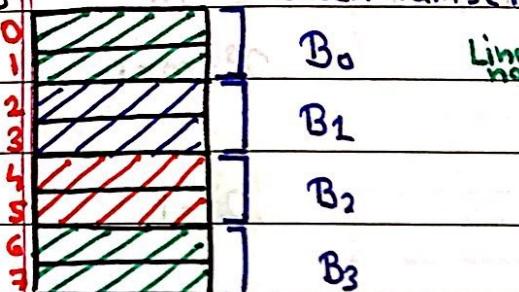
Load MM block to cache

Deliver RA word to CPU

Done

★ What is cache line?

Byte Address



Block number

Line no. Tag

Block

1

2

3

★ Here we have considered cache to have 3 lines (Hypothetical)

Block length

(2 words (bytes))

1



[] B0 (equivalent to line 1 above)

★ Cache Design

- ① Cache Size
- ② Cache Mapping Function
- ③ Replacement Algo
- ④ Write Policy
- ⑤ Line Size
- ⑥ No. of caches (Levels)

① Direct Mapping

- ① Each Block of main memory maps to only one cache line
- ② Address is in 2 parts

Least sig w bits identify unique word (block size = 2^W)

Most sig s bits specify one memory block

Tag : s-r

r : Cache line

s-r	r	w
Tag	Line	Word
8	12	4

24 bit address

16 byte block (2^4)

No two blocks in same line have same tag field

* Cache line

MM blocks assigned

0, m, 2m, 3m, ..., $2^s - m$

1, m+1, ..., $2^s - m + 1$

$m-1, 2m-1, \dots, 2^s - 1$

* MM Block j i.e. B_j assigned to Cache line i

then $i \equiv j \pmod{m}$

m : no. of cache lines

① Address length = $(s+w)$ bits (Addressable units 2^{s+w} words or bytes)

② Size of tag = $(s-r)$ bits

③ Block size (MM) = line size (Cache) = 2^w words

④ No. of lines in Cache = $m = 2^r$

⑤ No. of blocks in MM = 2^s

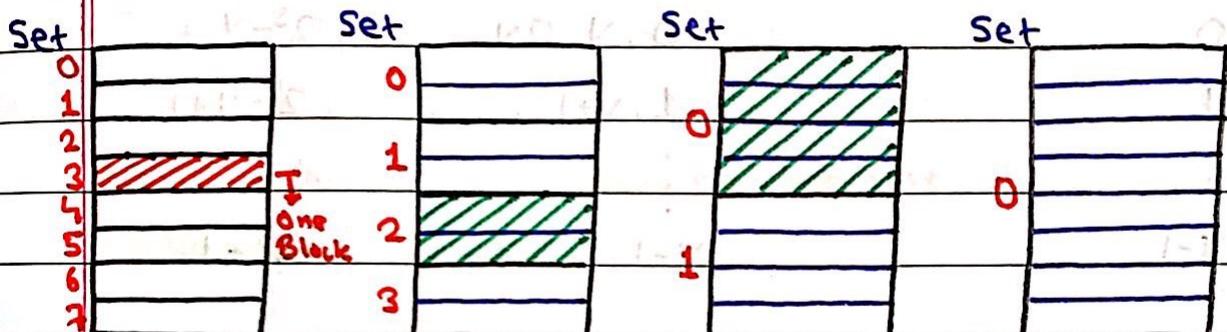
* Simple

(* Inexpensive)

* Fixed location for a given block (MM)

If program accesses 2 memory blocks that map to same line repeatedly, cache misses are very high
(i.e. memory blocks with index same value modulo m)

② Associative (Set Associative)



1-way

8 sets

1 block each

2-way

4 sets

2 blocks each

4-way

2 sets

4 blocks each

8-way

1 set

8 blocks

(Direct mapped)

(Fully associative)

a) Fully Associative:

① MM block can load into any line of cache

② Memory address as tag & word

③ Tag uniquely identifies the block of memory

(Every line tag is examined for match)

④ Cache searching gets expensive

b) Set Associative

(K lines/block per set \Rightarrow k -way set associative mapping)

① Cache divided into v sets

② Each set has k lines

③ No of lines in cache $m = v \times k$

④ B_j goes to set i

when $i \equiv j \pmod{v}$

★ Set no.

MM blocks held

0

 $0, V, 2V, \dots, 2^S - V$

1

 $1, V+1, \dots, 2^S - V + 1$

⋮

⋮

V-1

 $V-1, \dots, 2^S - 1$

★ tag | set | word

s-d

d

w

0000

0000

0000

① Address length = $(S+d+w)$ bits② Addressable units = 2^{S+w} words③ Block size = line size = 2^w words / bytes④ No. of blocks in MM = 2^S

⑤ No. of lines / set = K

⑥ No. of sets = V = 2^d ⑦ No. of lines in cache = KV = $K \times 2^d$ ⑧ Size of tag = $(S-d)$ bits

* Replacement Algo:

① Direct mapping:

① No choice (with address width only enough)

② Each block only maps to one line

③ Replace that line

→ Write Policy:

- ① Mustn't overwrite a cache block unless MM is up-to-date
- ② Multiple CPU's may have individual caches
- ③ I/O may address MM directly.

① Write through

→ All writes go to MM as well as cache

→ Lot of traffic (but no storage limitation)

→ Slows down write

→ Remember bugs (writes through cache)

(Multiple CPUs can monitor MM traffic to keep local cache up-to-date)

② Write back

→ Updates initially made in cache only

→ Update bit for cache is set when update occurs

→ If block to be replaced: write to MM only if update bit set,

(I/O must access main memory through cache)

② Associative mapping

- Least Recently Used (LRU) (which one block is LRU)
- FIFO (replace the block which is in cache for longest)
- Least frequently used (replace block with fewest hit)
- Random

* LRU

① Consider 4 line set in cache.

→ Control bits:

Tag bit

* 2-bit counter for each line (to track LRU block)

d-bit: Dirty Bit

f-bit: For occupied bit to find out if address is valid

② Initially reset all the counters (d-bit & f-bit)

③ Cache hit occurs:

→ Set counter value to 0 at this cache line.

→ for other counters if value less than the reference line increment the counter value.

else do not change counter value.

Cache miss:

① Set is not full

Set counter value of cache line to 0

Set f-bit to 1

increment counter of other lines whose f-bit is 1

② set is ~~not~~ full exactly

→ The line with ↑ counter value is removed, write back it

d-bit is 1 means new and old count is same (A)

→ new block is transferred to this line (B)

→ reset d-bit & counter to zero and repeat (C)

→ other counter values ↑ by 1

★ External Memory:

→ Performance:

* Recovery:
memory to
recovery before
next access

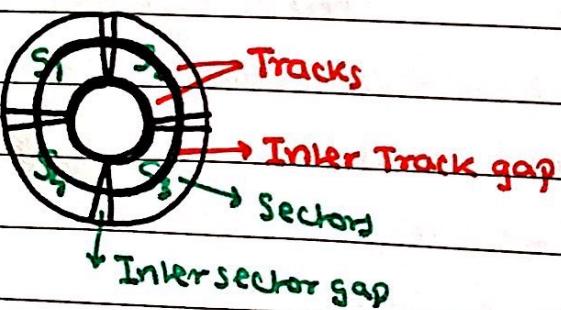
- ① Access Time: (Time req between presenting address & getting valid)
- ② Memory Cycle Time: Access + Recovery
- ③ Transfer Rate: (Rate at which data can be moved)

→ Physical Types

- ① Semiconductor : RAM
- ② Magnetic: Disk & Tape
- ③ Optical: CD & DVD

★ Disk Data Layout

a) Constant
Angular
Velocity

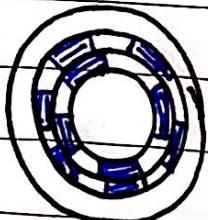


* Each track has same width as head.

* Each track same no. of bits

* ∵ bit density ↑ from outer ring to inner (track)
(Hence ↓ wastage of space)

b) Multiple zone recording



* Each zone fixed bits per track (Complex circuitry)

* Head

① Fixed (one head per track)

② Movable (one head for complete disk side)

* Removable / fixed

→ Single / double sided

→ Single / multiple platter

→ Head Mech

- Contact (Floppy)

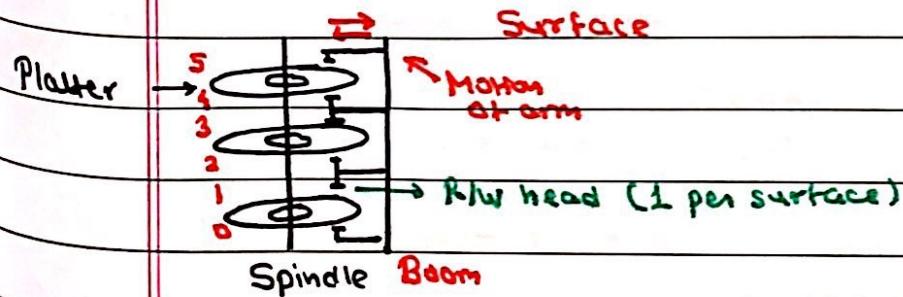
- Fixed gap

* Multiple Platter

→ One head per side

→ Heads joined & aligned

→ Aligned tracks on each platter form a cylinder



* Disk addressing:

- ① Sector no.
- ② Track no.
- ③ Surface no.

* Speed

① Seek time (T_s)

→ Moving head to correct track

② (Rotational) Latency ($\frac{1}{2}r$)

→ Moving to correct sector

$$\text{Access Time} = \text{Seek time} + \text{Latency}$$

③ Transfer Rate (T)

* Transfer Time T (Transfer Rate)

$$T = b/rN$$

b: No. of bytes to be transferred

N: No. of bytes on the track

r: rotational speed (revolution per second)

$$\therefore \text{Total time} = T_s + \frac{1}{2}r + T$$

$$= T_s + \frac{1}{2}r + \frac{b}{rN}$$

* Performance Issues:

→ Execution Time

- ① Elapsed Time = CPU time + wait time (I/O, other programs)
- ② CPU time = user CPU time + system CPU time

* Our focus: user CPU time (time spent in executing the lines of code that are in our program)

* For a machine X:

$$\text{① Performance}_X = \frac{1}{\text{Execution time}_X}$$

$$\star \text{Seconds} = \frac{\text{Cycles}}{\text{program}} \times \frac{\text{seconds}}{\text{cycles}}$$

$$\star \text{Clock freq} = \frac{\text{cycles}}{\text{seconds}} \quad (\text{clock rate})$$

$$\star \text{Cycle time} = \frac{1}{\text{clock freq}}$$

* Performance equation:

$$\text{seconds} = \frac{\text{cycles}}{\text{program}} \times \frac{\text{seconds}}{\text{cycle}}$$

$$\therefore \text{CPU execution time} = \text{CPU clock cycles} \times \text{clock cycle time}$$

for program

$$\star \text{CPI} = \text{cycles/instruction}$$

$$\star \text{MIPS} = \text{millions of instructions/second}$$

* Performance Equation 2:

CPU execution time = Instruction count \times Avg. CPI \times Clock cycle time

Time spent in L1 cache + wait time + wait time = total time

→ Cache

① Hit rate: Fraction of memory access that are hits *

② Miss rate = 1 - Hit rate

③ Hit time: time to determine access is a hit +

time to access & deliver the data to CPU

④ Miss penalty: time to determine access is a miss +

time to replace the block in upper level of cache

**EXTRA
time**

time to deliver data to CPU.

* Avg. memory access time = Hit time + (miss rate \times miss penalty)

Memory stall cycles = Memory accesses \times miss rate \times miss penalty

CPU time = (Execution cycles + memory stall cycles) \times cycle time