

L 25 / 18.10.2023 /Cache Memory

⌚ What, Why, How to use & Benefits?

⇒ Semiconductor memory.

⇒ Must faster than RAM (Main memory)

⇒ Expensive, cost per bit is very high.

⇒ Because of that, capacity of cache is kept very small than RAM.

⌚ Why ?

⇒ In a processing of an instruction, maximum time spend on Fetch and reading Data.

⇒ Time needs for fetch \Rightarrow depends on RAM and BUS.

Seacal-D

Calcium Carbonate (From Coral Source) and
Vitamin D₃ (Colecalciferol)

Seacal-DX

Calcium Carbonate (From Coral Source)
and Vitamin D₃ (Colecalciferol)

Let assume a program have 100 instruction.

All are register mode.

⇒ CPU can read only one instruction from RAM at a time.

Then, CPU needs to access RAM 100 times.

∴ Number of access = 100 times.

⊗ Memory access time is 100 usec.

⇒ then total access time = 100×100 usec

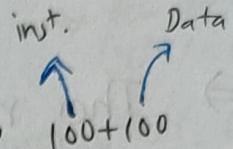
$$\therefore \text{Average memory access time} = \frac{\text{total access}}{\text{Total instruction}} = \frac{100 \times 100 \text{ usec}}{100} = \text{single access time}$$

⊗ If we don't have cache, then memory access time is the average memory access time. & depends on RAM.

$$\therefore \text{Average Access time} = \frac{\# \text{access} \times \text{Access time}}{\text{total instruction}}$$

⊗ Memory mode instruction:

OP Rx M Rz

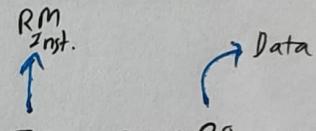


This time CPU needs to access memory
= 200 times.

$$\therefore \text{Average Access time} = \frac{200 \times 100 \mu\text{sec}}{100} = 200 \mu\text{sec.}$$

OP Rx Ry Rz \Rightarrow 70A.

OP Rx M Rz \Rightarrow 30A.



$$\text{then number of access time} = 70 + 30 + 30$$

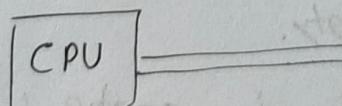
$$= 130 \text{ times.}$$

OP Rx M, Mz \Rightarrow 100 instruction

$$\therefore \text{number of memory access} = 100 + 200 = 300 \text{ times.}$$

inst.
Data Read

Locality of Reference:



RAM	
107	2-1
104	2
105	3
106	

⊗ Currently CPU is processing instruction-1, then,

which instruction will be read next?

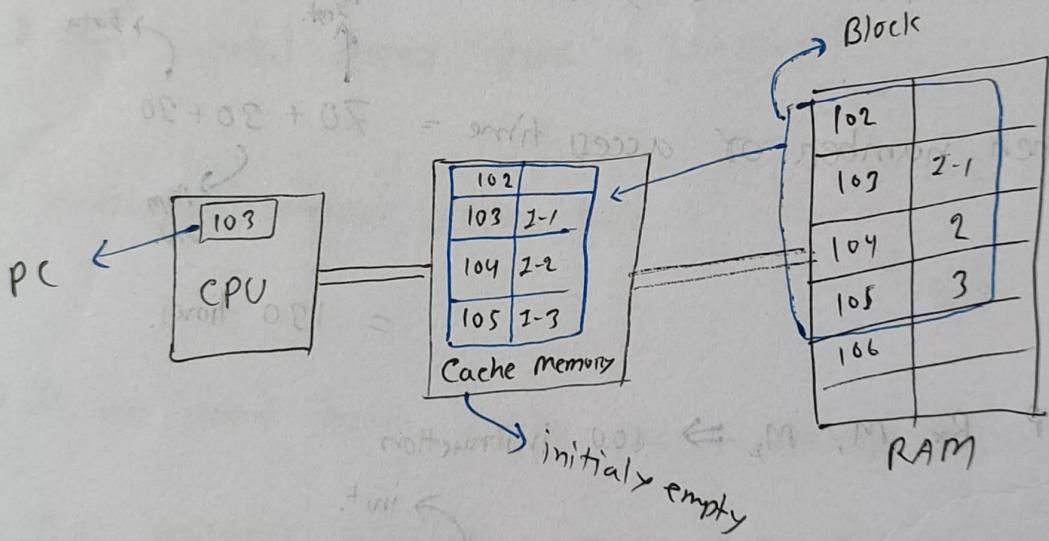
⇒ Then next memory address.

⇒ Instruction save next/close to current one.

↳ Spatial locality reference

⇒ Recently executed instruction

↳ Temporal locality of reference.



⊗ Steps:

i. Processor will search I_1 in cache by RAM address

103

→ initially empty.

ii. CPU does not find I_1 in cache: known as

cache miss.

- iii. In the event of cache miss, CPU will access RAM.
- (4/2/16/132)
- a section of RAM - capacity only few bytes.
- iv. A Block of RAM (including the instruction, CPU searching for) is transfer from RAM to Cache.

v. CPU Reads I_r from Cache memory.

- ⊗ And the PC will updated by the address of next instruction.
- ⊗ Whenever anything found in cache, known as cache hit.
- and in the event of cache hits CPU will simply read from cache memory.
- ⊗ Block Transfer depends on locality of Reference.

Seacal-D

Calcium Carbonate (From Coral Source) and Vitamin D₃ (Colecalciferol)

Seacal-DX

Calcium Carbonate (From Coral Source) and Vitamin D₃ (Colecalciferol)

* Let's assume that a program have 100 instruction.

80 instruction \Rightarrow Cache hit

Then, hit ratio = $\frac{80}{100} \times 100\% = 80\%$

= 0.8 (Normalized)

usually calculated
in %.

* If Hit ratio is 80%.

then miss ratio is 20%.

= 0.2 (Normalized)

* Then, without cache,

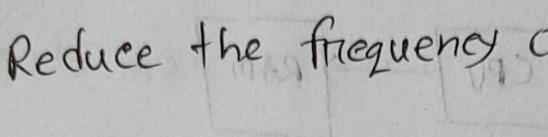
number of memory access is 100 times.

And with cache,

number of memory access is 20 times

* Reduce the frequency of access to RAM.

Cache Memory


 ☀ Objective: Reduce the frequency of Access to Main Memory.

☀ I = 200

70% \Rightarrow Register Mode

30% \Rightarrow LOAD & STORE

$$\begin{aligned} \text{∴ Number of access} &= 140 + 60 + 60 \\ &= 260 \text{ times.} \end{aligned}$$

↳ 70% Read Ins.
30% Read Ins.
30% Read/Write memory

☀ If cache memory used with Hit ratio 80%.

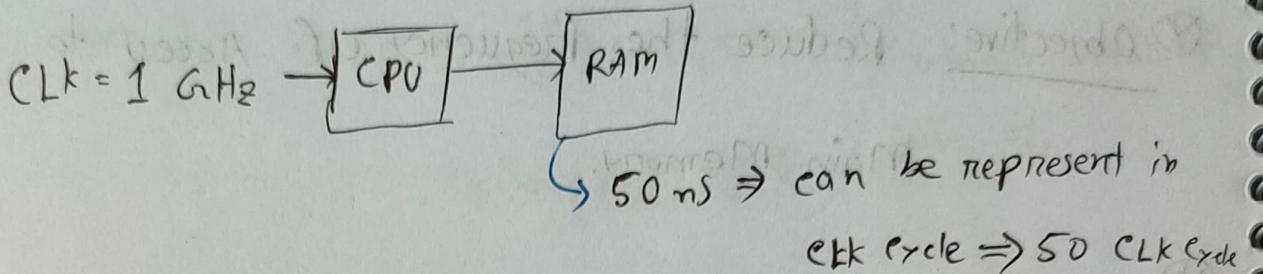
$$\text{Then, } 260 \times \frac{80}{100} = 208 \text{ (read from only cache)}$$

$$\begin{aligned} \text{∴ number of access memory} &= 260 - 208 \\ &= 52 \end{aligned}$$

$$\hookrightarrow \text{Alternative way} \Rightarrow 260 \times \frac{20}{100} = 52$$

{ Hit Ratio ↑
number of access ↓

* For using cache, Avg memory access time will be reduced.



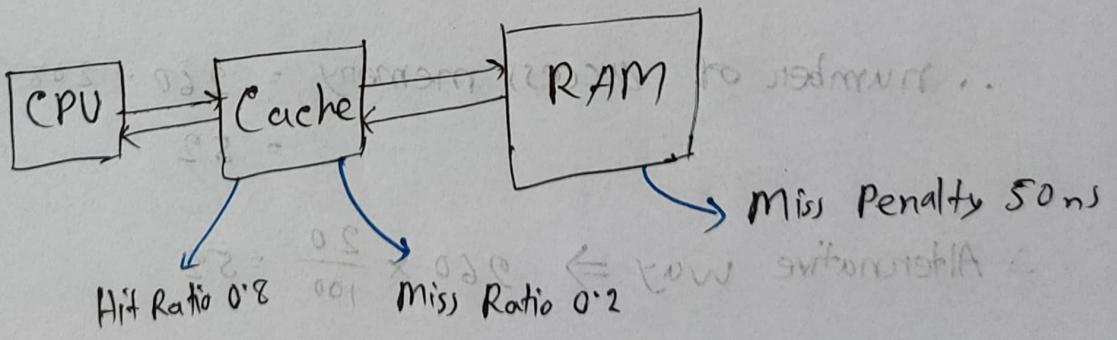
∴ Avg access time = 50 ns (Without Cache)

⇒ if cache memory used with 80% hit ratio and access time (Hit Time) 5 ns.

total inst. 1

$$\text{Avg normalized form} = 0.8 \times 5 + 0.2 \times 50 + 0.2 \times 5 \\ = 15 \text{ ns}$$

$$(\text{cache miss time}) \text{ 80\%} = \frac{0.8}{0.1} \times 0.2 \text{ ns}$$



$$\therefore \text{AMAT} = \text{Hit Time} + \text{Miss Rate} \times \text{Miss Penalty}$$

$$= 5 \text{ ns} + 0.2 \times 50 \text{ ns}$$

$$= 15 \text{ ns}$$

In terms of CLK \Rightarrow Avg CPI or base CPI

Effective CPI

\otimes Memory stall cycles = Memory Access \times miss penalty \times miss Ratio
 $=$ total for a program

\otimes Difference between, Memory stall cycles

vs AMAT

\otimes Base CPI = Hit Time

\otimes Avg CPI = Base CPI + Avg Instruction miss cycle

Avg instruction miss cycle + Avg data miss cycle

$$= 2 + 0.12 \times 10 + 0.2 \times 0.6 \times 10$$

$$= 13.38$$

Seacal-D

Calcium Carbonate (From Coral Source) and
Vitamin D₃ (Colecalciferol)

Seacal-DX

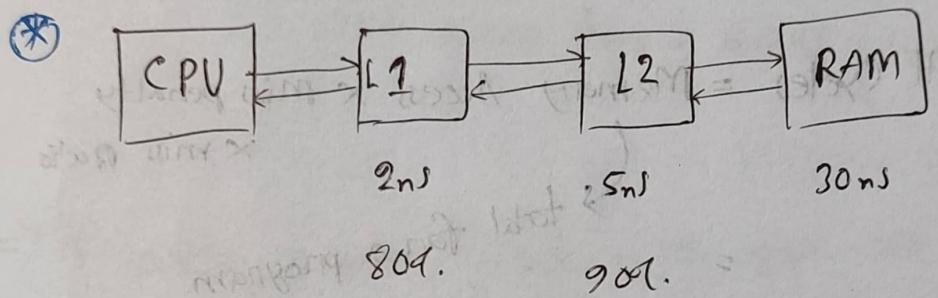
Calcium Carbonate (From Coral Source)
and Vitamin D₃ (Colecalciferol)

Multi level Cache!

⇒ CPU will always read from level 1.

Memory size: $L1 < 2 < 3 < RAM$

Speed: $L1 > 2 > 3 > RAM$



$$AMAT = 0.8 \times 2 + 0.2 \times 0.9 \times (2+5) + 0.02 \times (30+5+2)$$

$$= 3.6$$

$$\begin{aligned} \text{missing from Level 1} &= 1 - 0.8 \\ &= 0.2 \end{aligned}$$

$$\begin{aligned} \text{missing from Level 2} &= 0.2 - 0.2 \times 0.9 \\ &= 0.02 \end{aligned}$$

$$AMAT = \frac{L1}{\text{Hit Time}} + \frac{L1}{\text{Miss Rate}} \times \frac{L1}{\text{Miss Penalty}}$$

$$= \frac{\text{Hit Time}}{L1} + \frac{\text{Miss Rate}}{L1} \left[\frac{\text{Hit time}}{L2} + \frac{\text{miss rate} \times \text{miss}}{\text{penalty}} \right]$$

$$= 2 + 0.2 [5 + 0.1 \times 30]$$

$$= 3.6$$

L-27/01.11.2023/

* From Slide,

$$CP2 \text{ Stall} = CP2 \text{ base} + L1 \text{ inst miss cycle} + L1 \text{ data miss cycle}$$

$$= 2 + 1 \times 0.03 \times 100 + 0.4 \times 1.0 \times 100$$

$$= 9 \text{ cycles}$$

* Global Miss Rate

\Rightarrow in case of multi-level cache, the overall global miss rate is called

Global miss rate.

CP2 Stall = CP2 base + L1 inst miss cycle + L1 data miss cycle

+ L2 inst miss cycle + L2 data miss cycle

PSOI = mitoziol sklerozibba

istek N = blokla to ssir

$$(222-0) \times 35 = \frac{PSOI}{P} = 420/8 \rightarrow \text{medium}$$

Seacal-D

Calcium Carbonate (From Coral Source) and
Vitamin D₃ (Colecalciferol)

Seacal-DX

Calcium Carbonate (From Coral Source)
and Vitamin D₃ (Colecalciferol)

Cache Mapping

Three Type

- (i) Direct
- (ii) Fully Associative
- (iii) Set- Associative

* When a block of RAM for instr. or Data is copied into cache, where it is copied in cache?

\Rightarrow Transfer depends on first address of the block.

* Block : RAM divided into some equal section

RAM Size = 1 kB

Addressable location = 1024

Size of Block = 4 Bytes

$$\text{Number of Block} = \frac{1024}{4} = 256 (0-255)$$

* Blocks are identified by Block number.

(*) Starting address of block - 0 ?

⇒ 0

(*) What are the addressable locations block - 0 contains?

⇒ 0-3 (Decimal)

(*) Starting address of block - 1 ?

⇒ 4

Contains: 4-7

(*) Starting Address = Block number × Block Size

$$= 1 \times 4 = 4$$

(*) For block - 15: $15 \times 4 = 60$

(*) Address contains: $j \times \text{size of Block} \rightarrow j \times \text{size of Block} + \text{Size} - 1$

Starting address → Starting Address + Size - 1

Seacal-D

Calcium Carbonate (From Coral Source) and
Vitamin D₃ (Colecalciferol)

Seacal-DX

Calcium Carbonate (From Coral Source)
and Vitamin D₃ (Colecalciferol)

(*) Given memory address, calculate Block number that address belongs?

$$\Rightarrow \text{Block Number} = \frac{\text{Memory Address}}{\text{Size of a Block}} \quad (\text{integer part only})$$

(*) Cache Block Size = RAM Block Size \Rightarrow Block Number

can be called as Line

Line number

(*) Cache Size = 64 Bytes

Line/Block Size = 4 Bytes

\ Total number of Line, $m = 16$ (index, 0-15)

(*) Each Block/Line of cache must be shared by many blocks of RAM.

(*) How does CPU identify which block of RAM is currently occupying a particular block/line of cache?

(*) There is a Tag field for each block/line in cache.

(*) The first address of block from RAM will be saved in Tag.

(*) Direct Mapping formula:

$$\text{Line, } i = j \bmod m$$

total number of line
Block numbers of RAM

Example:

$$\begin{aligned}j &= 0 \Rightarrow i = 0 \\j &= 16 \Rightarrow i = 0 \\j &= 32 \Rightarrow i = 0\end{aligned}$$

multiple of m are fixed for line 0.

Seacal-D

Calcium Carbonate (From Coral Source) and Vitamin D₃ (Colecalciferol)

Seacal-DX

Calcium Carbonate (From Coral Source) and Vitamin D₃ (Colecalciferol)