

8/1/24

Computer Architecture : A Quantitative Approach (5th edition)
→ Hennessy

UNIT 1

FUNDAMENTALS OF QUANTITATIVE DESIGN & ANALYSIS

Classes of computer

PERSONAL MOBILE DEVICE

- should be cost effective, battery operated → minimal power consumption desired
- portable device, compact → area efficient
- key characteristics expected : responsiveness → minimal response time
predictability → maximum efficiency

DESKTOP COMPUTING

- used for performing repeated operations
- run some general-purpose applications

Performance {
 → computer performance } microprocessor choice
 → graphics performance }

SERVERS

- large scale enterprise computing
- availability, scalability, efficient throughput
- efficiency & cost effectiveness
- responsiveness to an individual request
⇒ key metric : how many requests can be handled in unit time

CLUSTERS

- group of computers connected through LAN

EMBEDDED COMPUTERS

- application specific
- all the requirements for the application are present within a given area (mantained on Si area)
- key characteristics : power, cost
- primary goal : meet performance requirements at a minimum price

Performance parameters

- area
- power
- speed

* larger area, more power is consumed in interacting with the components and the speed reduces (takes longer due to larger distance)

**PARALLELISM
IN APPLICATIONS****Data-level (DLP)**

arises because there are many data items that can be operated on at the same time

**Task/
Instruction-level (TLP)**

arises because tasks of work are created that can operate independently and largely in parallel

Pseudo parallelism \Rightarrow Pipelining

TASK - Smallest operation carried out by a device in order to meet the requirements of an application

12/22/2024

classmate
Date _____
Page 3

Computer hardware exploits application parallelism

INSTRUCTION-LEVEL PARALLELISM (ILP)
VECTOR ARCHITECTURES & GPU's (Unit 3)

THREAD-LEVEL PARALLELISM (Unit 4)

REQUEST-LEVEL PARALLELISM (Unit 4)
modularity
when stringent execution expected

DLP - Different instructions being performed on the same data items

ILP - Same instruction is given to multiple execution units simultaneously after reading the instruction only once - uses pipelining and speculative execution (Unit 2)

THREAD - Smallest active element that can be identified by the OS

VECTORIZATION - Try to perform similar operations for a set of data items for a large amount of data

TIGHTLY COUPLED HARDWARE - Expected response time is exactly matched; no relaxation

DECOPLED TASK - Break a task and handle simultaneously based on request

9/1/24

Harvard architecture

Von Neumann architecture

MICHAEL FLYNN Class 8: Parallelism & Parallel Architectures

Instruction Pool

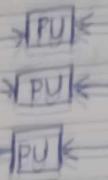
Data Pool



SISD Single Instruction, Single Data Stream

Instruction Pool

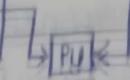
Data Pool



SIMD Single Instruction, Multiple Data stream

Instruction Pool

Data Pool



MISD Multiple Instruction, Single Data stream

Instruction Pool

Data Pool



MIMD Multiple Instruction, Multiple Data stream

Each processor fetches its own instructions & operates on its own & it targets task-level parallelism

* CISC \rightarrow Efficient pipelining is not possible

CA involves many aspects, including

- instruction set design
- functional organization, logic design, microarchitecture
- implementation

INSTRUCTION SET ARCHITECTURE (ISA)

\hookrightarrow Serves as the boundary between the software & hardware

(1) CLASS OF ISA

Register-memory ISAs which can access memory as part of many instructions
Eg: CISC 80x86

load-store ISAs which can access memory only with load/store instructions
Eg: RISC

(2) MEMORY ADDRESSING

Byte addressable

Word addressable An access of object of size s bytes at byte address A is aligned if $A \bmod s = 0$

RISC check alignment

- * Pure RISC machine \Rightarrow Every instruction is orthogonal to each other \Downarrow
 \Rightarrow Only 32 instructions possible misaligned data is not desired
- * ARM \Rightarrow 53 instructions possible

(3) ADDRESSING MODE
 \hookrightarrow the way in which an operand is being accessed by an instruction

In RISC/V/MIPS

- Register mov A, B
- Immediate (for constants) mov A, #50h
- Displacement Base address \pm offset (constant offset added) (BA)

In 80x86

- Absolute (no register) \rightarrow BA + di
- based indexed with displacement (2 registers)
- based with scaled index and displacement (2 registers; 1 reg \times size of operand in bytes)

(4) TYPES & SIZES OF OPERANDS

- (i) 8-bit ASCII character
- (ii) 16-bit Unicode character or half word
- (iii) 32-bit Integer or word
- (iv) 64-bit Double word or long integer
- (v) IEEE 754 floating point in 32-bit (single precision)
- (vi) IEEE 754 floating point in 64-bit (double precision)
- (vii) 80x86 supports 80-bit floating

(5) OPERATIONS

- data transfer
- arithmetic & logical
- control
- floating-point

(6) CONTROL FLOW INSTRUCTIONS

- (use PC-relative addressing)
- Conditional branches
 - Unconditional jumps
 - Procedure calls and returns

link register Holds return address

stack (in 80x86) for servicing interrupt

- push PC
- pop return address

10/11/24

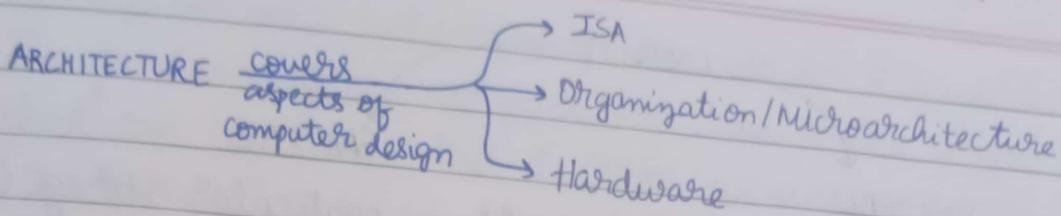
(7)

ENCODING AN ISA

FIXED LENGTH All ARM & MIPS instructions are 32 bits long, which simplifies instruction decoding

VARIABLE LENGTH 80x86 encoding is variable length, ranging from 1 to bytes

Advantageous : less space



Once ISA is fixed

- ⇒ ORGANIZATION : includes high-level aspects of a computer's design
- memory system
 - memory interconnect
 - design of internal processor (CPU)

* For single ISA, multiple microarchitecture possible

⇒ AMD Opteron Intel Core i7 implement x86 instruction set, have different cache

⇒ HARDWARE

Cache levels L1 → for processor
 L2, L3 → for outside processor

DESIGN OF EFFECTIVE ARCHITECTURE

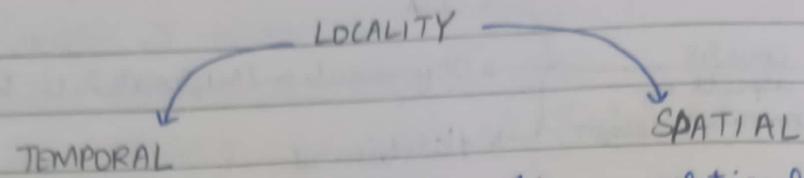
(1) Take advantage of parallelism

- ↳ Bit-level
- ↳ Instruction-level
- ↳ Data-level
- ↳ Thread-level

step, US WR
is observed
analyzed

(2) Principle of Locality

LOCALITY ⇒ predict with reasonable accuracy what instructions & data a program will use in the near future based on its accesses in the recent past



If an item is referenced then it will tend to be again referenced soon within relatively small time durations
⇒ loops in programs

If a particular storage location is referenced at a particular time then it is likely that nearby memory locations will be referenced in the near future
⇒ sequential code execution

(3) Focus on the Common Case

Make common case fast → performance enhanced

AMDAHL'S LAW

Defines the speedup that can be gained by using a particular feature

11/1/24

Speedup = Performance for entire task using the enhancement when possible

Performance for entire task without using the enhancement

Speedup = Execution time for entire task without using the enhancement

Execution time for entire task using the enhancement when possible

- Speedup for some enhancement depends on 2 factors:
- 1) The fraction of the computation time in the original computer that can be converted to take advantage of the (Fraction enhanced), is always less than or equal to 1.
 - 2) The improvement gained by enhanced execution mode, that is, how much faster the task would run if the enhanced mode were used for the entire program $\Rightarrow \text{Speedup}_{\text{enhanced}} \geq 1$

$$\text{Execution time}_{\text{new}} = \text{Execution time}_{\text{old}} \times \left\{ (1 - \text{Fraction}_{\text{enhanced}}) + \frac{\text{Fraction}_{\text{enhanced}}}{\text{Speedup}_{\text{enhanced}}} \right\}$$

Overall speedup \rightarrow Ratio of execution times

$$\text{Speedup}_{\text{overall}} = \frac{\text{Execution time}_{\text{old}}}{\text{Execution time}_{\text{new}}} = \frac{1}{(1 - \text{Fraction}_{\text{enhanced}}) + \frac{\text{Fraction}_{\text{enhanced}}}{\text{Speedup}_{\text{enhanced}}}}$$

Example 1 Suppose, we want to enhance the processor used for Web serving. The new processor is 10 times faster on computation in Web serving application than the original processor. Assuming the original processor is busy with computation 40% of the time & is waiting for I/O 60% of the time, what is the overall Speedup gained by incorporating the enhancement?

$$\text{Speedup}_{\text{enhanced}} = 10$$

\nearrow only computation performance can be enhanced

$$\text{Fraction}_{\text{enhanced}} = \frac{40}{100} = 0.4$$

$$\text{Speedup}_{\text{overall}} = \frac{1}{(1 - 0.4) + 0.4} = \frac{1}{0.6 + 0.04} = \frac{1}{0.64}$$

$$\boxed{\text{Speedup}_{\text{overall}} = 1.5625}$$

Example 2 Floating-point square root (FPSQR) \rightarrow 20% of execution time
 2 design alternatives proposed
 (i) Enhance FPSQR hardware and speedup the operation factor of 10

(ii) Try to make all FP instructions in the graphics processor run faster by a factor of 1.6; FP instructions are responsible for half of the execution time

$$(i) \text{ Speedup}_{\text{FPSQR}} = \frac{1}{\frac{(1+0.2)}{10} + 0.2} = \frac{1}{0.82} = 1.22$$

$$(ii) \text{ Speedup}_{\text{FP}} = \frac{1}{\frac{(1+0.55)}{1.6} + 0.55} = \frac{1}{0.8125} = 1.23$$

An

Improving the performance of FP operations overall is slightly better because of the higher frequency

11 Example 3 Assume a disk subsystem with the following components

MTTF:

- 10 disks, each rated at 1,000,000-hour MTTF
- 1 ATA controller, 500,000-hour MTTF
- 1 power supply, 300,000-hour MTTF
- 1 fan, 200,000-hour MTTF
- 1 ATA cable, 11,000,000-hour MTTF

Using the simplifying assumptions that the lifetimes are exponentially distributed & that failures are independent, compute the MTTF of the system as a whole

$$\text{Failure rate}_{\text{system}} = \frac{1}{\text{MTTF}_{\text{system}}} = \frac{1}{\frac{1}{1000000} + \frac{1}{500000} + \frac{1}{300000} + \frac{1}{200000} + \frac{1}{11000000}} = \frac{1}{0.23} = 4.35$$

Taking a reference of 1 billion wait time

Failure rate system = $\frac{0.3 \text{ BPP}}{1000000000 \text{ hours}}$ $= 3.000 \text{ FIT}$ Fail-in-time

$$\text{MTTF} = \frac{1}{\text{Failure rate}} = \frac{1}{3.000 \text{ FIT}} = 33333.3 \text{ hours} \approx 5 \text{ years}$$

16/12/24

Mean Time To Failure (MTTF)

Mean Time Between Failure (MTBF)

Mean Time to Repair (MTTR)

Example 4 Disk subsystems often have redundant power supplies to improve dependability. Using the components of MTTF from Example 3, calculate the reliability of redundant power supplies. Assume one power supply is sufficient to run the disk subsystem & that one redundant power supply is being added.

MTTF for redundant power supplies is the mean time until one power supply fails divided by the chance that the other will fail before the first one is replaced.

Mean time until one disk fails is $\frac{\text{MTTF}_{\text{power supply}}}{2}$

$$\text{MTTF}_{\text{power supply pair}} = \frac{\text{MTTF}_{\text{power supply}} / 2}{\text{MTTF}_{\text{power supply}}} = \frac{\text{MTTF}_{\text{power supply}}}{2 \times 0.9}$$

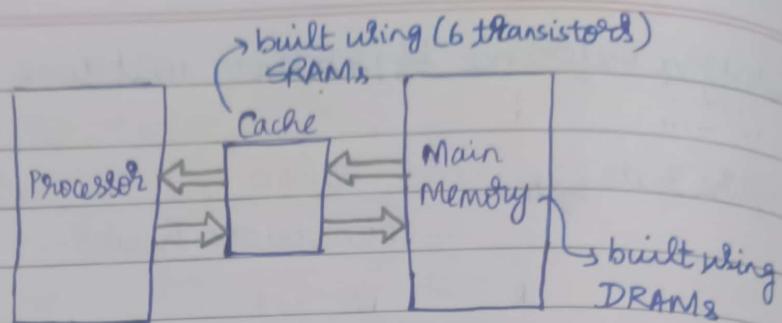
Assuming it takes the operator 24 hours to detect & repair the power supply after failure

$$= \frac{2000000}{2 \times 0.9} = 999999.9$$

$$\frac{833333333}{2000000} = 4166$$

No refreshing in SRAMs
→ faster access

MEMORY



As we move away from the Processor:

- access time increases
- storage capacity increases
- cost of storage per byte reduces

Operations

- Read (load) • contents unaltered
- Write (store) • contents altered

Processor places the address of the location on which the operation is to be performed on the address bus

If the operation is successful \Rightarrow HIT

failure \Rightarrow MISS \rightarrow more time consuming

↳ not found in cache

Memory Management Unit (MMU)

↓
go to Main Memory

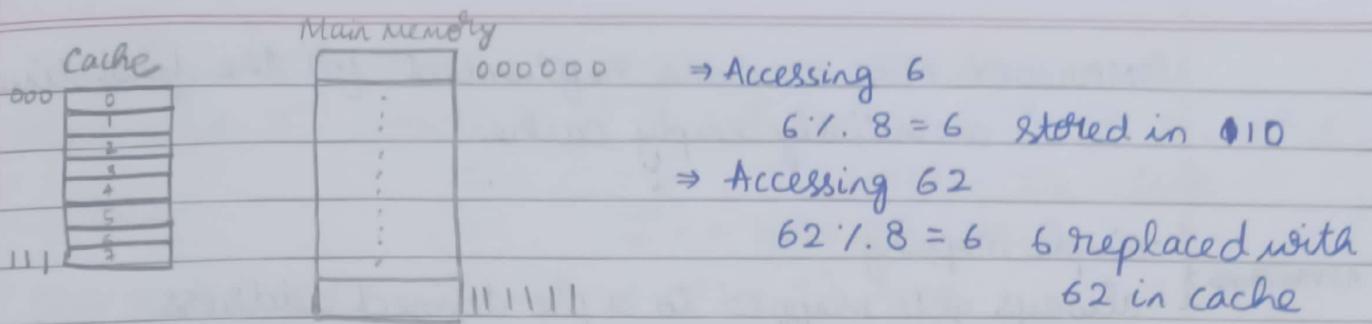
Access time

* Read operation takes lesser amount of time than write operation

assumed to have the content pertaining to the specific application

Mapping

↳ Direct Mapping (mod operation)



6 bits

3 bits	3 bits
--------	--------

 ↓ ↓
 Tag Index
 ↗ have to be matched
 in order to get a HIT

Valid bit → 0 by default

Protection bits
 ↗ ensures that there is
 no compromise in the data

Assume a cache with 8 blocks and the following references are made to the cache : 12d, 13d, 16d, 0d, 10d, 25d, 25d, 15d, 07d, 23d.

Identify the number of hits and misses using direct mapped cache

Initially empty cache (garbage values)

Cache	
0	16, 0, 16
1	25, 25
2	
3	
4	(12d)
5	(13d)
6	
7	(15d, 7d, 23d)

$$12 \bmod 8 = 4$$

placed at 100 ⇒ miss

$$13 \bmod 8 = 5$$

placed at 101 ⇒ miss

$$16 \bmod 8 = 0$$

placed at 000 ⇒ miss

$$0 \bmod 8 = 0$$

replaces 16 present at 000 ⇒ miss

$$16 \bmod 8 = 0$$

replaces 0's content ⇒ miss

000	000
010	000

Tag bits

16d content taken away & stored

in main memory

1 hit
9 misses

$$25 \bmod 8 = 1$$

25d content placed at 001 ⇒ miss

$$25 \bmod 8 = 1$$

already present at 001 ⇒ hit

$$15 \bmod 8 = 7$$

15d content placed at 111 ⇒ miss

$$7 \bmod 8 = 7$$

replaces 15d content at 111 ⇒ miss

$$23 \bmod 8 = 7$$

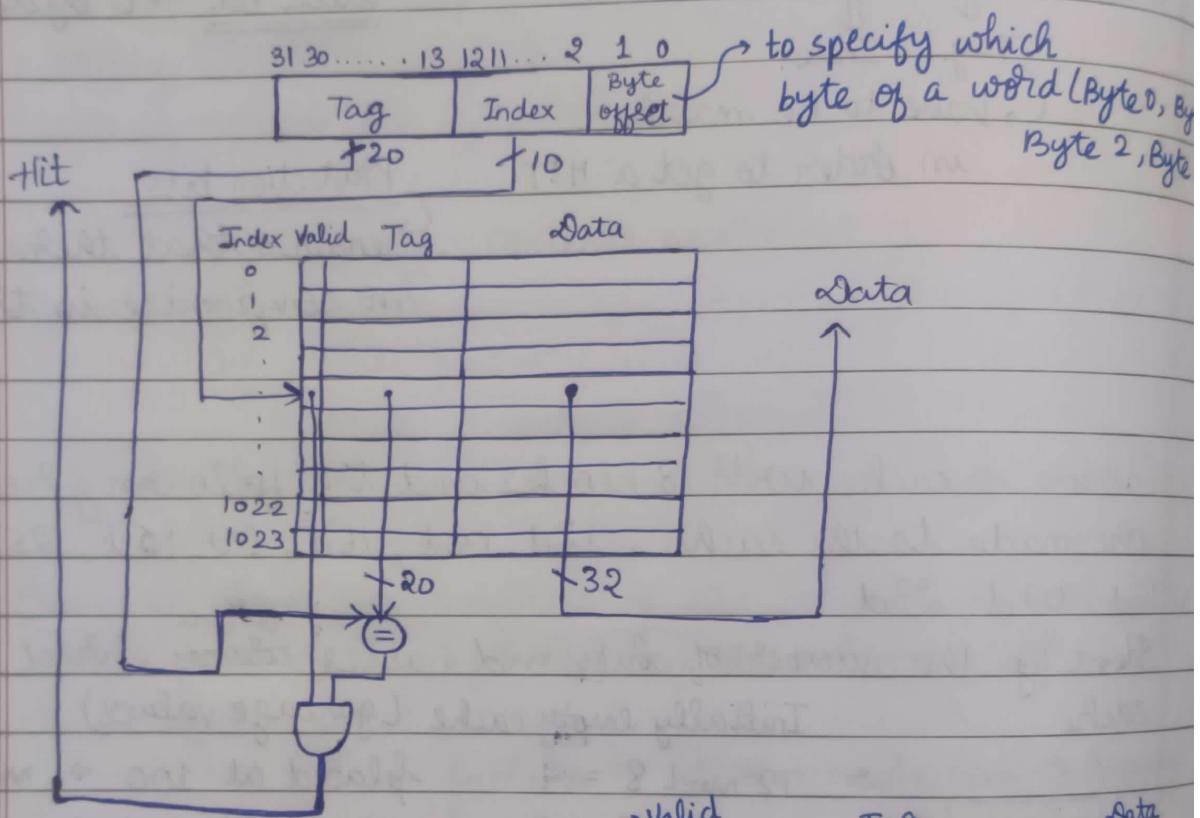
replaces 7d content at 111 ⇒ miss

* Even though there were other available free spaces, several misses encountered

COMPULSORY MISS \rightarrow Data referenced for the first time
in initially empty cache

- ADVANTAGES**
- simple mapping
 - always gets mapped to a pre-defined address
- DISADVANTAGE**
- available cache space is not efficiently utilized

17/1/24



Size of Cache

$$\begin{aligned}
 & \cancel{10 \times 1024} + 1 \times 1024 + 20 \times 1024 + 32 \times 1024 \\
 & = 54272 \text{ bits} \\
 & = 6784 \text{ bytes} \\
 & = 6.625 \text{ KB}
 \end{aligned}$$

$\hookrightarrow 2^{10}$
 $\Rightarrow 10.6$
 $\text{for } 10$

Actual Data Size = 32×1024 bits *
= 4 KB

Cache size 2^n blocks $\Rightarrow n$ bits used for index

Block size 2^m words = 2^{m+2} bytes $\Rightarrow m$ bits for the word within the block
 $2^2 = 4$ bytes per word

Size of tag field $32 - (n + m + 2)$ bits
 \downarrow
 byte offset

Total number of bits = $2^n \times (\text{block size} + \text{tag size} + \text{valid field size})$
 in direct-mapped cache

Q1. How many total bits are required for a direct-mapped cache with 16 KiB of data and 4-word blocks, assuming a 32-bit address?

~~Words = $4 - 2^2 \Rightarrow m = 2$ Block size = 4 words = $2^2 \Rightarrow m = 2$~~
 ~~$\frac{16}{4} = 4$ kilo words = 4096 words $\Rightarrow 2^{12}$ words~~

4 words per block

Total number of blocks = $\frac{2^{12}}{2^2} = 2^{10} = 1024$ blocks

$\Rightarrow 10$ bits used for index

Size of tag field = $32 - (10 + 2 + 2) = 32 - 14$
 $= 18$ bits

Total number of bits = $2^{10} \times (4 \times 32 + (32 - 10 - 2 - 2) + 1)$
 in direct-mapped cache
 $= 2^{10} \times 147$
 $= 147$ Kibibits $= 147/8$
 $= \underline{18.375}$ KiB for a 16 KiB cache

$\frac{18.375}{16} = 1.15$ times

Q2. Consider a cache with 64 blocks and a block size of 16 bytes.
To what block address does the byte address 1200 map?

$$64 \text{ blocks} = 2^6 \text{ blocks} \Rightarrow m = 6$$

$$\text{Block size } 16 \text{ bytes} = 4 \text{ words} = 2^2 \text{ words} \Rightarrow m = 2$$

$$\text{Size of tag field} = 32 - (6+2+2) = 22 \text{ bits}$$

Byte address 1200

$$\text{Block address } \frac{1200}{4} = 300$$

$$\text{Byte address in block } \frac{300}{4} = 75$$

per block per word
4 words 4 bytes

A1 64 blocks present in the cache

$$\Rightarrow 75 \bmod 64 = 11$$

in cache

∴ The byte address 1200 maps to block address 11.

* $(\text{Block address}) \bmod (\text{Number of blocks in cache})$

* Block address = $\frac{\text{Byte address}}{\text{Bytes per block}}$

* Block address is the block containing all addresses between

$$\left[\frac{\text{Byte address}}{\text{Bytes per block}} \right] \times \text{Bytes per block}$$

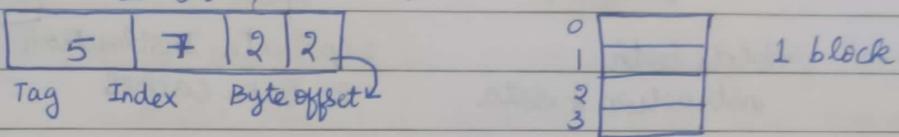
18/1/24

- Q3. Consider a cache with 128 blocks each with 16 bytes. Find the number of bits required to represent tag field, index field, block offset and byte field a) For a byte address 8200, find the mapped index value in cache. Assume 16-bit addressing.
(using direct-mapping)

$$128 \text{ blocks} = 2^7 \text{ blocks} \Rightarrow n = 7 \text{ Index field bits}$$

$$\text{Block size: } 16 \text{ bytes} = 4 \text{ words} = 2^2 \text{ words} \Rightarrow m = 2$$

16-bit \rightarrow Block offset



Byte offset \rightarrow 2 bits to specify which byte within a word

Block offset \rightarrow 2 bits to specify which of the 4 words in the block is being referred to

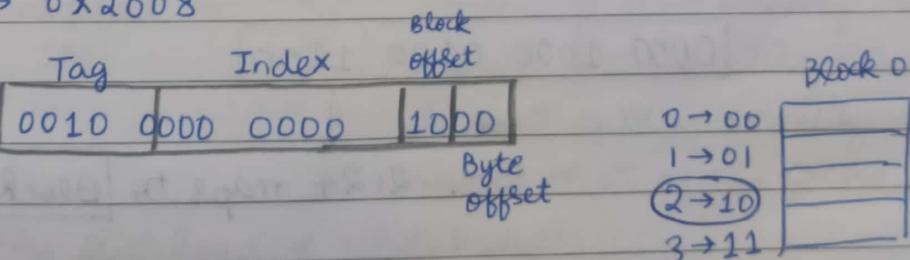
$$\text{Index} \rightarrow 7 \text{ bits} \because 128 \text{ blocks} = 2^7 \text{ blocks}$$

$$\text{Tag field} \rightarrow 16 - (7 + 2 + 2) = 5 \text{ bits}$$

$$\begin{array}{l} \text{8200 byte address} = 512 \cdot 5 \\ \text{16 bytes} \end{array} \quad 0.5_{10} \rightarrow 10_2$$

$$512 \bmod 128 = 0 \Rightarrow \text{Block 0}$$

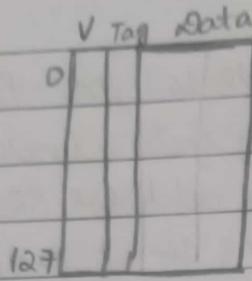
$$8200 \rightarrow 0x2008$$



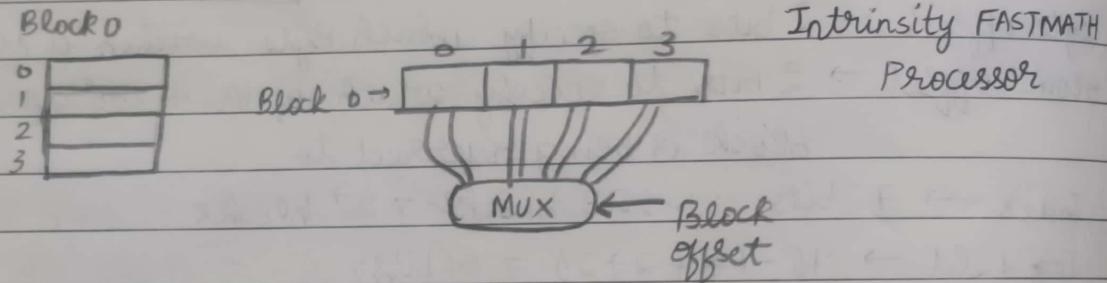
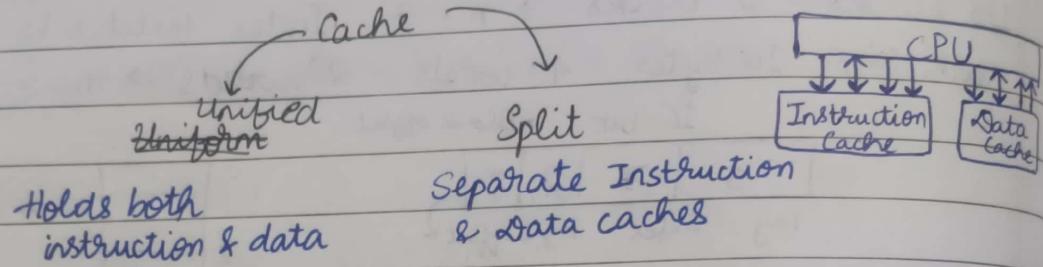
$$\text{Block offset value} = 10_2 \Rightarrow 2$$

8200 byte address maps to Block 0, Word 2

128 Validity bits (1 valid bit per block)



$$\begin{aligned}
 \text{Cache size} &= 2^7(1 + 5 + 4 \times 32) \\
 &= 2^7(134) \\
 &= 17152 \text{ bits} = 2144 \text{ bytes} \\
 &= \underline{\underline{2.09375 \text{ KiB}}}
 \end{aligned}$$

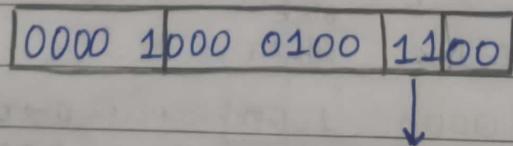


Q3. b) Byte address 2124 maps to ?

$$\frac{2124}{16} = 132.75$$

$$0.75_{10} \rightarrow 0.11_2$$

$$132 \bmod 128 = 4 \quad \text{Block 4}$$



$\therefore 2124 \text{ maps to } \boxed{\text{Block 4, Word 3}}$

22/1/24

Handling Writes

Store instruction

↳ Data written only into data cache → Data ^{in main memory} unchanged

INCONSISTENCY : Memory & cache have different data stored

Write Through Policy Change data in at one memory location, data changed ~~at~~ across ~~at~~ the memory hierarchy ; every change is accounted

WRITE BUFFER . Placed between processor and main memory

- Contains data to be written into cache
- Operations on cache need not be stalled
- Main memory data updated using data present in write buffer

Write Back Policy When a write occurs , the new value is written only to the block in the cache .

Dirty/Modified bit → 0 : Content of main memory matches (D) (M) with that of the cache

→ 1 : Content of Main memory mismatched with that of the cache

- no extra overhead of write buffer
- an extra bit in the cache is introduced
→ extra time due to miss

$$\text{Write-stall clock cycles} = \left(\frac{\text{Memory access}}{\text{Program}} \right) \times \text{Miss rate} \times \text{Miss penalty}$$

Separate instruction & data caches \rightarrow Harvard architecture

Example 1 Assume the miss rate of an instruction cache is 2% and miss rate of data cache is 4%. If a processor has CPI of 2, without any memory stalls & the miss penalty is 10 cycles for all misses, and the frequency of all loads & stores is 36%.

- 1) determine how much faster a processor would run w/ perfect cache that never missed
- 2) what happens if we speed up the machine by reducing CPI to 1 without changing the clock
- 3) what happens if we speed up the machine by doubling its clock rate but absolute time for miss penalty remains the same.

$$\text{CPU time} = \left(\frac{\text{CPU execution}}{\text{clock cycles}} + \frac{\text{Memory-stall}}{\text{clock cycles}} \right) \times \frac{\text{Clock cycle time}}{\text{time}}$$

$$\text{Memory-stall clock cycles} = \text{Read-stall cycles} + \text{Write-stall cycles}$$

	$\text{MOV } R_0, R_1 \rightarrow 2 \text{ cycles}$	Read contents of R_1 & write into R_0
store	$\text{MOV } [10], R_1 \rightarrow 2 \text{ cycles}$	
load	$\text{MOV } R_1, [10] \rightarrow 2 \text{ cycles}$	

- 1) Assuming instruction count = I
2 cycles per instruction

$$\text{Instruction miss cycles} = I \times 2\% \times 100 = 2I$$

\nearrow Miss Rate \searrow miss penalty

$$\text{Data miss cycles} = (I \times 36\%) \times 4\% \times 100 = 1.44I$$

\nearrow only a percentage of instructions access memory

$$\text{Total number of memory-stall cycles} = (2 + 1.44)I = 3.44I$$

$$\text{CPU time}_{\text{perfect}} = (2 \times I + 0) \text{ clock cycles} = 2I \text{ clock cycles}$$

↑ no stall

$$\text{CPU time}_{\text{stall}} = (2 \times I + 3.44 \times I) \text{ clock cycles} = 5.44I \text{ clock cycles}$$

$$\frac{\text{CPU time}_{\text{stall}}}{\text{CPU time}_{\text{perfect}}} = \frac{I \times \text{CPI}_{\text{stall}} \times \text{clock cycle}}{I \times \text{CPI}_{\text{perfect}} \times \text{clock cycle}}$$

$$= \frac{\text{CPI}_{\text{stall}}}{\text{CPI}_{\text{perfect}}} = \frac{5.44}{2}$$

$$= \underline{\underline{2.72}}$$

∴ Perfect cache's performance is better by 2.72.

- 2) Processor made faster (CPI changed to 1)

$$\begin{aligned}\text{CPU execution time} &= (1 \times I + 3.44 \times I) \text{ clock cycles} \\ &= 4.44I \text{ clock cycles}\end{aligned}$$

System with perfect cache would be $\frac{4.44}{1} = 4.44$ times as fast

- 3) Amount of execution time spent on memory stalls would have risen from $\frac{3.44}{5.44} = 63\%$. to $\frac{3.44}{4.44} = 77\%$.

Example 2 Assume we have a computer where CPI is 1 when all memory accesses hit in the cache. The only data accesses are loads and stores and these total 50% of the instructions. If the miss penalty is 25 clock cycles, and miss rate is 2%, how much faster would the computer be if all instructions were cache hits?

$$\text{CPI} = 1 \quad \text{Miss rate} = 2\% \quad \text{Miss penalty} = 25 \text{ clock cycles}$$

$$\begin{aligned}\text{Total number of } &= (I \times 50\%) \times 2\% \times 25 = 0.75I \\ \text{memory-stall cycles} &\quad \begin{matrix} \cancel{I} \\ \text{for instruction} \end{matrix} \quad \begin{matrix} \cancel{I} \\ \text{data} \end{matrix}\end{aligned}$$

$$\text{CPU time}_{\text{stall}} = (1 \times I + 0.75 \times I) \text{ clock cycles} = 1.75I \times \text{clock cycles}$$

$$\text{CPU time}_{\text{perfect}} = (1 \times I + 0) = 1I \times \text{clock cycles}$$

(all hits)

* In memory stall cycles $\frac{(1+0.5) \times IC}{\downarrow \text{data}}$

instruction

$$\frac{\text{CPU time}_{\text{stall}}}{\text{CPU time}_{\text{perfect}}} = \frac{1.75 \times I \times \text{Clock cycle}}{1 \times I \times \text{Clock cycle}} = 1.75$$

∴ Computer with no cache misses is 1.75 times faster

23/1/24

Problem 1 As

ble

all

of

See

FL

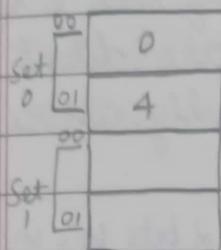
FULLY ASSOCIATIVE (FA) MAPPING

- Data can be placed anywhere in the cache
- The number of hits goes up, reducing the access time
- Searching for the desired data is harder due to random placement of the data → overhead
 - * search time increases
 - * hardware complexity increases
 - ↳ for example, more comparators are included to perform parallel checks
- makes use of cache memory efficiently
- hard to implement

0, 4, 8, 12, 0, 4, 8, 12
M M M H H H H

0
4
8
12

SET ASSOCIATIVITY MAPPING (2-way Set Associative)



- Divide the blocks into sets
 - Perform $Data \% \# \text{Sets}$ to obtain Set index
 - Place in one of the blocks of the set
 - reduced hardware complexity
- $\ast \# \text{Blocks in a set} = \text{power of } 2$
- $0, 4, 0, 4$
 M ↓ M ↓ M ↓ M ↓
 $0 \% 2 = \text{Set}0$ $4 \% 2 = \text{Set}0$ H H
 Place in 00 Place in 01

Cache Mapping

* FA is a special case of SA, where all of the blocks are considered to be a part of the same set

Problem 1 Assume there are 3 small caches, each consisting of 4 one-word blocks. One cache is fully associative, a second is 2-way set associative, and the third is direct-mapped. Find the number of misses for each cache organization given the following sequence of block addresses: 0, 8, 0, 6, 8.

FULLY ASSOCIATIVE

0
8
6

2-WAY SET ASSOCIATIVE

Set 0	0'6
0	8
Set 1	
1	

DIRECT MAPPING

0'8	0	8
	6	

0, 8, 0, 6, 8
M M H M H

Misses = 3

Hits = 2

0, 8, 0, 6, 8
M M H M H

$0 \bmod 2 = 0 \text{ set}$

$8 \bmod 2 = 0 \text{ set}$

$6 \bmod 2 = 0 \text{ set}$

Misses = 3

Hits = 2

0, 8, 0, 6, 8
M M M M M

$0 \bmod 4 = 0$

$8 \bmod 4 = 0$

$6 \bmod 4 = 2$

Misses = 5

Hits = 0

Problem 2 Assuming a cache of 4096 blocks, a 4-word block size & a 32-bit address, find the total number of sets & the total number of tag bits for caches that are direct mapped, 2-way & 4-way set associative & fully associative.

$$\text{Number of blocks} = 4096 = 4 \times 2^{10} = 2^{12} \Rightarrow m = 12 \text{ bits for index}$$

$$\text{Number of words per block} = 4 = 2^2$$

DIRECT MAPPING

$$\text{Number of tag bits} = 32 = (12 + 2 + 2) = 16 \text{ bits}$$

↑ 2 bits for block offset
↓ 2 bits for byte offset

$$\begin{aligned}\text{Total number of tag bits required} &= 16 \times \text{Number of blocks} \\ &= \frac{16 \times 4096}{1024} \\ &= \underline{\underline{64 \text{ kbits}}}\end{aligned}$$

2-way SA

$$\text{Number of sets} = \frac{4096}{2} = 2048 \text{ sets} \Rightarrow 11 \text{ bits}$$

$$\text{Number of tag bits} = 32 = (11 + 2 + 2) = 17 \text{ bits}$$

↑ byte offset
↓ block offset

$$\begin{aligned}\text{Total number of tag bits required} &= 17 \times \text{Number of sets} \times 2 \\ &= \frac{17 \times 2048 \times 2}{1024} \\ &= \underline{\underline{68 \text{ kbits}}}\end{aligned}$$

4-way SA

$$\text{Number of sets} = \frac{4096}{4} = 1024 \text{ sets} \Rightarrow 10 \text{ bits}$$

$$\text{Number of tag bits} = 32 = (10 + 2 + 2) = 18 \text{ bits}$$

$$\text{Total number of tag bits required} = \frac{18 \times 1024 \times 4}{1024} = \underline{\underline{72 \text{ kbits}}}$$

FA

$$\text{Number of sets} = \frac{4096}{4096} = 1 \Rightarrow 0 \text{ bits}$$

$$\text{Number of tag bits} = 32(2+2) = 28 \text{ bits}$$

$$\text{Total number of tag bits} = \frac{28 \times 4096 \times 1}{1024} = 112 \text{ Kbits}$$

LEAST RECENTLY USED (LRU) Computer Organisation by Carl Armature

0, 8, 0, 6, 8

Ref.	loaded pages ②	Referenced String
0	0 -	0,
8	0 8	0, 8
0	0 8	0, 8, 0
6	0 6	0, 8, 0, 6
8	8 6	0, 8, 0, 6, 8

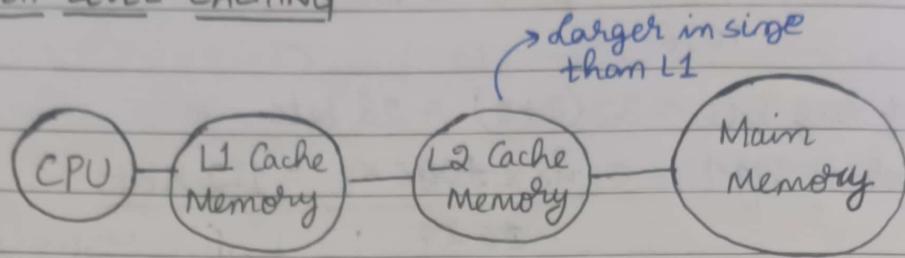
Example Consider the following memory references and find the number of page faults. Assume, the number of loaded pages is 4.

1, 1, 6, 9, 6, 8, 4, 1, 6, 9, 15, 17, 9, 1, 8

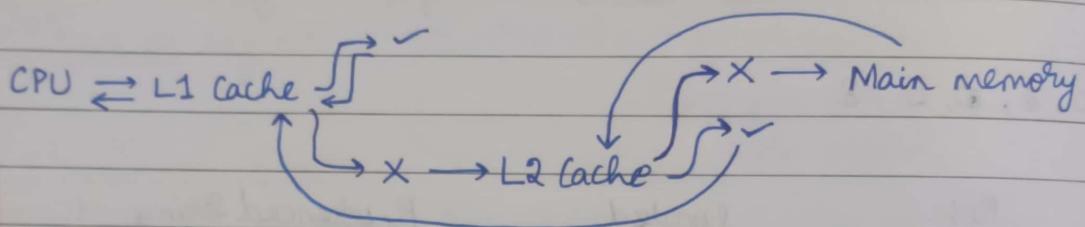
Page faults	loaded pages ④	Referenced String
①	1 1	1,
	1 1	1, 1
②	6 1 6	1, 1, 6
③	9 1 6 9	1, 1, 6, 9
④	6 1 6 9	1, 1, 6, 9, 6
⑤	8 1 6 9 8	1, 1, 6, 9, 6, 8
⑥	4 4 8 9 8	1, 1, 6, 9, 6, 8, 4



24/11/24

MULTI-LEVEL CACHING

MEMORY CONTROLLER UNIT manages the maintenance and accessing of the different levels of memory



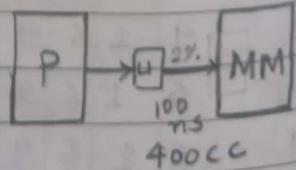
Problem Suppose we have a processor with a base CPI of 1, assuming all references hit in the primary cache, and a clock rate of 4GHz. Assume a main memory access time of 100 ns, including all the miss handling. Suppose, the miss rate per instruction at the primary cache is 2%. How much faster will the processor be if we add a secondary cache that has an access time for either a hit or miss & is large enough to reduce the miss rate to main memory to 0.5%?

$$\text{Clock cycle time} = \frac{1}{\text{clock rate}} = \frac{1}{4 \times 10^9} = 0.25 \text{ ns/clock cycles}$$

1 level of caching Miss penalty to main memory = $\frac{100 \text{ ns}}{0.25 \text{ ns}} = 400 \text{ clock cycles}$

$$\text{Total CPI} = 1 + \underbrace{2\% \times 400}_{\text{Memory-stall cycles per instruction}} = 9$$

Base CPI
Memory-stall cycles per instruction



2 levelsof caching Miss penalty for access to 2nd level cache

$$= \frac{5 \text{ ns}}{0.25 \text{ ns}} = 20 \text{ clock cycles}$$

clock cycle

1° stalls per instr.

$$\begin{aligned} \text{Total CPI} &= 1 + 2\% \times 20 + \\ &\quad 0.5\% \times 400 = 1 + 0.4 + 2 \\ &= 3.4 \quad \rightarrow 2^{\circ} \text{ stalls per instr.} \end{aligned}$$

P	$\xrightarrow{5 \text{ ns}}$	$\xrightarrow[2\%]{100 \text{ ns}}$	$\xrightarrow[0.5\%]{400 \text{ ns}}$	MM
	20 CC	400 CC		

Processor with secondary cache is faster by $\frac{9}{3.4} = \underline{2.64}$

(OR)

$$(2\% - 0.5\%) 20 = 0.3 \quad \text{Stall cycles of those references that hit in the secondary cache}$$

$$\begin{aligned} + \quad 0.5\% (20 + 400) &= 2.1 \quad \text{References that go to main memory (inclusive of 2° cache + main memory access time)} \\ + \quad \underline{1} \quad \text{Base CPI} \end{aligned}$$

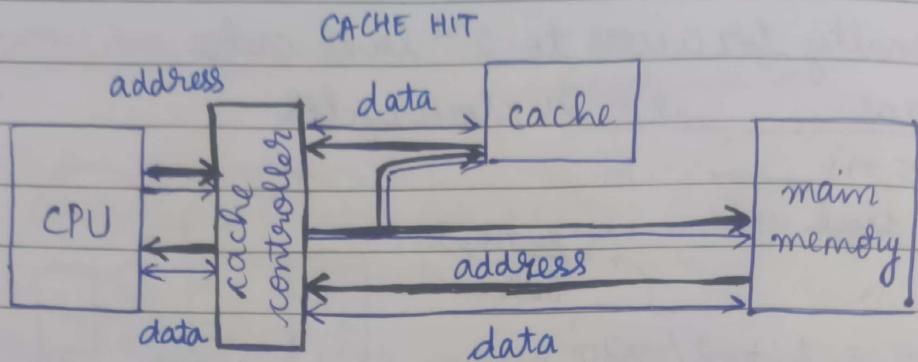
Advantages of Multilevel Cache Organisation

- Reduced access time
- Improved system performance
- Lower cost
- Energy efficiency

Disadvantages of Multilevel Cache Organisation

- Complexity
- Higher latency for cache misses
- Higher cost
- Cache coherence issues

29/1/24



- conflict (Collision)
- multiple memory locations mapped to the same cache location

CACHE MISS (3 Cs)

→ Compulsory: first access to a block
(cold start or process migration, first reference)

Capacity

- cache cannot contain all blocks accessed by the program

SOLUTION → increase cache size

SOLUTIONS

- increase cache size
- increase associativity

$$\text{Average Memory Access Time} = \text{Hit Time} + \text{Miss Rate} \times \text{Miss Penalty}$$

GOALS

- * Reducing Miss Rate
- * Reducing Miss Penalty
- * Reducing Hit Time

BASIC APPROACHES

Larger block size, larger cache size & higher associativity
Multilevel caches & higher read priority over writes

Avoid address translation when in the cache
↳ taken care of by the cache controller

Basic Cache Optimizations

(i) LARGER BLOCK SIZE

↓ compulsory misses; ↑ capacity, ↑ conflict misses, ↑ miss penalty

(ii) LARGER TOTAL CACHE CAPACITY TO REDUCE MISS RATE

↑ hit time, ↑ power consumption

(iii) HIGHER ASSOCIATIVITY

↓ conflict misses, ↑ hit time, ↑ power consumption

(iv) HIGHER NUMBER OF CACHE LEVELS

↓ overall memory access time

(v) GIVING PRIORITY TO READ MISSES OVER WRITES

↓ miss penalty

(vi) AVOIDING ADDRESS TRANSLATION IN CACHE INDEXING

↓ hit time

WAR Write After Read

RAW Read After Write

WAW White After Write

RAR Read After Read

MISS RATE IN MULTILEVEL CACHES

Local Miss Rate

→ Number of accesses in a cache

Total number of memory
accesses generated by processor

Global Miss Rate

→ Number of misses in cache

Total number of memory
accesses generated by processor

⇒ for L1 : MissrateL1

⇒ for L2 : Miss rate L1 × Miss rate L2

29)

- Q1. Using the data in Figure B.8 (Appendix B) determine whether a 32 KB 1-way set associative L1 cache is faster memory access time than a 32 KB 2-way set associative L1 cache. Assume the miss penalty to L2 is the access time of the faster L1 cache. Ignore misses beyond L2. Which has the faster average memory access time.

	2-way	4-way
Miss rate	0.038	0.037

Assume the access time on 4-way cache is 1.4 times longer

$$\text{Average Memory Access Time} = \text{Hit Rate} + \text{Miss Rate} \times \text{Miss penalty}$$

Access Time

Not mentioned in the question \Rightarrow Assume Hit rate is 1

2-way

$$1 + 0.038 \times 15 = 1.57 \text{ time units}$$

4-way

$$1.4 + 0.037 \times 15 = 1.807 \text{ time units}$$

$$\begin{array}{r} 1.4 \\ \times 15 \\ \hline 1.807 \end{array}$$

$$\frac{1.807}{1.57} = 1.15$$

- Q2. A system has 4GB of main memory, cache of 4KB where the cache has 64 Bytes per block. Write address fields for the following mapping techniques:

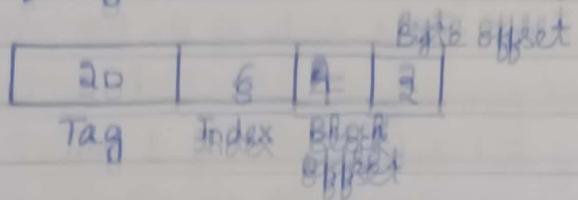
- a) Direct Mapping b) Fully Associative c) 4-way Set Associative
(Assume 32-bit addressing) d) 2-way Set Associative

$$\text{Block size} = 64 \text{ Bytes} = 2^6 \text{ Bytes} \Rightarrow \text{Index bits} = 6$$

$$\text{Number of blocks} = \frac{4 \times 1024}{64} = 64 \text{ blocks} = 2^6 \text{ blocks} \Rightarrow 6 \text{ index bits}$$

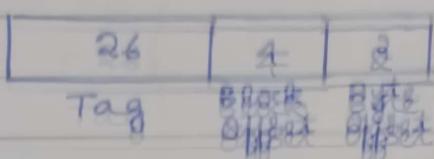
$$\text{Number of words per block} = \frac{64}{4} = 16 \text{ words} = 2^4 \text{ words} \Rightarrow 4 \text{ word bits}$$

a) Number of tag bits = $32 = (8+4+2) = 26$ bits



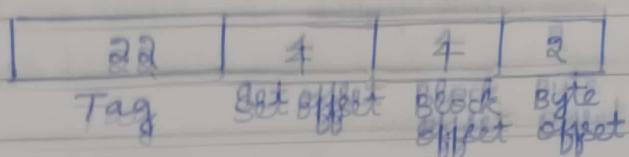
b) Number of sets = $\frac{64}{64} = 1 \Rightarrow 2^0$ 0 Index bits

Number of tag bits = $32 = (4+2) = 26$ bits



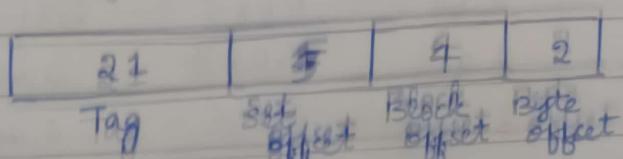
c) Number of sets = $\frac{64}{4} = 16$ sets $\Rightarrow 2^4$ 4 bits

Number of tag bits = $32 = (4+4+2) = 22$ bits



d) Number of sets = $\frac{64}{3} = 32$ sets $\Rightarrow 2^5$ 5 bits

Number of tag bits = $32 = (4+5+2) = 21$ bits



2- 30/1/24

- Q1. Find the average memory access time for a processor with 2ns clock cycle time, a miss rate of 0.04 misses per instruction, a missed penalty of 25 clock cycles, and a access time (including hit detection) of 1 clock cycle. Assume that the read & write miss penalties are the same and ignore other write stalls.

$$\text{Average memory access time} = 1 + 0.04 \times 25 = 2 \text{ clock cycles}$$

$$1 \text{ clock cycle} = 2 \text{ ns}$$

$$\therefore 2 \text{ clock cycles} = \underline{\underline{4 \text{ ns}}}$$

- Q2. Assume the memory system takes 80 clock cycles of overhead and then delivers 16 bytes every 2 clock cycles. Thus, supply 16 bytes in 82 clock cycles, 32 bytes in 84 clock cycles & so on. Which block size has the smallest average memory access time for each cache size?

CACHE SIZE

BLOCK SIZE	4K	16K	64K	256K
16	8.57%	3.94%	2.04%	1.09%
32	7.24%	2.87%	1.35%	0.71%
64	7%	2.64%	1.06%	0.51%
128	7.78%	2.77%	1.02%	0.49%
256	9.51%	3.29%	1.15%	0.49%

Block Size	16	32	64	128	256
Miss penalty	$\frac{16}{16} \times 2 + 80$	$\frac{32}{16} \times 2 + 80$	$\frac{64}{16} \times 2 + 80$	$\frac{128}{16} \times 2 + 80$	$\frac{256}{16} \times 2 + 80$
	= 82	= 84	= 88	= 96	= 112

Average memory access time = Hit time + Miss rate \times Miss penalty
 $= 1 + (8.57\% \times 82) = 8.0274$ clock cycles

for a 256 block size of 16 bytes & cache size of 4KB

For 256 B block size & 256 KB cache size

Average memory access time = $1 + (0.49\% \times 112)$
 $= 1.5488$ clock cycles $\frac{256}{16} \times 2 + 80$
 overhead

For 128 B Block size & 16KB cache size

~~AMAT = $1 + \left\{ 2.77\% \times \left(\frac{128}{16} \times 2 + 80 \right) \right\}$~~
 = $1 + \left\{ 2.77\% \times \left(\frac{128}{16} \times 2 + 80 \right) \right\} = 1 + (2.77\% \times 96)$
 $= 1 + 2.6592 = 3.6592$ clock cycles

For 32B Block size & 64 KB cache size

AMAT = $1 + \left\{ 1.35\% \times 84 \right\} = 2.134$ clock cycles

For 64B Block size & 4KB cache size

AMAT = $1 + (7\% \times 88) = 7.16$ clock cycles

Q3. Assume that higher associativity would increase clock cycle time as:

Clock-cycle time 2-way = $1.36 \times$ clock cycle time, way

The clock cycle time for direct mapped cache on hit is 1cc & for 2-way it increases to $1.36 \times$ direct mapped...
 4-way = $1.44 \times$ clock cycle time, way
 8-way = $1.52 \times$ clock cycle time, way

Assume that hit time = 1 clock cycle

Miss penalty for direct mapped case = 25 clock cycles to a L2 cache that never misses

Miss rate penalty need not be

rounded to an integral number. Which statement is true?

AMAT_{8-way} < AMAT_{4-way}

AMAT_{4-way} < AMAT_{2-way}

AMAT_{2-way} < AMAT_{1-way}

(only hit time is getting affected, not miss - penalty)

4KB cache size

1-way $1 + (0.098 \times 25) = 3.45$

2-way $1 + (0.076 \times 25) \times 1.36 = 3.584$

4-way $1 + (0.071 \times 25) \times 1.44 = 3.556$

8-way $1 + (0.071 \times 25) \times 1.52 =$

4KB cache size

1-way $1 + (0.098 \times 25) = 3.45$

2-way $1.36 + (0.076 \times 25) = 3.26$

4-way $1.44 + (0.071 \times 25) = 3.215$

8-way $1.52 + (0.071 \times 25) = 3.295$

For 4KB cache size, $\text{AMAT}_{8\text{-way}} < \text{AMAT}_{4\text{-way}}$ X
 $\text{AMAT}_{4\text{-way}} < \text{AMAT}_{2\text{-way}}$ ✓
 $\text{AMAT}_{2\text{-way}} < \text{AMAT}_{1\text{-way}}$ ✓

- Q4. Suppose that in 1000 memory references, there are 40 misses in L1 cache & 20 misses in L2 cache. What are the various miss rates? Assume miss penalty from L2 cache to memory is 200 clock cycles. Hit time of L2 cache is 10 clock cycles, hit time of L1 is 1 clock cycle & there are 1.5 memory references per instruction. What is average memory access time & average stall cycles per instruction? Ignore impact of writes.

	L1	L2
Miss rate	$\frac{40}{1000} = 4\%$	$\frac{20}{1000} = 2\% \text{ (global)}$

$$\frac{20}{40} = 50\% \text{ (local)}$$

Average memory access time = hit time_{L1} + miss rate_{L1} × (hit rate_{L2} + miss rate_{L2} × miss penalty_{L2})

$$= 1 + \frac{4}{100} \left(10 + \frac{50}{100} \times 200 \right) = \underline{\underline{5.4 \text{ clock cycles}}}$$

1000 memory ref. ≈ 667 instructions

1.5 mem ref per instr.

$$40 \times 1.5 = 60 \text{ L1 misses}$$

$$20 \times 1.5 = 30 \text{ L2 misses}$$

per 1000 instructions

$$\begin{aligned} \text{Average memory} &= \text{Misses per instr.}_{L_1} \times \text{Hit time}_{L_2} + \text{Misses per instr.}_{L_2} \\ \text{stalls per instruction} &\quad \times \text{Miss penalty}_{L_2} \\ &= \frac{60}{1000} \times 10 + \frac{30}{1000} \times 200 \\ &= 6.6 \text{ clock cycles} \end{aligned}$$

→ AMAT
 $(5.4 - 1) \times 1.5 = 4.4 \times 1.5 = 6.6 \text{ clock cycles}$

L1 hit avg. number of
time mem ref per instr.

31/1/24

Memory-stall clock cycles = (Read-stall cycles + Write-stall cycles)

$$\text{Read-stall cycles} = \frac{\text{Reads}}{\text{Program}} \times \frac{\text{Read miss rate}}{} \times \frac{\text{Read miss penalty}}{}$$

$$\begin{aligned} \text{Write-stall cycles} &= \left(\frac{\text{Writes}}{\text{Program}} \times \frac{\text{Write miss rate}}{} \times \frac{\text{Write miss penalty}}{} \right) \\ &\quad + \text{Write buffer stalls} \end{aligned}$$

$$\text{Memory-stall clock cycles} = \frac{\text{Memory accesses}}{\text{Program}} \times \frac{\text{Miss rate}}{} \times \frac{\text{Miss penalty}}{}$$

$$= \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Misses}}{\text{Instruction}} \times \frac{\text{Miss penalty}}{}$$

$$\begin{aligned} \text{Memory stall clock cycles} &= 1.5 \times \frac{40}{1000} \times 10 + 1.5 \times \frac{60}{1000} \times 200 = 6.6 \text{ clock cycles} \end{aligned}$$

Associativity → increases hit time
→ decreases miss rate

CLASSMATE
Date _____
Page 36

Q5. Given the data below, what is the impact of L2 cache associativity on its miss penalty? 1/2/24

- hit time_{L2} for direct mapped = 10 clock cycles
- 2-way set associativity increases hit time by 0.1 clock cycle, 2)
10.1 clock cycles
- Local miss rate_{L2} for direct mapped = 25%.
- Local miss rate_{L2} for 2-way set associative = 20%.
- Miss penalty_{L2} = 200 clock cycles

For a direct mapped L2 cache

$$\text{L1 cache miss penalty} \quad \text{Miss penalty}_{\text{2-way } L2} = 10 + 25\% \times 200 \\ = 60 \text{ clock cycles}$$

$$\text{Miss penalty}_{\text{2-way } L2} = 10.1 + 20\% \times 200 \\ = 50.1 \text{ clock cycles}$$

L2 hit time has to be an integral number of clock cycles
→ 10.1 can be approximated to 10 $\Rightarrow 10 + 20\% \times 200 = 50 \text{ CC}$
 $11 \Rightarrow 11 + 20\% \times 200 = 51 \text{ CC}$

either
an improvement
direct-mapped cache

1/2/24

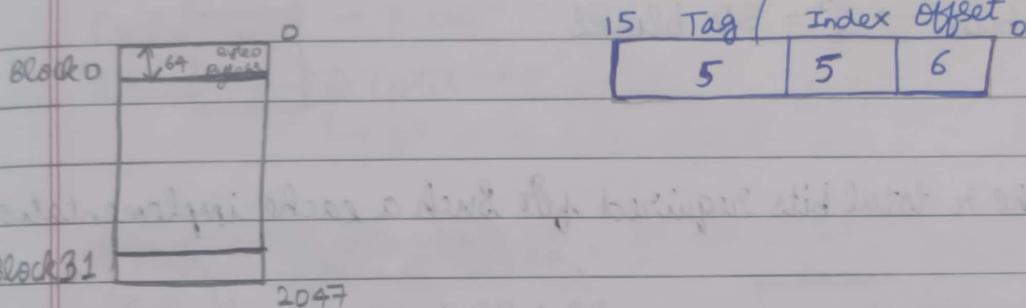
Problems on Cache (Question Bank)

2) 16-bit memory addresses

2KB cache ; 64 bytes per cache block.

Size of each memory word = 1 byte

a) Number of blocks = $\frac{2 \times 1024}{64} = 32 = 2^5 \Rightarrow 5 \text{ bits for index}$
 $16 - (5+6) = 16 - 11 = 5$



b) Sequential access of 128, 144, 2176, 2180, 128, 2176

Initially empty cache

Block

~~128 mod 32 = 0~~

M

Block 0

~~128~~

~~2176~~

~~2176~~

~~63~~

~~127~~

~~191~~

~~144 mod 32 = 16~~

M

Block 4

~~2180~~

~~2180 mod 32 = 4~~

M

Block 16

~~144~~

~~128 mod 32 = 0~~

M

Block 31

~~144~~

★ $128 \Rightarrow 0 \times 80$

0000 0000 1000 0000

(M)

Block 2

★ $144 \Rightarrow 0 \times 90$

0000 0000 1001 0000

(H)

Block 2

★ $2176 \Rightarrow 0 \times 880$

0000 1000 1000 0000

(M)

Block 2

Same offset as
128 \Rightarrow Block 2 removed
& 2176 stored (2176-2239)

★ $2180 \Rightarrow 0 \times 884$

0000 1000 1000 0100

(H)

Block 2

Block 2 already present
 \Rightarrow offset + empty
(2176-2239)

★ $128 \Rightarrow 0 \times 80$

0000 1000 1000 0000

(M)

Block 2

∴ Currently has 2176-2239 memory addresses

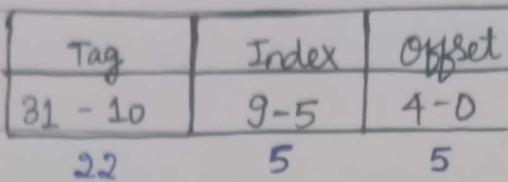
★ $2176 \Rightarrow 0 \times 880$

0000 1000 1000 0100

(M)

Block 2

∴ Currently has 128-191 memory addresses

1) 

Tag 31 - 10	Index 9-5	Offset 4-0
22	5	5

Direct-mapped cache
32-bit address

a) Cache block size (in words)

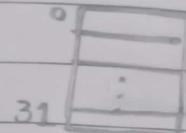
$$\text{Offset} = 5 \text{ bits} \rightarrow 2^5 = 32 \text{ bytes}$$

$$= 8 \text{ words}$$

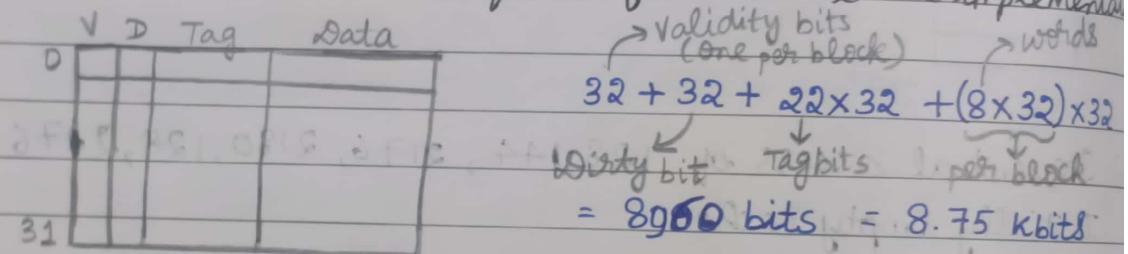
Index = 5 bits

b) How many entries does the cache have

$$\text{Index} = 5 \text{ bits} \rightarrow 2^5 = 32 \text{ blocks}$$



c) What is the total bits required for such a cache implementation?



3) For 2), consider 2-way set associativity

$$\text{Number of sets} = \frac{\text{Number of blocks}}{2} = \frac{32}{2} = 16 \Rightarrow 2^4$$

Tag	Set offset	2
6	4	6

Block offset → Set 2

- (M) 128 $\Rightarrow 0 \times 80$
- (H) 144
- (M) 2176
- (H) 2180
- (H) 128
- (H) 2176

0000	0000	1000	0000
------	------	------	------

Set 2

128
191
2176
2239

5/2/24

~~For a direct mapped cache, 32-bit addressing~~

- 5) 64-bit processor has a 8 MB, 4-way set associative cache with 32-byte blocks. How is the address arranged in terms of set, block and offset bits.

* Processor capacity is determined by data bus width and not address bus width

Cache size \rightarrow 8 MB

1 block \rightarrow 32 bytes \rightarrow 8 words

$$\hookrightarrow 2^5 \Rightarrow 5 \text{ bits offset}$$

$$\frac{8 \text{ MB}}{32 \text{ B}} = \frac{8 \times 2^{20}}{2^5} = 2 \times 2^{15} = 2^3 \times 2^5 \times 2^{10} = 256 \times 2^{10}$$

= 256 K Blocks

$$= \frac{256 \times 2^{10}}{4}$$

= 64 K sets

Tag	Index	Offset
11	16	3 2

$$2^6 \times 2^{10} \Rightarrow 16 \text{ bits for Index}$$

$$\text{Tag} = 32 - (16 + 3 + 2) = \underline{\underline{11 \text{ bits}}}$$

- 7) Consider a direct mapped cache of size 16 KB with block size 256 bytes. The size of main memory is 128 KB. Find the number of tag bits in tag field.

Cache size = 16 KB

$$\text{Number of blocks} = \frac{16 \text{ KB}}{256 \text{ B}} = \frac{2^4 \times 2^{10}}{2^8} = 2^6 \text{ blocks} \Rightarrow 6 \text{ bits for Index}$$

Block size = 256 bytes = 64 words

$$2^8 \Rightarrow \text{offset} = 8 \text{ bits}$$

Offset
6 2

Main memory size = 128 KB = $2^7 \cdot 2^{10} \Rightarrow 17$ -bit addressing

$$\text{Tag bits} = 17 - (6 + 8) = \underline{\underline{3 \text{ bits}}}$$

Tag	Index	Offset
16	3	8

- 8) Consider a direct mapped cache of size 512 KB with block size 1KB. There are 7 bits in the tag. Find size of main memory.

Block size = 1KB = 1024 bytes = 2^{10} \Rightarrow 10 bits for tag
 Number of blocks = $\frac{512\text{ KB}}{1\text{ KB}} = 2^9$ \Rightarrow 9 bits for index

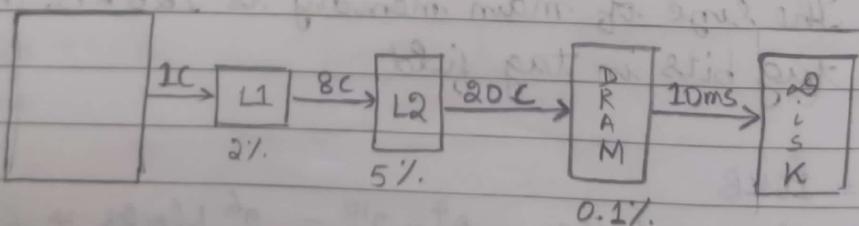
Total number of bits = $7 + 9 + 10 = 26$ bits

25		
Tag	Index	Offset
7	9	10

Size of main memory = $2^{26} = 2^6 \cdot 2^{20} = 64\text{ MB}$

- 9) Given the following data and assuming a clock rate of 1000 MHz.

Memory	Hit Time	Miss Rate	Calculate the average memory access time
L1	1 cycle	2%	
L2	8 cycles	5%	
DRAM	20 cycles	0.1%	
Disk	10ms	0.1%	



Clock rate = 1000 MHz

Clock cycles = $\frac{1}{1000 \times 10^6} = 1\text{ ns} \Rightarrow 1\text{ clock cycle}$

1ns \rightarrow 1 CC

10ms \rightarrow ? $\Rightarrow \frac{10\text{ ms}}{1\text{ ns}} = \frac{10 \times 10^{-3}}{10^{-9}} = 10^7 = 10\text{ M clock cycles}$

$$\begin{aligned}\text{Average memory access time} &= 1 + 2\% \cdot (8 + 5\%) \cdot (20 + 0.1/(10 \times 10^6)) \\ &= \underline{\underline{11.18 \text{ cycles}}}\end{aligned}$$