

A memory system has 1M words. Each word is an 8-bit byte. The memory is organized into blocks of 8 words each. The cache has 256K words, organized into cache lines of 8 words each. The memory cache is organized into 4-way set associative cache. Which of the following statements are true?

(a) 13 bits are needed for cache set address.

Ans: 快取(cache)有256K個 words，並且是4-way set associative cache。每個快取線(cache line)有8個 words。這意味著有 $256K/8 = 32K$ 個快取線，分為4路，因此會得到8K為1組cache set address。要地址8K組所需的位數是 $\log_2(8K) = 13$ bits

(b) 4 bits are for Tag.

Ans: Tag 用於區分可以映射到同一組的不同區塊(block)。在這種情況下，內存大小為1M words，每個區塊是8個 words，我們有 $1M/8 = 128K$ 個區塊。

同時，我們有8K組，每個區塊映射到8K組的其中一個。每個組的唯一區塊數 = $128K$ 區塊 / $8K$ 組 = 每組16個區塊。然而，每個組只能容納4個區塊 (4-way set associative cache)，所以每個組的其他12個區塊需要使用標籤來區分。

在這種情況下，標籤需要對16個blocks進行位址，這需要 $\log_2(16) = 4$ bits。

(c) 21 bits are needed to address all words.

Ans: 內存中的words總數為1M。要位址1M的位置，我們需要 $\log_2(1M) = 20$ bits，不是21 ($1M \approx 2^{20}$)

(d) 17 bits are needed to address all blocks.

Ans: blocks of 8 words each, 所以，內存中的區塊(block)總數是 $1M / 8 = 128K$ 。要位址128K的位置，我們需要 $\log_2(128K) = 17$ bits

(e) None of the above.

For a typical disk without considering queuing delay, the average seek time is 5 milliseconds, and the transfer rate is 2M-bytes per second. The disk rotates at 600 RPM (Revolution Per Minute), and the controller overhead is 0.5 millisecond. Determine the average time to read a 1792-bytes sector.

Answer:

1. Seek Time: This is the time it takes for the read/write head to move to the track where the data is located. It is given as 5 milliseconds.

2. Rotational Delay: This is the time it takes for the disk to rotate until the desired sector is under the read/write head. On average, it takes half a rotation for the desired sector to be accessible. Given that the disk rotates at 600 RPM (Revolutions Per Minute), the time for one full rotation is $(1 \text{ minute} / 600) = 0.1 \text{ seconds} = 100 \text{ milliseconds}$. Therefore, the average rotational delay (half a rotation) is $100/2 = 50 \text{ milliseconds}$

3. Transfer Time: This is the time it takes to actually read the data. Given that the sector size is 1792 bytes and the transfer rate is 2M bytes per second, the transfer time can be computed as $(1792 \text{ bytes} / 2\text{M bytes/second}) = 0.000896 \text{ seconds} = 0.896 \text{ milliseconds}$.

4. Controller Overhead: This is the time spent on tasks like error checking and correction, which is given as 0.5 milliseconds.

Total Time = Seek Time + Rotational Delay + Transfer Time + Controller Overhead

= 5 ms + 50 ms + 0.896 ms + 0.5 ms

= 56.396 milliseconds

Assume that the miss rate of an instruction cache is 3% and the miss rate of the data cache is 6%. If a processor has a CPI of 4 without any memory stalls and the miss penalty is 100 cycles for all misses. Assume the frequency of all loads and stores is 30%. How much faster a processor will run with a perfect cache that never missed.

Answer:

- Original:

- Instruction fetches: $0.03 * 100 = 3$ cycles(stall)
- Data loads / stores: $(stall) = 0.30 * 0.06 * 100 = 1.8$ cycles(stall)

Total miss penalty: 3 cycles + 1.8 cycles = 4.8 cycles

CPI with memory stalls: 4 cycles + 4.8 cycles = 8.8 cycles

- Perfect:

Only base CPI: 4 cycles

Speedup = Total CPI with original cache / Total CPI with perfect cache

$8.8 / 4 = 2.2$

2. (25 points) The execution of an instruction can be divided into five parts: instruction fetch (IF), register read (RR), ALU operation (EX), data access (MEM), and register write (RW). The following Table 1 shows the execution time of each part for several types of instructions, assuming that the multiplexors, and control unit have no delay.

	Instruction fetch	Register read	ALU operation	Data access	Register write
Load word (lw)	200 ps	100 ps	200 ps	200 ps	100 ps
Store word (sw)	200 ps	100 ps	200 ps	200 ps	
R-format (add, sub, and, or, slt)	200 ps	100 ps	200 ps		100 ps
Branch (beq)	200 ps	100 ps	200 ps		

If instructions are to be executed in a pipelined CPU with five pipeline stages, IF, RR, EX, MEM, RW where the pipeline stages execute the corresponding operations mentioned above.

- What is the cycle time of the pipelined CPU? What is the maximum working frequency?
- What is the latency of executing the load-word instruction (lw) in the pipelined CPU?
- What is the latency of executing the add instruction (add) in the pipelined CPU?
- What is the maximum throughput of the pipelined CPU?

(i) The cycle time of the pipelined CPU is determined by the longest delay of the pipeline stages. In this case, it is 200 ps. The maximum working frequency is the reciprocal of the cycle time, which is $1/(200 \text{ ps}) = 5 \text{ GHz}$.

(ii) The latency of executing the load-word instruction (lw) in the pipelined CPU is the time it takes for the instruction to go through all five pipeline stages. In this case, it is $200 \text{ ps} + 100 \text{ ps} + 200 \text{ ps} + 200 \text{ ps} + 100 \text{ ps} = 800 \text{ ps}$.

(iii) The latency of executing the add instruction (add) in the pipelined CPU is also the time it takes for the instruction to go through all five pipeline stages. In this case, it is $200 \text{ ps} + 100 \text{ ps} + 200 \text{ ps} + 0 \text{ ps} + 100 \text{ ps} = 600 \text{ ps}$.

(iv) The maximum throughput of the pipelined CPU is one instruction per cycle time, which is $1/(200 \text{ ps}) = 5 \text{ billion instructions per second}$.

A virtual memory system often implements a TLB to speed up the virtual-to-physical address translation. A TLB has the following characteristics. Assume each TLB entry has a valid bit, a dirty bit, a tag, and the page number. Determine the exact total number of bits to implement this TLB.

- It is direct-mapped.
- It has 16 entries.
- The page size is 4 Kbytes.
- The virtual address space is 4 Gbytes.
- The physical memory is 1 Gbytes.

Answer:

- The TLB is direct-mapped and has 16 entries.
- The page size is 4 Kbytes = 2^{12} bytes.
- The virtual address space is 4 Gbytes = 2^{32} bytes.
- The physical memory is 1 Gbyte = 2^{30} bytes.

1. Tag:

- The tag is used to determine which virtual page number is currently stored in a TLB entry. Given a virtual address space of 2^{32} bytes and a page size of 2^{12} bytes, we have a total of $(2^{32} / 2^{12}) = 2^{20}$ pages.
- Since the TLB is direct-mapped and has 16 entries, each TLB entry can potentially map to one of $(2^{20} / 2^4) = 2^{16}$ pages.
- The number of bits required for the tag is $\log_2(2^{16}) = 16$ bits per entry

2. Page Number:

- The page number refers to the physical page number, given a physical memory size of 2^{30} and a page size of 2^{12} , we have $(2^{30} / 2^{12}) = 2^{18}$ physical pages
- Therefore, the number of bits required for the PPN is $\log_2(2^{18}) = 18$ bits per entry

3. Valid and Dirty bit: Since each TLB entry has a valid bit and a dirty bit, we need 2 bits

4. We need $(16 + 18 + 2) * 16(\text{entries}) = 576$ bits to implement this TLB

Total number of bits:
32 bits (valid and dirty bits) + 256 bits (tag) + 288 bits (Page Number) = 576 bits

Problems

- What is F8EF - 1A0B when these values represent signed 16-bit hex numbers stored in sign-magnitude format?

```
F8EF = 1111 1000 1110 1111 -> negative => -30959_ten
1A0B = 0001 1010 0000 1011 => 6667_ten
F8EF - 1A0B = -30959 - 6667 = -37626
=> 37626 = 1001 0010 1111 1010
=> -37626 = 1 1001 0010 1111 1010 (overflow)
Ans: 1001 0010 1111 1010 = 92FA
```

- Assume 168 and 78 are signed 8-bit decimal integers stored in two's complement format. Calculate 168+78 using saturating arithmetic.

```
168_ten = 1010 1000 =>(minus 1) 1010 0111 => (Reverse) 0101 1000 = 88 => -88
-88 + 78 = -10
```

