

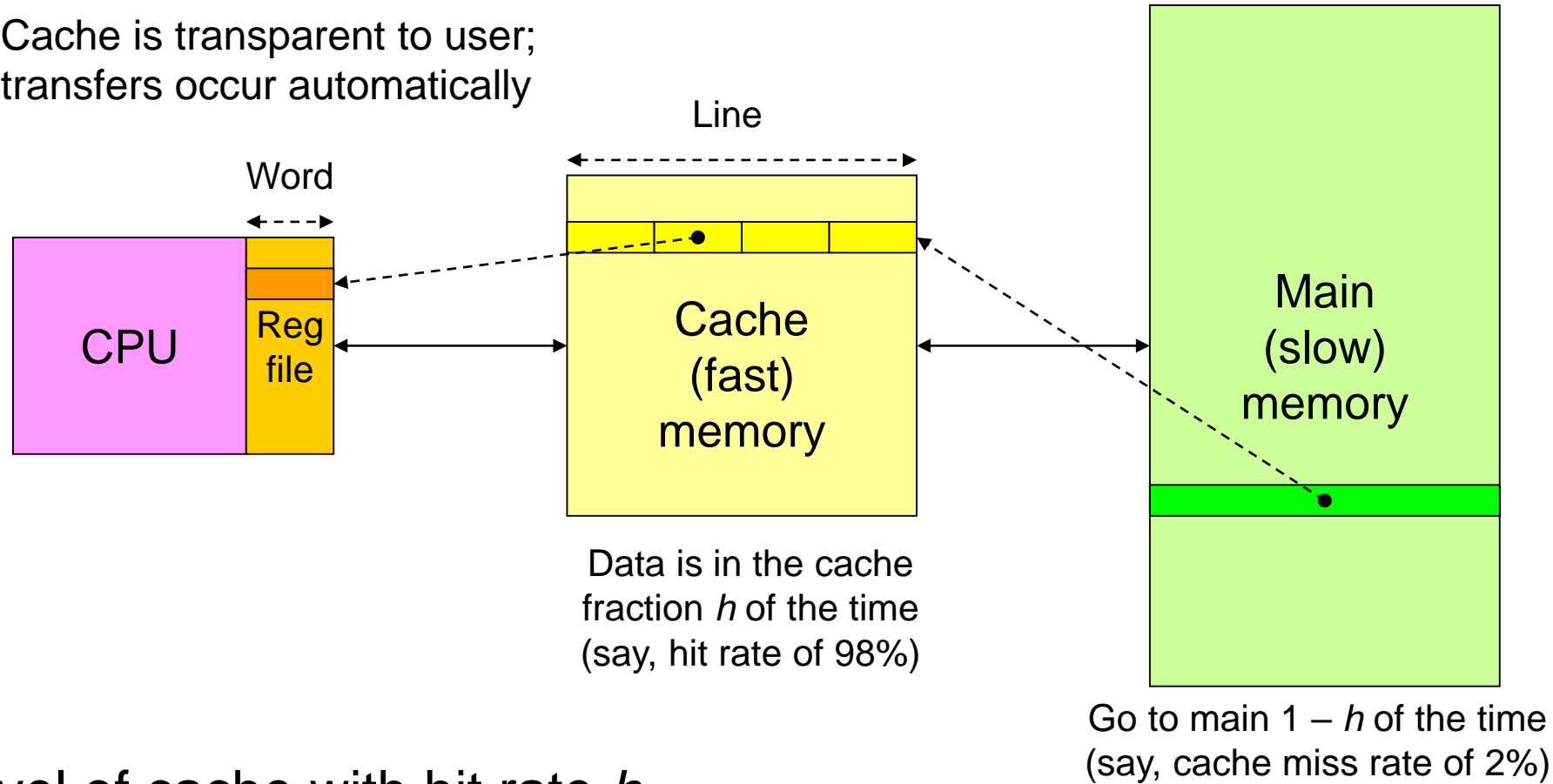
Computer Organization and Architecture

Tutorial on Memory Module

14.09.2023

Cache, Hit/Miss Rate, and Effective Access Time

Cache is transparent to user;
transfers occur automatically



One level of cache with hit rate h

$$C_{\text{eff}} = hC_{\text{fast}} + (1 - h)(C_{\text{slow}} + C_{\text{fast}}) = C_{\text{fast}} + (1 - h)C_{\text{slow}}$$

Example 1: Performance of Hierarchical cache

A system with L1 and L2 caches has a CPI of 1.2 with no cache miss. There are 1.1 memory accesses on average per instruction.

What is the effective CPI with cache misses factored in?

What are the effective hit rate and miss penalty overall if L1 and L2 caches are modeled as a single cache?

Level	Local hit rate	Miss penalty
L1	95 %	8 cycles
L2	80 %	60 cycles

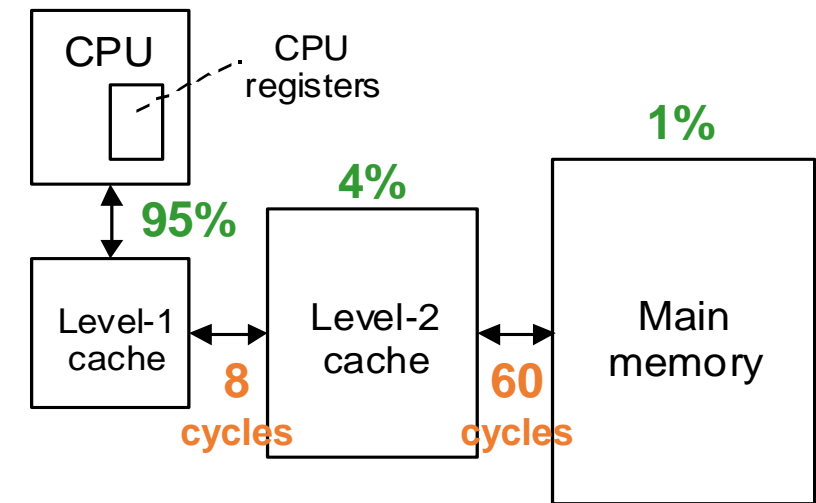
Solution

$$C_{\text{eff}} = C_{\text{fast}} + (1 - h_1)[C_{\text{medium}} + (1 - h_2)C_{\text{slow}}]$$

Because C_{fast} is included in the CPI of 1.2, we must account for the rest

$$\text{CPI} = 1.2 + 1.1(1 - 0.95)[8 + (1 - 0.8)60] = 1.2 + 1.1 \times 0.05 \times 20 = 2.3$$

Overall: hit rate 99% (95% + 80% of 5%), miss penalty 60 cycles



Example 2:

- Consider cache memory having hit ratios for read and write are 60% and 80%. Access Time of cache memory = 40ns. Main memory access time is 400ns. When there is a miss in cache, 4-word block is copied from main memory to cache. CPU generates 60% read and 40 % write requests. What is the bandwidth using write through cache?

- Write Through Cache

Avg. Time to read

$$\begin{aligned}T_{\text{avg_read}} &= T_{\text{cache}} + (1 - \text{hit_read}) T_{\text{mm}} \\&= 40 + 0.4 * 400 \\&= 200 \text{ ns}\end{aligned}$$

$$T_{\text{avg_write}} = 400/4 = 100 \text{ ns}$$

$$\begin{aligned}T_{\text{avg_write_through}} &= \% \text{ read_req} * T_{\text{avg_read}} + \% \text{ write_req} * T_{\text{avg_write}} \\&= 0.6 * 200 + 0.4 * 100 \\&= 160 \text{ ns}\end{aligned}$$

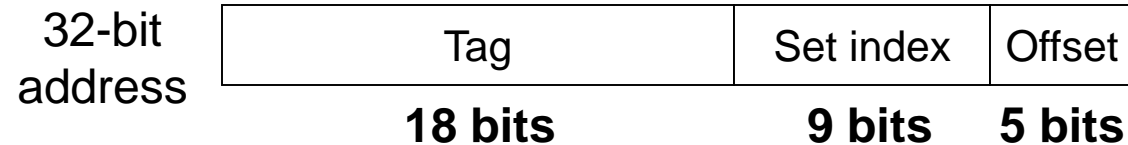
Example 3:

- A 64 KB four-way set-associative cache is byte-addressable and contains 32 B lines. Memory addresses are 32 b wide.
 - a. How wide are the tags in this cache?
 - b. Which main memory addresses are mapped to set number 5?

Solution

The number of sets in the cache = $64\text{KB}/(4 \times 32\text{B}) = 512$.

- a. Address (32 b) = 5 b byte offset + 9 b set index + 18 b tag
- b. Addresses that have their 9-bit set index equal to 5. These are of the general form $2^{14}a + 2^5 \times 5 + b$; e.g., 160-191, 16 554-16 575, . . .



$$\begin{aligned}\text{Tag width} &= \\ 32 - 5 - 9 &= 18\end{aligned}$$

$$\begin{aligned}\text{Set size} &= 4 \times 32 \text{ B} = 128 \text{ B} \\ \text{Number of sets} &= 2^{16}/2^7 = 2^9\end{aligned}$$

$$\begin{aligned}\text{Line width} &= \\ 32 \text{ B} &= 2^5 \text{ B}\end{aligned}$$

Example 4

- Consider a CPU with average CPI of 1.1.
 - Assume an instruction mix: ALU – 50%, LOAD – 15%, STORE – 15%, BRANCH – 20%
 - Assume a cache miss rate of 1.5%, and miss penalty of 50 cycles ($= t_{MM}$).
 - Calculate the effective CPI for a unified L1 cache, using *write through and no write allocate*, with:
 - a) No write buffer
 - b) Perfect write buffer
 - c) Realistic write buffer that eliminates 85% of write stalls.

Number of memory accesses per instruction = $1 + 0.15 + 0.15 = 1.3$

% Reads = $(1 + 0.15) / 1.3 = 88.5\%$ % Writes = $0.15 / 1.3 = 11.5\%$

- **Solution:**

a) With no write buffer (i.e. *stall on all writes*)

- Memory stalls / instr. = $1.3 \times 50 \times (88.5\% \times 1.5\% + 11.5\%) = 8.33$ cycles
- $CPI = CPI_{avg} + \text{Memory stalls / instr.} = 1.1 + 8.33 = 9.43$

b) With perfect write buffer (i.e. *all write stalls are eliminated*)

- Memory stalls / instr. = $1.3 \times 50 \times (88.5\% \times 1.5\%) = 0.86$ cycles
- $CPI = 1.1 + 0.86 = 1.96$

c) With realistic write buffer (*85% of write stalls are eliminated*)

- Memory stalls / instr. = $1.3 \times 50 \times (88.5\% \times 1.5\% + 15\% \times 11.5\%) = 1.98$ cycles
- $CPI = 1.1 + 1.98 = 3.08$

Example 5

- Consider a CPU with average CPI of 1.1.
 - Assume the instruction mix: ALU – 50%, LOAD – 15%, STORE – 15%, BRANCH – 20%
 - Assume a cache miss rate of 1.5%, and miss penalty of 50 cycles ($= t_{MM}$).
 - Calculate the effective CPI for a unified L1 cache, using *write back and write allocate*, with the probability of a cache block being dirty is 10%.

Number of memory accesses per instruction = $1 + 0.15 + 0.15 = 1.3$

- **Solution:**

- Memory accesses per instruction = 1.3
- Stalls / access for writes = $(1 - H_{L1}) \cdot (t_{MM} \times \% \text{ clean} + 2t_{MM} \times \% \text{ dirty})$
= $1.5\% \times (50 \times 90\% + 100 \times 10\%) = 0.825 \text{ cycles}$
- Memory stalls / instr. for write = $1.3 \times 0.825 = 1.07 \text{ cycles}$
- Thus, effective CPI = $1.1 + 1.07 = 2.17$