

Lecture 13 Cache and Memory Review Exercises

Q1:

Consider two computers A and B:

Parameters	A	B
Base CPI	1.0	2.0
Percentage of load/store instructions	50%	50%
L1 caches hit time	1 cycle	2 cycles
L1 instruction cache miss rate	0%	0%
L1 data cache miss rate	4%	5%
L2 cache hit time	15 cycles	10 cycles
L2 cache local miss rate	75%	80%
Main memory access time	250 cycles	100 cycles

What is the CPI for each computer? What is the Average Memory Access Time (AMAT) for each computer?

Solution:

Given L1 instr and data caches hit time for A and B => all L1 caches hit time fits within the base CPI

So $CPI = CPI_{base} + (\%cache_access * \%miss * Cycles_miss)$, where $Cycles_miss = Cycles_hit + \%miss * Cycles_miss$ at the next level

$$CPI \text{ of A} = 1.0 + (50\% * 4\% * (15 + 75\% * 250)) = 5.05$$

$$CPI \text{ of B} = 2.0 + (50\% * 5\% * (10 + 80\% * 100)) = 4.25$$

$$AMAT = t_{avg}(L1) = CPI_base + (\%cache_access * \%miss * Cycles_miss)$$

$$AMAT \text{ of A} = 1.0 + (50\% * 4\% * (15 + 75\% * 250)) = 5.05 \text{ cycles}$$

$$AMAT \text{ of B} = 2.0 + (50\% * 5\% * (10 + 80\% * 100)) = 4.25 \text{ cycles}$$

Q2:

Given a system with:

- 2 memory channels
- 2 DRAM DIMMS (1 DIMMS per channel)

Each DIMM has:

- 1rank
- 8 chips per rank
- 8-bit column size
- 4 banks per chip
- 32,768 rows per bank
- 2,048 columns per bank
- 8-byte bus

Assume a minimum number of bits needed to cover the physical address space are used for physical addresses. Physical addresses are assigned to use the row interleaving scheme. Also assume that the upper bit(s) of the physical address are used to select the channel. Please determine the address mapping of the physical address, i.e., which portions of the physical address bits are used for: column, channel, bank, row, byte in bus offset and in what order?

Solution:

Memory per chip = banks per chip * rows per bank * columns per bank * bits per column = $4 * 32768 * 2048 * 8 = 2\text{Gbits}$

Memory per DIMM = memory per chip * ranks per DIMM * chips per rank = $2\text{Gbit} * 1 * 8 = 16\text{Gbit}$

Total physical memory = memory per DIMM * DIMMs = $16\text{Gbit} * 2 = 32\text{Gbits} = 4\text{GBytes}$

Minimum of physical address bits needed = $\log_2(4\text{G}) = \log_2(4 * 1024 * 1024 * 1024) = 32 \text{ bits}$

Row interleaving is used, So the order of the IDs is as follows:

| Channel ID | Row ID | Bank ID | Column ID | Byte in Bus Offset |

2 channels → Channel ID is 1 bit

4 banks per chip → 4 banks per rank, because all the chips in a rank has the same number of banks → Bank ID is 2 bits

32,768 rows per bank → Row ID is $\log_2(32768) = 15 \text{ bits}$

2048 columns per bank → Column ID is $\log_2(2048) = 11 \text{ bits}$

8B bus → Byte in Bus offset is $\log_2(8) = 3 \text{ bits}$

So, address mapping is as follows:

| Channel (31) | Row (30:16) | Bank (15:14) | Column (13:3) | Byte in Bus Offset (2:0) |

Q3 (Challenging problem, not required):

Consider two computers A and B:

Parameters	A	B
ISA	MIPS	x86

Clock rate	2 GHz	3 GHz
Base CPI	1	2
Number of pipeline stages	5	10
Percentage of branch instructions	10%	10%
Branch miss predictions	10%	5%
Branch misprediction penalty	1 cycle	3 cycles
Percentage of load/store instructions	30%	30%
L1 instruction cache hit time	1 cycle	1 cycle
L1 instruction cache miss rate	2%	2%
L1 data cache hit time	1 cycle	2 cycles
L1 data cache miss rate	8%	5%
L2 cache hit time	15 cycles	12 cycles
L2 cache global miss rate	2% for instruction fetch, 3% for data access	1% for instruction fetch, 4% for data access
Main memory access time	125 ns	100 ns

What is the CPI for each computer? What is the Average Memory Access Time (AMAT) for each computer? Which computer is faster and by how much, on average, if programs execute 1.25 times as many MIPS instructions as x86 instructions (hint: please compare based on **execution time**)?

Solution:

Given the clock rates of A and B,

Main memory access time for A = 125 ns = 250 cycles

Main memory access time for B = 100 ns = 300 cycles

Programs execute 1.25 times as many MIPS instructions as x86 instructions =>

Assume instr count of B is N, then instr count of A is 1.25N

Given L1 instr and data caches hit time for A and B => all L1 caches hit time fits within the base CPI

$$\begin{aligned}
 \text{So CPI} &= \text{CPI_base} + (\% \text{cache_access} * \% \text{miss} * \text{Cycles_miss}) + \\
 &\quad (\% \text{branch} * \% \text{miss_prediction} * \text{miss_penalty}) \\
 &= \text{AMAT} + (\% \text{branch} * \% \text{miss_prediction} * \text{miss_penalty})
 \end{aligned}$$

For computer A:

First, let's calculate AMAT_A:

$$\begin{aligned}\text{AMAT} &= \text{CPI_base} + (\% \text{cache_access} * \% \text{miss} * \text{Cycles_miss}) \\ &= \text{CPI_base} + (\% \text{instr_cache_access} * \% \text{instr_miss} * \text{Cycles_instr_miss}) + \\ &\quad (\% \text{data_cache_access} * \% \text{data_miss} * \text{Cycles_data_miss})\end{aligned}$$

Note that miss rates in the above equation (i.e., %instr_miss and %data_miss) are “local” miss rates, where

Local miss rate of L2 = number of misses in L2 / number of accesses to L2

Global miss rate of L2 = number of misses in L2 / total number of load and stores
= number of misses in L2 / number of accesses to L1

So,

L2 instr local miss rate for A = 2%/2% = 100%

L2 data local miss rate for A = 3%/8% = 37.5%

$$\text{AMAT_A} = 1 + 100\% * 2\% * (15 + 100\% * 250) + 30\% * 8\% * (15 + 37.5\% * 250) = 8.91 \text{ cycles}$$

Calculate CPI_A:

$$\begin{aligned}\text{CPI_A} &= \text{AMAT_A} + (\% \text{branch} * \% \text{miss_prediction} * \text{miss_penalty}) \\ &= 8.91 + (10\% * 10\% * 1) = 8.92\end{aligned}$$

$$\text{Execution time for A} = \text{CPI_A} * \text{Instr_A} * \text{Cycle_time_A} = 8.92 * 1.25N * 0.5\text{ns} = 5.575N \text{ ns}$$

For computer B:

L2 instr local miss rate for B = 1%/2% = 50%

L2 data local miss rate for B = 4%/5% = 80%

$$\text{AMAT_B} = 2 + 100\% * 2\% * (12 + 50\% * 300) + 30\% * 5\% * (12 + 80\% * 300) = 9.02 \text{ cycles}$$

$$\text{CPI_B} = \text{AMAT_B} + (10\% * 5\% * 3) = 9.02 + (10\% * 5\% * 3) = 9.035$$

$$\text{Execution time for B} = \text{CPI_B} * \text{Instr_B} * \text{Cycle_time_B} = 9.035 * N * 0.33\text{ns} = 2.98N \text{ ns}$$

Therefore, Speedup of B over A = 5.575N / 2.98N = 1.87 => B is faster