# ECE 201 Spring 2022 - Homework 3

Assigned by Prof. Wujie Wen
**Due Apr. 29th, 2022 23:59pm**
*Please type/write your answers on a single document and submit through coursesite*

**1. (30 pts) Cache Average Memory Access Time (AMAT) Calculation.**

**5.9** Cache block size (B) can affect both miss rate and miss latency. Assuming a machine with a base CPI of 1, and an average of 1.35 references (both instruction and data) per instruction, find the block size that minimizes the total miss latency given the following miss rates for various block sizes.

| 8: 4% | 16: 3% | 32: 2% | 64: 1.5% | 128: 1% |
|-------|--------|--------|----------|---------|

**5.9.1** [10] <§5.3> What is the optimal block size for a miss latency of $20 \times B$ cycles?

**5.9.2** [10] <§5.3> What is the optimal block size for a miss latency of $24 + B$ cycles?

**5.9.3** [10] <§5.3> For constant miss latency, what is the optimal block size?

**2. (40 pts) Set associative cache questions, including 4 sub questions.**

**5.11** This exercise examines the effect of different cache designs, specifically comparing associative caches to the direct-mapped caches from Section 5.4. For these exercises, refer to the sequence of word address shown below.

```
0x03, 0xb4, 0x2b, 0x02, 0xbe, 0x58, 0xbf, 0x0e, 0x1f,
0xb5, 0xbf, 0xba, 0x2e, 0xce
```

**5.11.1** [10] <§5.4> Sketch the organization of a three-way set associative cache with two-word blocks and a total size of 48 words. Your sketch should have a style similar to Figure 5.18, but clearly show the width of the tag and data fields.

**5.11.2** [10] <§5.4> Trace the behavior of the cache from Exercise 5.11.1. Assume a true LRU replacement policy. For each reference, identify

- the binary word address,
- the tag,
- the index,
- the offset
- whether the reference is a hit or a miss, and
- which tags are in each way of the cache after the reference has been handled.

**5.11.5** [5] <§5.4> Sketch the organization of a fully associative cache with two-word blocks and a total size of eight words. Your sketch should have a style similar to Figure 5.18, but clearly show the width of the tag and data fields.

**5.11.6** [10] <§5.4> Trace the behavior of the cache from Exercise 5.11.5. Assume an LRU replacement policy. For each reference, identify

- the binary word address,

- the tag,

- the index,

- the offset

- whether the reference is a hit or a miss, and

- the contents of the cache after each reference has been handled.

## 3. (30 pts) Multi-level Caching Performance (Question 5.12 on textbook)

**5.12** Multilevel caching is an important technique to overcome the limited amount of space that a first-level cache can provide while still maintaining its speed. Consider a processor with the following parameters:

| Base CPI, No Memory Stalls | Processor Speed | Main Memory Access Time | First-Level Cache Miss Rate per Instruction** | Second-Level Cache, Direct-Mapped Speed | Miss Rate with Second-Level Cache, Direct-Mapped | Second-Level Cache, Eight-Way Set Associative Speed | Miss Rate with Second-Level Cache, Eight-Way Set Associative |
|---|---|---|---|---|---|---|---|
| 1.5 | 2 GHz | 100 ns | 7% | 12 cycles | 3.5% | 28 cycles | 1.5% |

**First Level Cache miss rate is per instruction. Assume the total number of L1 cache misses (instruction and data combined) is equal to 7% of the number of instructions.

**5.12.1** [10] <§5.4> Calculate the CPI for the processor in the table using: 1) only a first-level cache, 2) a second-level direct-mapped cache, and 3) a second-level eight-way set associative cache. How do these numbers change if main memory access time doubles? (Give each change as both an absolute CPI and a percent change.) Notice the extent to which an L2 cache can hide the effects of a slow memory.

**5.12.2** [10] <§5.4> It is possible to have an even greater cache hierarchy than two levels? Given the processor above with a second-level, direct-mapped cache, a designer wants to add a third-level cache that takes 50 cycles to access and will have a 13% miss rate. Would this provide better performance? In general, what are the advantages and disadvantages of adding a third-level cache?

**5.12.3** [20] <§5.4> In older processors, such as the Intel Pentium or Alpha 21264, the second level of cache was external (located on a different chip) from the main processor and the first-level cache. While this allowed for large second-level caches, the latency to access the cache was much higher, and the bandwidth was typically lower because the second-level cache ran at a lower frequency. Assume a 512 KiB off-chip second-level cache has a miss rate of 4%. If each additional 512 KiB of cache lowered miss rates by 0.7%, and the cache had a total access time of 50 cycles, how big would the cache have to be to match the performance of the second-level direct-mapped cache listed above?