

## Ve370 Introduction to Computer Organization

## Homework 7

1. Exercise 5.2.1 (5 points)
2. Exercise 5.2.2(5 points)
3. Exercise 5.2.3 (5 points)

**Exercise 5.2**

In this exercise we look at memory locality properties of matrix computation. The following code is written in C, where elements within the same row are stored contiguously.

<b>a.</b>	<pre>for (I=0; I&lt;8; I++)   for (J=0; J&lt;8000; J++)     A[I][J]=B[I][0]+A[J][I];</pre>
<b>b.</b>	<pre>for (J=0; J&lt;8000; J++)   for (I=0; I&lt;8; I++)     A[I][J]=B[I][0]+A[J][I];</pre>

**5.2.1** [5] <5.1> How many 32-bit integers can be stored in a 16-byte cache line?

**5.2.2** [5] <5.1> References to which variables exhibit temporal locality?

**5.2.3** [5] <5.1> References to which variables exhibit spatial locality?

4. Exercise 5.3.1 (5 points)
5. Exercise 5.3.2 (5 points)
6. Exercise 5.3.3 (10 points)
7. Exercise 5.3.4 (10 points)

## Exercise 5.3

Caches are important to providing a high-performance memory hierarchy to processors. Below is a list of 32-bit memory address references, given as word addresses.

a.	3, 180, 43, 2, 191, 88, 190, 14, 181, 44, 186, 253
b.	21, 166, 201, 143, 61, 166, 62, 133, 111, 143, 144, 61

**5.3.1** [10] <5.2> For each of these references, identify the binary address, the tag, and the index given a direct-mapped cache with 16 one-word blocks. Also list if each reference is a hit or a miss, assuming the cache is initially empty.

**5.3.2** [10] <5.2> For each of these references, identify the binary address, the tag, and the index given a direct-mapped cache with two-word blocks and a total size of 8 blocks. Also list if each reference is a hit or a miss, assuming the cache is initially empty.

**5.3.3** [20] <5.2, 5.3> You are asked to optimize a cache design for the given references. There are three direct-mapped cache designs possible, all with a total of 8 words of data: C1 has 1-word blocks, C2 has 2-word blocks, and C3 has 4-word blocks. In terms of miss rate, which cache design is the best? If the miss stall time is 25 cycles, and C1 has an access time of 2 cycles, C2 takes 3 cycles, and C3 takes 5 cycles, which is the best cache design?

There are many different design parameters that are important to a cache's overall performance. The table below lists parameters for different direct-mapped cache designs.

	Cache Data Size	Cache Block Size	Cache Access Time
a.	32 KB	2 words	1 cycle
b.	32 KB	4 words	2 cycle

**5.3.4** [15] <5.2> Calculate the total number of bits required for the cache listed in the table, assuming a 32-bit address. Given that total size, find the total size

of the closest direct-mapped cache with 16-word blocks of equal size or greater. Explain why the second cache, despite its larger data size, might provide slower performance than the first cache.

8. Exercise 5.4.1 (5 points)
9. Exercise 5.4.2 (5 points)

10. Exercise 5.4.3 (5 points)
11. Exercise 5.4.4 (5 points)
12. Exercise 5.4.5 (10 points)
13. Exercise 5.4.6 (10 points)

### Exercise 5.4

For a direct-mapped cache design with a 32-bit address, the following bits of the address are used to access the cache.

	Tag	Index	Offset
a.	31-10	9-5	4-0
b.	31-12	11-6	5-0

**5.4.1** [5] <5.2> What is the cache line size (in words)?

**5.4.2** [5] <5.2> How many entries does the cache have?

**5.4.3** [5] <5.2> What is the ratio between total bits required for such a cache implementation over the data storage bits?

Starting from power on, the following byte-addressed cache references are recorded.

Address	0	4	16	132	232	160	1024	30	140	3100	180	2180
---------	---	---	----	-----	-----	-----	------	----	-----	------	-----	------

**5.4.4** [10] <5.2> How many blocks are replaced?

**5.4.5** [10] <5.2> What is the hit ratio?

**5.4.6** [20] <5.2> List the final state of the cache, with each valid entry represented as a record of <index, tag, data>.

14. Exercise 5.6.1 (5 points)
15. Exercise 5.6.2 (10 points)

### Exercise 5.6

Media applications that play audio or video files are part of a class of workloads called “streaming” workloads; i.e., they bring in large amounts of data but do not reuse much of it. Consider a video streaming workload that accesses a 512 KB working set sequentially with the following address stream:

0, 2, 4, 6, 8, 10, 12, 14, 16, ...

**5.6.1** [5] <5.5, 5.3> Assume a 64 KB direct-mapped cache with a 32-byte line. What is the miss rate for the address stream above? How is this miss rate sensitive to the size of the cache or the working set? How would you categorize the misses this workload is experiencing, based on the 3C model?

**5.6.2** [5] <5.5, 5.1> Re-compute the miss rate when the cache line size is 16 bytes, 64 bytes, and 128 bytes. What kind of locality is this workload exploiting?