

## HW5: 5.1, 5.2, 5.3, 5.4, 5.5, 5.6, 5.7

more easily  
different  
be much  
he is not  
ows us to  
at may be

memory  
ilities is a  
explored.  
al locality.  
ys as the  
e, such as  
improved

ought into  
hardware  
notice.  
memory  
ter 6 use  
ory on a  
imization

**5.1** 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. Assume each word is a 32-bit integer.

```
for (I=0; I<8; I++)
  for (J=0; J<8000; J++)
    A[I][J]=B[I][0]+A[J][I];
```

**5.1.1** [5] <§5.1> How many 32-bit integers can be stored in a 16-byte cache block?

**5.1.2** [5] <§5.1> References to which variables exhibit temporal locality?

**5.1.3** [5] <§5.1> References to which variables exhibit spatial locality?

Locality is affected by both the reference order and data layout. The same computation can also be written below in Matlab, which differs from C by storing matrix elements within the same column contiguously in memory.

```
for I=1:8
  for J=1:8000
    A(I,J)=B(I,0)+A(J,I);
  end
end
```

472

Chapter 5 Large and Fast: Exploiting Memory Hierarchy

**5.1.4** [10] <§5.1> How many 16-byte cache blocks are needed to store all 32-bit matrix elements being referenced?

**5.1.5** [5] <§5.1> References to which variables exhibit temporal locality?

**5.1.6** [5] <§5.1> References to which variables exhibit spatial locality?

**5.2** 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.

```
3, 180, 43, 2, 191, 88, 190, 14, 181, 44, 186, 253
```

**5.2.1** [10] <§5.3> 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.

**5.2.2** [10] <§5.3> 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.

**5.2.3** [20] <§§5.3, 5.4> You are asked to optimize a cache design for the given references. There are three direct-mapped cache designs...

index  
possibl  
change

**5.3** F  
adres

31

**5.3.1**

**5.3.2**

**5.3.3**  
imple  
Starti

**5.1.4** [10] <§5.1> How many 16-byte cache blocks are needed to store all 32-bit matrix elements being referenced?

**5.1.5** [5] <§5.1> References to which variables exhibit temporal locality?

**5.1.6** [5] <§5.1> References to which variables exhibit spatial locality?

**5.2** 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.

3, 180, 43, 2, 191, 88, 190, 14, 181, 44, 186, 253

**5.2.1** [10] <§5.3> 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.

**5.2.2** [10] <§5.3> 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.

**5.2.3** [20] <§§5.3, 5.4> 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. Below are listed parameters for different direct-mapped cache designs.

**Cache Data Size:** 32 KiB

**Cache Block Size:** 2 words

**Cache Access Time:** 1 cycle

**5.2.4** [15] <§5.3> Calculate the total number of bits required for the cache listed above, 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.

**5.2.5** [20] <§§5.3, 5.4> Generate a series of read requests that have a lower miss rate on a 2 KiB 2-way set associative cache than the cache listed above. Identify one possible solution that would make the cache listed have an equal or lower miss rate than the 2 KiB cache. Discuss the advantages and disadvantages of such a solution.

**5.2.6** [15] <§5.3> The formula shown in Section 5.3 shows the typical method to index a direct-mapped cache, specifically (Block address) modulo (Number of blocks in the cache). Assuming a 32-bit address and 1024 blocks in the cache, consider a different

indexing function, specifically (Block address[31:27] XOR Block address[26:22]). Is it possible to use this to index a direct-mapped cache? If so, explain why and discuss any changes that might need to be made to the cache. If it is not possible, explain why.

**5.3** 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 |
|-------|-------|--------|
| 31-10 | 9-5   | 4-0    |

**5.3.1** [5] <§5.3> What is the cache block size (in words)?

**5.3.2** [5] <§5.3> How many entries does the cache have?

**5.3.3** [5] <§5.3> 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 |

**5.3.4** [10] <§5.3> How many blocks are replaced?

**5.3.5** [10] <§5.3> What is the hit ratio?

**5.3.6** [20] <§5.3> List the final state of the cache, with each valid entry represented a record of <index, tag, data>.

**5.4** Recall that we have two write policies and write allocation policies, and the combinations can be implemented either in L1 or L2 cache. Assume the following choices for L1 and L2 caches:

| L1                                | L2                         |
|-----------------------------------|----------------------------|
| Write through, non-write allocate | Write back, write allocate |

**5.4.1** [5] <§5.3, 5.8> Buffers are employed between different levels of memory hierarchy to reduce access latency. For this given configuration, list the possible buffers needed between L1 and L2 caches, as well as L2 cache and memory.

Describe the procedure of handling an L1 write-r

cycles, which is overall signs.

5.3.5 [10] <§5.3> What is the hit ratio?

5.3.6 [20] <§5.3> List the final state of the cache, with each valid entry represented as a record of <index, tag, data>.

5.4 Recall that we have two write policies and write allocation policies, and their combinations can be implemented either in L1 or L2 cache. Assume the following choices for L1 and L2 caches:

| L1                                | L2                         |
|-----------------------------------|----------------------------|
| Write through, non-write allocate | Write back, write allocate |

5.4.1 [5] <§§5.3, 5.8> Buffers are employed between different levels of memory hierarchy to reduce access latency. For this given configuration, list the possible buffers needed between L1 and L2 caches, as well as L2 cache and memory.

5.4.2 [20] <§§5.3, 5.8> Describe the procedure of handling an L1 write-miss, considering the component involved and the possibility of replacing a dirty block.

5.4.3 [20] <§§5.3, 5.8> For a multilevel exclusive cache (a block can only reside in one of the L1 and L2 caches), configuration, describe the procedure of handling an L1 write-miss, considering the component involved and the possibility of replacing a dirty block.

listed  
osest  
the  
the  
rate  
ible  
e 2  
to  
in  
nt

Chapter 5 Large and Fast: Exploiting Memory Hierarchy

Consider the following program and cache behaviors.

| Data Reads per 1000 Instructions | Data Writes per 1000 Instructions | Instruction Cache Miss Rate | Data Cache Miss Rate | Block Size (byte) |
|----------------------------------|-----------------------------------|-----------------------------|----------------------|-------------------|
| 250                              | 100                               | 0.30%                       | 2%                   | 64                |

5.4.4 [5] <§§5.3, 5.8> For a write-through, write-allocate cache, what are the minimum read and write bandwidths (measured by byte per cycle) needed to achieve a CPI of 2?

5.4.5 [5] <§§5.3, 5.8> For a write-back, write-allocate cache, assuming 30% of replaced data cache blocks are dirty, what are the minimal read and write bandwidths needed for a CPI of 2?

5.4.6 [5] <§§5.3, 5.8> What are the minimal bandwidths needed to achieve the performance of CPI=1.5?

5.5 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 KiB working set

5.6  
per  
ma  
The  
P2.

P  
P

5  
P

5

5  
t  
P

**5.5** Media applications that are 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 KiB working set sequentially with the following address stream:

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

**5.5.1** [5] <§5.4, 5.8> Assume a 64 KiB direct-mapped cache with a 32-byte block. 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?

**5.5.2** [5] <§5.1, 5.8> Re-compute the miss rate when the cache block size is 16 bytes, 64 bytes, and 128 bytes. What kind of locality is this workload exploiting?

**5.5.3** [10] <§5.13> “Prefetching” is a technique that leverages predictable address patterns to speculatively bring in additional cache blocks when a particular cache block is accessed. One example of prefetching is a stream buffer that prefetches sequentially adjacent cache blocks into a separate buffer when a particular cache block is brought in. If the data is found in the prefetch buffer, it is considered as a hit and moved into the cache and the next cache block is prefetched. Assume a two-entry stream buffer and assume that the cache latency is such that a cache block can be loaded before the computation on the previous cache block is completed. What is the miss rate for the address stream above?

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

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

64 bytes, and 128 bytes. What kind of locality is this workload exploiting?

**5.5.3** [10] <§5.13> “Prefetching” is a technique that leverages predictable address patterns to speculatively bring in additional cache blocks when a particular cache block is accessed. One example of prefetching is a stream buffer that prefetches sequentially adjacent cache blocks into a separate buffer when a particular cache block is brought in. If the data is found in the prefetch buffer, it is considered as a hit and moved into the cache and the next cache block is prefetched. Assume a two-entry stream buffer and assume that the cache latency is such that a cache block can be loaded before the computation on the previous cache block is completed. What is the miss rate for the address stream above?

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

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

**5.5.4** [10] <§5.3> What is the optimal block size for a miss latency of  $20 \times B$  cycles?

**5.5.5** [10] <§5.3> What is the optimal block size for a miss latency of  $24 + B$  cycles?

**5.5.6** [10] <§5.3> For constant miss latency, what is the optimal block size?

For the next three exercises, presumably making use of the AMAT and hit times from the previous exercises, the miss rates indicated is its local miss rate.

L2 Size

1 MiB

**5.6.4** [10] <§5.4> Assume a 1-MiB L2 cache with an AMAT better than the L1 cache.

**5.6.5** [5] <§5.4> Assume a 1-MiB L2 cache with a total CPI for the L2 cache of 1.5.

**5.6.6** [10] <§5.4> Assume a 1-MiB L2 cache with a total CPI for the L2 cache of 1.5. For each of the two cache configurations, what is the faster, what is the slower?

**5.7** This exercise is for comparing and contrasting the two exercises, re-compute the miss rate.

**5.7.1** [10] Assume a 1-MiB L2 cache with a cache content size of 24 words and 24 tag bits, the miss rate is 1.5%.

**5.6.5** [5] Assume a 1-MiB L2 cache with a total CPI for the L2 cache of 1.5.

**5.6.6** [10] Assume a 1-MiB L2 cache with a total CPI for the L2 cache of 1.5. For each of the two cache configurations, what is the faster, what is the slower?

**5.7** This exercise is for comparing and contrasting the two exercises, re-compute the miss rate.

**5.7.1** [10] Assume a 1-MiB L2 cache with a cache content size of 24 words and 24 tag bits, the miss rate is 1.5%.

**5.7.2** [10] Assume a 1-MiB L2 cache with a cache content size of 24 words and 24 tag bits, the miss rate is 1.5%. Use LRU as the replacement policy. Is a hit or a miss?

| Block Size<br>(byte) |
|----------------------|
| 64                   |

**5.6** In this exercise, we will look at the different ways capacity affects overall performance. In general, cache access time is proportional to capacity. Assume that main memory accesses take 70 ns and that memory accesses are 36% of all instructions. The following table shows data for L1 caches attached to each of two processors, P1 and P2.

|    | L1 Size | L1 Miss Rate | L1 Hit Time |
|----|---------|--------------|-------------|
| P1 | 2 KiB   | 8.0%         | 0.66 ns     |
| P2 | 4 KiB   | 6.0%         | 0.90 ns     |

**5.6.1** [5] <§5.4> Assuming that the L1 hit time determines the cycle times for P1 and P2, what are their respective clock rates?

**5.6.2** [5] <§5.4> What is the Average Memory Access Time for P1 and P2?

**5.6.3** [5] <§5.4> Assuming a base CPI of 1.0 without any memory stalls, what is the total CPI for P1 and P2? Which processor is faster?

For the next three problems, we will consider the addition of an L2 cache to P1 to presumably make up for its limited L1 cache capacity. Use the L1 cache capacities and hit times from the previous table when solving these problems. The L2 miss rate indicated is its local miss rate.

| L2 Size | L2 Miss Rate | L2 Hit Time |
|---------|--------------|-------------|
| 1 MiB   | 95%          | 5.62 ns     |

| L2 Size | L2 Miss Rate | L2 Hit Time |
|---------|--------------|-------------|
| 1 MiB   | 95%          | 5.62 ns     |

**5.6.4** [10] <§5.4> What is the AMAT for P1 with the addition of an L2 cache? Is the AMAT better or worse with the L2 cache?

**5.6.5** [5] <§5.4> Assuming a base CPI of 1.0 without any memory stalls, what is the total CPI for P1 with the addition of an L2 cache?

**5.6.6** [10] <§5.4> Which processor is faster, now that P1 has an L2 cache? If P1 is faster, what miss rate would P2 need in its L1 cache to match P1's performance? If P2 is faster, what miss rate would P1 need in its L1 cache to match P2's performance?

**5.7** This exercise examines the impact of different cache designs, specifically comparing associative caches to the direct-mapped caches from Section 5.4. For these exercises, refer to the address stream shown in Exercise 5.2.

**5.7.1** [10] <§5.4> Using the sequence of references from Exercise 5.2, show the final cache contents for a three-way set associative cache with two-word blocks and a total size of 24 words. Use LRU replacement. For each reference identify the index bits, the tag bits, the block offset bits, and if it is a hit or a miss.

**5.7.2** [10] <§5.4> Using the references from Exercise 5.2, show the final cache contents for a fully associative cache with one-word blocks and a total size of 8 words. Use LRU replacement. For each reference identify the index bits, the tag bits, and if it is a hit or a miss.

**5.7.3** [15] <§5.4> Using the references from Exercise 5.2, what is the miss rate for a fully associative cache with two-word blocks and a total size of 8 words, using LRU replacement? What is the miss rate using MRU (most recently used) replacement? Finally what is the best possible miss rate for this cache, given any replacement policy?

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 MissRate per Instruction | Second Level Cache, Direct-Mapped Speed | Global Miss Rate with Second Level Cache, Direct-Mapped | Second Level Cache, Eight-Way Set Associative Speed | Global 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%  |

**5.7.4** [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 is doubled? If it is cut in half?

| Base CPI, No Memory Stalls | Processor Speed | Main Memory Access Time | First Level Cache MissRate per Instruction | Second Level Cache, Direct-Mapped Speed | Global Miss Rate with Second Level Cache, Direct-Mapped | Second Level Cache, Eight-Way Set Associative Speed | Global 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%  |

**5.7.4** [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 is doubled? If it is cut in half?

**5.7.5** [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 reduce the global miss rate to 1.3%. Would this provide better performance? In general, what are the advantages and disadvantages of adding a third level cache?

**5.7.6** [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 global miss rate of 4%. If each additional 512 KiB of cache lowered global 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? Of the eight-way set associative cache?

**5.8.1** [5]

**5.8.2** [5]

**5.8.3** [5]  
realistic sit

**5.8.4** [5]  
device is d

**5.9** This  
DED) Har

**5.9.1** [5]  
128-bit wo

**5.9.2** [5]  
employ S

performa

number of

**5.9.3** C  
valu

**5.9.1** [5]  
128-bit word using

**5.9.2** [5] <§5.5>  
employ SEC/DED

performance ratio

number of parity b

can be corrected.

**5.9.3** Consider a  
value 0x375, is the

**5.10** For a high  
size is determin

average a B-tree

its B-tree depth,

entries, and a 10

optimal page size

Page Size (Ki

2

4

8

16

32

64