

Two hours - online

The exam will be taken on line.
This paper version is made available as a backup

**UNIVERSITY OF MANCHESTER
SCHOOL OF COMPUTER SCIENCE**

System Architecture

Date: Monday 21st May 2018

Time: 14:00 - 16:00

Please answer all Questions.

Use a SEPARATE answerbook for each SECTION

© The University of Manchester, 2018

This is a CLOSED book examination

The use of electronic calculators is permitted provided they
are not programmable and do not store text

[PTO]

Section A1. **Caches**

- a) What are the three categories of cache misses? Explain under which circumstances each type of cache miss will occur. For each type of cache miss, identify a cache configuration where they would be the most likely to occur. (5 marks)
- b) Describe the forms of locality exploited by caches. For each type of locality, propose a program optimization that would increase the availability of such locality. (4 marks)
- c) The Intel manual for the Skylake architecture describes the following cache organization.
First level data cache (L1D): 32 kilobytes (KB), 8-way set-associative with 64 bytes (B) cache line size and 4 cycles access latency (time required to access this cache).
Second level cache (L2): 256KB, 8-way set-associative with 64B cache line size and 12 cycles access latency.
We will assume that the RAM access time is 200 cycles.
- i) Considering that you are developing an application for which you have measured a 90% cache hit rate in L1D and also 90% cache hit rate in L2, what is the average memory access time of your processor for that application? Show your reasoning. (3 marks)
- ii) After optimizing your code to better exploit the first level of cache, you are now measuring an improved 95% L1D hit rate and still the same 90% L2 hit rate. What is the new average memory access time for your application ? Explain your reasoning. (2 marks)
- iii) One of your colleagues suggests a different optimization that does not improve the L1D hit rate, which remains at the original 90%, but eliminates all L2 cache misses. Is that better or worse than the previous optimization? Explain why. (2 mark)
- iv) The Intel Skylake can additionally include a large level three (L3) cache with an access time of 44 cycles. Assuming that due to its large size the L3 cache achieves 100% hit rate on your application, would this additional level of cache change which of the two optimizations in questions (c.2) and (c.3) is the most beneficial? Explain and detail your reasoning. (4 marks)

2. Virtualization and Storage

- a) Your company is in the process of renewing the staff desktop computers. The supplier is offering a 1 terabyte (TB) hard disk drive (HDD) as standard on each desktop or an upgrade to 2TB for £40. Alternatively, you have the choice of equipping each desktop with an extra 1TB HDD (so having two 1TB HDDs) also for an extra £40. Considering that there is no room or energy constraint in the desktop enclosure and both upgrades are within budget, which option would you choose? Detail your reasoning. (3 marks)
- b) A high-end enterprise server is equipped with a storage system comprised of 16 identical HDDs, each with a 2TB capacity and a sustained bandwidth of 200 MB/s. The HDDs are organised as a RAID 160.
- i) Considering that the storage system is symmetrically structured (i.e., each level of the system is uniform), what are the possible configurations of this RAID system (i.e., distribution of HDDs across the different RAID levels)? Explain your reasoning. (2 marks)
- ii) What is the overall storage capacity of this storage system? Explain. (1 mark)
- iii) What is the maximum number of HDDs that can fail before any data is lost? And what is the minimum number of HDDs that can always safely fail within this RAID configuration? Explain each answer. (4 marks)
- iv) Considering that all of the disks from the first scenario in the question (b.3) above have failed, how long will the recovery of this RAID system take once all failed HDDs have been replaced with new, empty ones? Consider that the machine is isolated and is not performing any other task. Detail your reasoning. (4 marks)
- c) Define the term “checkpointing and restoring”, also called “snapshot and rollback”, explain how it can be implemented and optimized using System Virtualization. (4 marks)
- d) What are possible uses for checkpointing and restoring? Cite at least two. (2 marks)

Section B3. **Processor architecture**

a) Consider the following set of instructions:

1. LDR R1, X
2. MUL R1, R1, R1
3. LDR R2, Y
4. MUL R2, R2, R2
5. LDR R3, Z
6. MUL R3, R3, R3
7. ADD R1, R1, R2
8. ADD R1, R1, R3

Identify the dependencies between instructions and discuss whether they are affecting the performance of this code.

How long will it take to run this code in a classic 5-stage pipeline? Provide an explanation for your answer.

Calculate the Cycles-per-Instruction (CPI) ratio.

Assume all forms of forwarding are implemented and all Instructions and Data are available in the L1 cache. (8 Marks)

b) Reorder the instructions so that you minimise the execution time. What would be the execution time now? And the CPI? (6 Marks)

c) Discuss the differences between the Scoreboard and the Tomasulo Out-of-order processor architectures. (6 Marks)

4. **Multicore architectures**

- a) Explain the concept of cache coherence in the context of multicore processors and its importance in shared memory systems. (2 Marks)
- b) Discuss the similarities and differences between Snoopy and Directory cache coherence protocols in terms of the variety protocols they support, where and how status information is stored, their communication infrastructure and their scalability. (8 Marks)
- c) We are planning to replace the school server supporting our teaching infrastructure. We contacted two separate providers and each of them provided the best configuration within our budget. The specifications of the servers are as follows:

	Provider #1	Provider #2
# of cores	8 out-of-order cores	36 in-order cores
Multithreading	2 threads - SMT	No Multithreading
Superscalar	4-way	3-way
Clock freq.	2 GHz	1 GHz
RAM	128GB	192GB
Cache Hierarchy	L1: 64KB dedicated L2: 256KB dedicated L3: 32MB Shared	L1: 64KB dedicated L2: 256KB dedicated L3: 24MB Shared

Compute the following metrics for each of the systems: (8 Marks)

- i) single-thread peak performance (instructions per second)
- ii) peak IPC (instructions per cycle) of the system
- iii) peak computing throughput of the system (instructions per second)
- iv) the total number of hardware threads supported by the system
- d) In order to decide between these two alternatives, we have analysed the workloads that run on top of the current system. We found out that the run time is divided roughly evenly into 2 different phases: i) a strictly sequential one and ii) a fully parallel one. Given the disparity of the two alternative systems we modelled these workloads in a full system simulator and obtained per-thread IPCs of 2.7 and 1.5 for phase i) and ii), respectively for Provider #1 solution. Similarly, for Provider #2 system, we got an IPC of 1.4 in both phases. Based on these estimations, discuss which system would be able to provide higher performance. Justify your decision numerically. (2 Marks)

END OF EXAMINATION