

※ 考生請注意：本試題不可使用計算機。請於答案卷(卡)作答，於本試題紙上作答者，不予計分。
共 5 題，請在答案卷作一表格如下，並清楚地填入這些題目的答案，否則不予計分。

題號	答案
1.	(1)
	(2)
	(3)
2.	(1)
	(2)
	(3)
3.	(1)
	(2)
	(3)
	(4)
	(5)
4.	(1)
	(2)
	(3)
	(4)
5.	(1)
	(2)
	(3)
	(4)

請勿在此作答

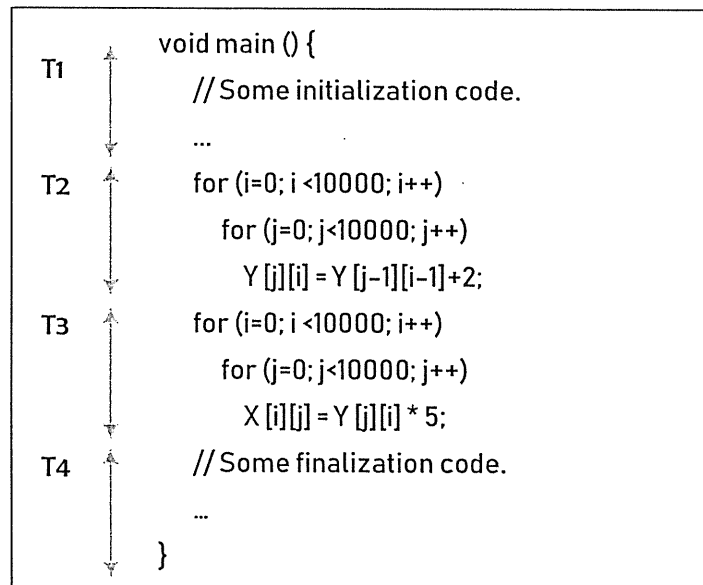
1. [20%] Caches are important to providing a high-performance memory hierarchy to processors. Below is a list of 32-bit memory address references, given as word addresses.

35, 149, 43, 90, 191, 91, 148, 14, 42, 190, 69, 15

- (1) [8%] For each of these references, identify the tag and the index given a direct-mapped L1 cache with two-word blocks and a total size of 16 words. Also list if each reference is a hit or a miss, assuming the cache is initially empty.
- (2) [8%] For each of these references, identify the tag and the index given a four-way set-associative L1 cache with two-word blocks and a total size of 16 words. Also list if each reference is a hit or a miss, assuming the cache is initially empty and LRU replacement is used.
- (3) [4%] Given the above memory address reference, what is the average memory access time (AMAT) if the hit time of the directed-mapped L1 is one cycle and main memory access takes 200 cycles. What is the AMAT of the four-way set-associative cache if the hit time is 2 cycles and main memory access also takes 200 cycles?

2. [15%] Consider a sequential C program as illustrated in the figure below. The program is divided into four code segments, T1, T2, T3, and T4. We profile the program execution time on the single-core MIPS processor and the execution time for the code segments is: 2ms for T1, 8ms for T2, 16ms for T3, and 2ms for T4. Moreover, two different processor variants are available: one is the quad-core MIPS processor (with the same instruction set architecture as the single-core processor), and the other one is the single-core MIPS processor with an 8-wide SIMD engine. Assume we have converted the program into parallel versions suitable for the two processors, respectively. Please answer the following questions.

- (1) [6%] Parallelize the program suitable for the multicore processor and calculate the speedup of the parallelized program.
- (2) [6%] Parallelize the program suitable for the processor with the SIMD engine and calculate the speedup of the parallelized program.
- (3) [3%] Based on the speedups calculated in (1) and (2), please determine which processor is faster? Why?



3. [15%] Determine whether each of the following statements is True (T) or False (F), and explain your answer.

- (1) [3%] Graphics applications deliver higher performance because the DRAM chips on GPU cards help reduce memory latency.
- (2) [3%] A shared memory can be created for computer clusters to exchange data.
- (3) [3%] A GPU is always faster than a CPU at the cost that the GPU has higher power consumption.
- (4) [3%] Given the two-level caches, the purpose of the first-level cache is more about miss rate, and the purpose of the second-level cache is more about hit time.
- (5) [3%] Associativity of the addition operations holds for both integer and floating-point numbers.

4. [20%] Consider a computer with a 64-entry TLB in its main processor. The computer has a hard disk spinning at 5400 RPM and holding 40 sectors per track with a sector size of 512 bytes. In addition, the access latency time for the hard drive is 10 milliseconds, and the seek time is 30 milliseconds.

- (1) [4%] Calculate TLB reach when the computer supports memory page sizes of 1 KB and 4 KB, respectively.
- (2) [4%] Compute the data transfer rate of the hard disk in KB/s.
- (3) [4%] Estimate the I/O times when page faults occur; give your I/O times in milliseconds for handling the page fault of a 1-KB page and a 4-KB page, respectively.
- (4) [8%] We assume that the data transfer time of the hard disk is proportional to the page size. What are the I/O times to transfer the same amount of data (i.e., 4KB data) with 1-KB and 4-KB pages, respectively? In such a case, if we want to minimize the I/O time on the computer, which page size (1 KB or 4KB) is desired?

5. [30%] Considering the real-time systems, which are waiting for events in runtime to occur, when an event occurs, the system must act and respond to the event as fast as possible. On such systems, event latency is used to denote the amount of the elapsed time from the event occurring time to the time where it is serviced.

- (1) [10%] What are the two types of latencies that affect the event latency of such systems?
- (2) [5%] Which process scheduling scheme (i.e., preemption- or non-preemption-based scheme) is preferred to provide less latency?
- (3) [10%] Suppose that the following processes arrive for execution upon at the arrival time. Each process will run the amount of time as listed below. Based on your answer given in (2), choose FCFS or RR (quantum = 20 milliseconds) to compute the average waiting time and the average turnaround time for the processes.

<u>Process</u>	<u>Arrival Time</u>	<u>Burst Time</u>
P ₁	0	53
P ₂	0	17
P ₃	0	68
P ₄	0	24

- (4) [5%] If FCFS is used to schedule the processes given in (3) on one processor, how many possible different schedules are there?