1. (a) (5%) A CPU contains a 32-kilobyte first-level cache, that is addressed by the 32-bit virtual address generated by the CPU. Assume that the block size of the cache is 32 bytes and tag field of the cache controller is 18 bits. What is the degree of set associativity of the cache? Please show how you get your answer.

(b) (10%) After the system is reset, the CPU accesses the memory with the following sequence:

1245C6C4Hx,
228046C5Hx,
228006D8Hx,
228046D9Hx,
1245C6C3Hx.

How many cache misses will occur? Describe the reason of each miss. Assume that one byte is the basic memory unit addressed by the CPU and the cache controller implements a LRU (Least Recently Used) replacement policy.

2. (10%) A CPU executes a register-register instruction, such as
ADD R1, R2,
in a 4-stage pipeline:
   IF: fetch the instruction
   ID: decode the instruction
   ALU: compute the result
   WR: write result to the destination register
The same CPU executes a register-memory instruction, such as
ADD R1, Memory[R2+10],
in a 6-stage pipeline.
   IF: fetch the instruction
   ID: decode the instruction
   ADD: compute effective address
   MEM: access memory
   ALU: compute the result
   WR: write the result to the destination register
In implementing the control of the CPU, does the designer need to consider the write-after-write data hazard? Please explain your reason.

3. (15%) Please describe and discuss Processor-memory bus, Backplane bus and I/O bus and the connection among them.

4. (10%) Please discuss two kinds of branch prediction mechanism.

接背面

5. CPU scheduling is for short-term scheduling of processes. Please answer the following questions:

    a. What are the differences between "non-preemptive scheduling" and "preemptive scheduling"? (10pts)

    b. Given three processes T1, T2, and T3 arriving at time 0, please draw the Gantt chart of their executions under the Round-Robin scheduling with the time slice equal to 2. Note that the CPU burst time of T1, T2, and T3 are 7, 4, and 5, respectively. Initially, T1 is in the front of the ready queue, and T3 is at the end of the queue. (5pts)

6. In the storage hierarchy, we have registers, cache, main memory, electronic disks, magnetic disks, optical disks, and magnetic tapes. Different storage levels have different performance and characteristics. Please answer the following questions:

    A. What is the difference between caching and buffering? (10pts)

    B. Virtual memory is often implemented by demand-paging. Consider an operating system with on-demand paging over a single-level page table, where the page table of each process always resides in the memory. Suppose that the TLB access time, the memory access time, and the average page-fault service time are 20ns, 100ns, and 25ms, respectively. Let the hit ratio of TLB and the page fault rate be 98% and 0.0001%, respectively. What is the effective access time? ($1ns = 10^{-9}$ seconds. $1ms = 10^{-3}$ seconds.) (10pts)

7. Process synchronization should avoid deadlocks and starvation.

    Critical regions are high-level language constructs in the following format:

        **region v when B do S;**

where v, B, and S are a variable accessed inside the region, the condition for entering the region, and the statement inside the region. Please implement the dining philosophers problem using critical regions. You may assume that the problem has five philosophers. There should not be any deadlock. (In the dining philosophers problem, five philosophers share a common circular table, and the table is laid with five chopsticks. A chopstick is between every two philosophers. No philosopher can eat until he picks up his left and right chopsticks simultaneously.) (15pts)