

CHAPTER 1 COMPUTER SYSTEM OVERVIEW

ANSWERS TO QUESTIONS

- 1.1** A processor, which controls the operation of the computer and performs its data processing functions ; a **main memory**, which stores both data and instructions; **I/O modules**, which move data between the computer and its external environment; and the system bus, which provides for communication among processors, main memory, and I/O modules.
- 1.2 User-visible registers:** Enable the machine- or assembly-language programmer to minimize main memory references by optimizing register use. For high-level languages, an optimizing compiler will attempt to make intelligent choices of which variables to assign to registers and which to main memory locations. Some high-level languages, such as C, allow the programmer to suggest to the compiler which variables should be held in registers. **Control and status registers:** Used by the processor to control the operation of the processor and by privileged, operating system routines to control the execution of programs.
- 1.3** These actions fall into four categories: **Processor-memory:** Data may be transferred from processor to memory or from memory to processor. **Processor-I/O:** Data may be transferred to or from a peripheral device by transferring between the processor and an I/O module. **Data processing:** The processor may perform some arithmetic or logic operation on data. **Control:** An instruction may specify that the sequence of execution be altered.
- 1.4** Multiple interrupts may be serviced by assigning different priorities to interrupts arising from different sources. This enables a higher-priority interrupt to be serviced first when multiple requests arrive simultaneously; it also allows a higher-priority interrupt to preempt a lower-priority interrupt. For example, suppose a system has assigned a higher priority to a communication line and a lower priority to a magnetic disk. When two simultaneous requests arrive, the computer services the communication line. Similarly, if some disk operations are ongoing when a request for the communication line arrives, the state of the disk is put in a stack and the communication line operations are catered to.

- 1.5** In interrupt-driven I/O, when data is available in the peripheral, an interrupt facility and special commands inform the interface to issue an interrupt request signal. In the meantime, the CPU can continue its other activities. When the CPU detects an external signal-interrupt, it momentarily stops the task it is processing, services the I/O transfer process, and then resumes the original task.
- 1.6** The characteristics observed while going up the memory hierarchy are **a.** increasing cost per bit, **b.** decreasing capacity, **c.** decreasing access time, and **d.** increasing frequency of access to the memory by the processor.
- 1.7** The main trade-offs for determining the cache size are the speed and the cost of the cache.
- 1.8** A multicore computer is a special case of a multiprocessor, in which all of the processors are on a single chip.
- 1.9** The cache write policies are as follows:
- a. Write through:** Whenever a block in the cache is altered, it is immediately written to the main memory.
 - b. Write back:** The contents of a cache block are written to the main memory only when that block is replaced from the cache.
- 1.10 Spatial locality** is generally exploited by using larger cache blocks and by incorporating prefetching mechanisms (fetching items of anticipated use) into the cache control logic. **Temporal locality** is exploited by keeping recently used instruction and data values in cache memory and by exploiting a cache hierarchy.

ANSWERS TO PROBLEMS

- 1.1** Memory (contents in hex): 300: 3005; 301: 5940; 302: 7006
Step 1: 3005 → IR; **Step 2:** 3 → AC
Step 3: 5940 → IR; **Step 4:** 3 + 2 = 5 → AC
Step 5: 7006 → IR; **Step 6:** AC → Device 6
- 1.2 1. a.** The PC contains 300, the address of the first instruction. This value is loaded in to the MAR.
b. The value in location 300 (which is the instruction with the value 1940 in hexadecimal) is loaded into the MBR, and the PC is incremented. These two steps can be done in parallel.
c. The value in the MBR is loaded into the IR.

2.
 - a. The address portion of the IR (940) is loaded into the MAR.
 - b. The value in location 940 is loaded into the MBR.
 - c. The value in the MBR is loaded into the AC.
3.
 - a. The value in the PC (301) is loaded in to the MAR.
 - b. The value in location 301 (which is the instruction with the value 5941) is loaded into the MBR, and the PC is incremented.
 - c. The value in the MBR is loaded into the IR.
4.
 - a. The address portion of the IR (941) is loaded into the MAR.
 - b. The value in location 941 is loaded into the MBR.
 - c. The old value of the AC and the value of location MBR are added and the result is stored in the AC.
5.
 - a. The value in the PC (302) is loaded in to the MAR.
 - b. The value in location 302 (which is the instruction with the value 2941) is loaded into the MBR, and the PC is incremented.
 - c. The value in the MBR is loaded into the IR.
6.
 - a. The address portion of the IR (941) is loaded into the MAR.
 - b. The value in the AC is loaded into the MBR.
 - c. The value in the MBR is stored in location 941.

1.3 a. $2^{24} = 16$ MBytes

- b. **(1)** If the local address bus is 32 bits, the whole address can be transferred at once and decoded in memory. However, since the data bus is only 16 bits, it will require 2 cycles to fetch a 32-bit instruction or operand.
- (2)** The 16 bits of the address placed on the address bus can't access the whole memory. Thus a more complex memory interface control is needed to latch the first part of the address and then the second part (since the microprocessor will end in two steps). For a 32-bit address, one may assume the first half will decode to access a "row" in memory, while the second half is sent later to access a "column" in memory. In addition to the two-step address operation, the microprocessor will need 2 cycles to fetch the 32 bit instruction/operand.
- c. The program counter must be at least 24 bits. Typically, a 32-bit microprocessor will have a 32-bit external address bus and a 32-bit program counter, unless on-chip segment registers are used that may work with a smaller program counter. If the instruction register is to contain the whole instruction, it will have to be 32-bits long; if it will contain only the op code (called the op code register) then it will have to be 8 bits long.

1.4 In cases **(a)** and **(b)**, the microprocessor will be able to access $2^{16} = 64K$ bytes; the only difference is that with an 8-bit memory each access will transfer a byte, while with a 16-bit memory an access may

transfer a byte or a 16-byte word. For case **(c)**, separate input and output instructions are needed, whose execution will generate separate "I/O signals" (different from the "memory signals" generated with the execution of memory-type instructions); at a minimum, one additional output pin will be required to carry this new signal. For case **(d)**, it can support $2^8 = 256$ input and $2^8 = 256$ output byte ports and the same number of input and output 16-bit ports; in either case, the distinction between an input and an output port is defined by the different signal that the executed input or output instruction generated.

1.5 Clock cycle = $\frac{1}{8 \text{ MHz}} = 125 \text{ ns}$

Bus cycle = $4 \times 125 \text{ ns} = 500 \text{ ns}$

2 bytes transferred every 500 ns; thus transfer rate = 4 MBytes/sec

Doubling the frequency may mean adopting a new chip manufacturing technology (assuming each instructions will have the same number of clock cycles); doubling the external data bus means wider (maybe newer) on-chip data bus drivers/latches and modifications to the bus control logic. In the first case, the speed of the memory chips will also need to double (roughly) not to slow down the microprocessor; in the second case, the "word length" of the memory will have to double to be able to send/receive 32-bit quantities.

1.6 a. Input from the Teletype is stored in INPR. The INPR will only accept data from the Teletype when FGI=0. When data arrives, it is stored in INPR, and FGI is set to 1. The CPU periodically checks FGI. If FGI =1, the CPU transfers the contents of INPR to the AC and sets FGI to 0.

When the CPU has data to send to the Teletype, it checks FGO. If FGO = 0, the CPU must wait. If FGO = 1, the CPU transfers the contents of the AC to OUTF and sets FGO to 0. The Teletype sets FGI to 1 after the word is printed.

b. The process described in **(a)** is very wasteful. The CPU, which is much faster than the Teletype, must repeatedly check FGI and FGO. If interrupts are used, the Teletype can issue an interrupt to the CPU whenever it is ready to accept or send data. The IEN register can be set by the CPU (under programmer control)

1.7 If a processor is held up in attempting to read or write memory, usually no damage occurs except a slight loss of time. However, a DMA transfer may be to or from a device that is receiving or sending data in a stream (e.g., disk or tape), and cannot be stopped. Thus, if the DMA module is held up (denied continuing access to main memory), data will be lost.

1.8 Let us ignore data read/write operations and assume the processor only fetches instructions. Then the processor needs access to main memory