



**King Fahd University of Petroleum and Minerals**  
**Department of Computer Engineering**

**COMPUTER ARCHITECTURE COE 308**

**Homework 4**

Student Name:.....

Student ID: .....

Problems	Grading
Question 1 / 10	
Question 2 / 10	
Question 3 / 10	
Question 4 / 10	
<b>TOTAL</b>	

## 1 – Hierarchical Memory Systems

A computer has a processor P with a hierarchical memory system that consists of a 4-ways set-associative cache memory CM and a main-memory MM. P generates a 32-bit memory address ADDR = (TAG, Index, Word). CM has 64-entries. The page size is 32 bytes, where each word addressed by P is one byte. MM has  $T_{mm}=18$  ns access time. Denote by  $(T_{cm})$  and  $(P_h)$  as CM access time and its hit probability, respectively.

Answer each of the following questions:

- Find the number of bits in “Index”, “Word”, and “TAG”.
- Determine the page size in bits.
- Determine the cache capacity in bits.
- Determine the main memory capacity in bits.
- CM can be designed in the following ways: (1) as a 4-way set-associative cache with  $T_{cm} = 1.25$  ns and  $P_h = 0.88$ , (2) as a fully-associative cache with  $T_{cm} = 1.5$  ns and  $P_h = 0.97$ , (3) as a direct-mapped cache with  $T_{cm} = 1$  ns and  $P_h = 0.80$ . Answer each of the above questions:
  - Which of the above three cache designs has the least conflict misses.
  - For what reason the direct-mapped cache has the least access time.
  - For what reason the cache with largest hit probability has the largest access time.
  - Determine which of the above three cache designs gives the best global performance of hierarchical memory system including CM and MM.
- The processor is running a benchmark with the following instruction distribution: (1) 55% of R-type, (2) 18% are loads, (3) 12% are stores, and (4) 15% are branching instructions. Suppose the retained cache solution is as a 4-way set-associative cache with  $T_{cm} = 1.25$  ns and  $P_h = 0.8$  which is used as an instruction cache and as a data cache. Answer each of the following questions:
  - Evaluate the average access time of the memory.
  - Evaluate the average stall time per instruction ( $T_{stall}$ ).
  - Evaluate the CPI for this processor if the ideal CPI is 2 clocks and processor clock rate is 1 GHz.

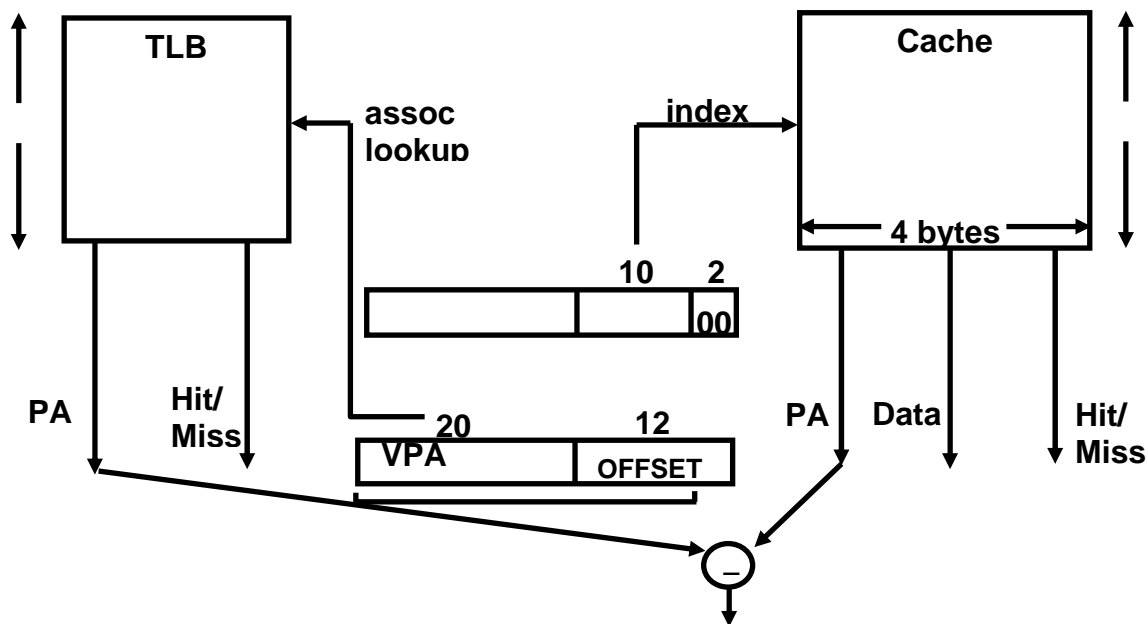
### **Solution:**

- Since CM has 64-entries, then Index is 6 bits. The page size is 32 bytes, then Word is 5 bits. As ADDR is 32 bits, then  $32 - (6+5) = 21$  bits which is the number of bits in TAG.
- The page size is  $32 \times 8 = 256$  bits
- The cache capacity is  $4 \times 64 \times (21 + 256) = 70912$  bits.
- The main memory capacity is  $2^{32} \times 8 = 2^{35}$  bits = 32 Gbits.
- First, considering the three cache designs, the fully-associative cache has the least conflict misses. Second, the direct-mapped cache has the least access time because it has the simplest access hardware. Third, the fully-associative cache has least hit probability and largest access time because it has a 4-to-1 MUX to select one page out of four. Fourth, the above three cache designs have the following average access times: (1)  $T_{av} = 1.25 + (1 - 0.88) \times 18 = 3.41$  ns, (2)  $T_{av} = 1.5 + (1 - 0.97) \times 18 = 2.04$  ns, (3)  $T_{av} = 1 + (1 - 0.8) \times 18 = 4.6$  ns. Therefore the least  $T_{av}$  corresponds to 2.04 (fully-associative) which give the best global performance of hierarchical memory system.
- First, the memory average access time is  $T_{av} = 1.25 + 0.8 \times 18 = 4.85$  ns. Second, the average stall time per instruction  $T_{stall} = 1 \times (1 - 0.8) \times 18 + (0.18 + 0.12) \times (1 - 0.8) \times 18 = 4.68$  ns. Third, the  $CPI = 2 + T_{stall} = 2 + 4.68 = 6.68$  ns.

## 2 – Virtual Memory System

A processor generates a Virtual Address  $VA = (VPA, Offset)$  where VPA is the virtual page address and Offset is the offset of the addressed word within a given VP. To shorten address translation time, the parallel translation scheme is used by simultaneously sending parts of the VA to a TBL and to a direct mapped cache memory.

- a. Give a block diagram of the parallel address translation by referring to CPU, VA, TLB, Cache, Main Memory, and Disk Memory.



**IF cache hit AND (cache tag = PA) then deliver data to CPU  
ELSE IF [cache miss OR (cache tag = PA)] and TLB hit THEN  
access memory with the PA from the TLB  
ELSE do standard VA translation**

- b. Answer only by writing text each of the following issues:
- What part of VA is used to access the TLB, The Virtual page address VPA.
  - how a hit or miss if found in TLB, TLB is an associative memory which is searched using VPA. A hit means VPA is found.
  - how to service a hit or a miss in TLB A hit leads to outputting the corresponding physical page address PPA. A miss lead to fetching a VP from Disk, make placement in MM and update the TBL for translation.
  - What is retrieved from TLB Once VPA is found in TLB, the corresponding physical page address PPA is retrieved from TLB to help computing the physical address.

- e) What part of VA is used to access the Cache  
Some part of OFFSET is used as index for accessing the cache.
- f) how a hit or miss if found in cache,  
A Hit if the TAG from (PPA,OFFSET) matches the stored TAG. Else a miss.
- g) How to service a hit or a miss in cache  
Hit: retrieve page and select word using WORD. Miss: retrieve page from MM, placement in cache, service the CPU and update cache entry.
- h) What is retrieved from the cache.  
Only when we have a hit we retrieve a page and select a word for the CPU.

### **3 - Multiprocessor System**

1. Shortly describe the main difference between SIMD and MIMD multiprocessors in terms of number of decoding units, and style of communication among processors.
2. Consider the following sequential program (SP):  
For i=0,N-1  
     $A(i) = (b(i) - c(i)) * b(i)$   
Endfor

Where N=1000.

Answer each of the following questions:

Translate program SP to an SIMD program (SP-SIMD) targeted to an SIMD multiprocessor with NP=20 as the number of processing elements. Each processing element consists of an ALUs and a small set of registers. Your SP-SIMD program should refer to loading array data, computing, and storing results. Assume local array addresses are already available in local registers.

### **Solution:**

1. The main difference between SIMD and MIMD multiprocessors are (1) SIMD has one Decoding unit (DU) while MIMD has multiple DUs, (2) SIMD operates as a lock step computer where all unit must start and complete in synchrony while MIMD processors are completely asynchronous, (3) SIMD PEs may communicate by either message passing or through sharing memory while in MIMD the processors may share a memory or communicate by exchanging packet messages.
2. Translating the program SP to an SIMD program (SP-SIMD):

```

For I = 0, N/NP-1 ; outer loop over vectorized body
    Load rb, b ; rb is some vector register having NP components
    Load rc, c ; rc is some vector register having NP components
    Subtract reg, rb, rc ; vector subtract
    Multiply reg, reg, rb ; vector multiply
    Store a, reg ; vector store
Endfor

```

The array addresses fetched are function of PEs numbers, i.e. vector data is fetched and each PE buffers its part in the vector. The same process is done in storing vector data.