

## Assignment 7 - Semil Jain (u1417989)

### 1. Large Caches (20 points)

**Soln:**

OS page size = 16KB

LLC size = 32MB

Blocksize = 64 B

#ways = 16 way set-associative

$$\begin{aligned}\text{\#Sets} &= (\text{LLC size}) / (\text{\#ways} * \text{Block size}) \\ &= 32 * 2^{20} / (16 * 64) \\ &= 2^{15} \text{ sets}\end{aligned}$$

Now assuming a 40 or X bit total address (Doesn't matter what length we assume as our offset and index bits would be the least significant bits of the address on the right hand side).

$$\begin{aligned}\text{Blk Offset bits} &= 6 && // \log_2(\text{Block size}) \\ \text{Index bits} &= 15 && // \log_2(\text{\#sets})\end{aligned}$$

Our physical memory will have the same 40 bit as assumed. Out of which:

Page offset = 14 bits  $// \log_2(\text{Page size})$

Remaining bits indicate the page number.

In order for OS to have full flexibility to select the tiles our bits that decide this need to be in the **Intersection of Page number bits and index bits =  $(15 + 6 - 14) = 7$  bits.**

The 7 leftmost bits in the index can be used to point to the tiles in order to have the control with the OS.

Since the question asks maximum tiles we can use all 7 bits. Hence,

$$\text{Total Max Tiles} = 2^7 = 128 \text{ tiles}$$

### 2. Virtually Indexed Cache (20 points)

**Soln:**

OS min page size: 16KB

Assuming L1 cache must be 2-way set-associative

To find the largest L1 cache that we can design.

Since our OS page size is 16KB the page offset bits are 14 bits.

Now if we want to virtually index our physically tagged cache then our index + block offset bits need to be in sync with our page offset bits. Hence total index + block offset bits can be 14 bits. These 14 bits can be partitioned into block sizes and #sets of any sizes and If we design a 1-way cache then the total cache size adds up to 16KB of size.

Now we know that the number of sets and no. of block sizes are kind of invariant. But we can change the no. of ways. The question mentions that L1 must be a 2-way cache.

**Hence, our L1 cache size to design must be  $2 * 16 = 32\text{KB}$ .**

### 3. Organizing Ranks (20 points)

**Soln:**

No. of processor sockets: 2

No. of memory channels: 6

No. Of ranks: 4

DRAM chip capacity: 2Gb, 4Gb or 8 Gb, To maximize our capacity we will take 8Gb chips.

Data output width per chip - 4, 8 or 16. To maximize our capacity we will select a lesser data output width (4).

Total chips per rank = 64b output / each chip width 4 = 16 chips

**Maximum Capacity** =  $[2 \text{ sockets} * 6 \text{ channels} * 4 \text{ ranks} * 16 \text{ chips} * 8\text{Gb}]$   
= 6144 Gb  
= **768 GB (Gigabytes)**

Each memory channel operates at 1.2GHz. Hence,

**Memory Bandwidth supported** =  $2 \text{ sockets} * 6 \text{ channels} * 1.2\text{G} * 2 \text{ (DDR, hence 2 transfers per cycle)} * 64 \text{ (bits per transfer)} = 1843.2 \text{ Gbps}$   
= **230.4 GBps**

### 4. Refresh (20 points)

**Soln:**

Total system capacity: 1 TB = 1024 GB

We have 32 ranks, Hence:

Each rank =  $1024 / 32 = 32 \text{ GB rank}$

Each rank has 16 banks:

Each bank =  $32 / 16 = 2\text{GB bank}$

A row has capacity: 8KB in a bank

Total rows:  $2\text{GB} / 8 \text{ KB} = 256 \text{ K rows each bank}$

Acc. To DDR standards there are 8K refresh commands

**Each refresh command handles** =  $256\text{K} / 8\text{K}$   
= **32 rows**

It takes 40ns to refresh one row, Hence:

Total refresh time =  $32 * 40 = 1280 \text{ ns}$

Refresh command is issued every 6.8 micro seconds. Hence,

**Memory system unavailable performing refresh:**  $1280 \text{ ns} / 7.8 \text{ micro s}$   
=  $1280 \text{ ns} / 7800 \text{ ns}$   
=  $0.1641 * 100 \%$   
= **16.41 %**

## 5. Row Buffers (20 points)

**Soln:**

Assumptions:

- Bus latencies to be 0
- Bank is already recharged at the start
- Row buffer hit - 20ns
- Row buffer conflict - 60ns
- Empty row/ recharged row - 40ns

Row being accessed	Arrival time at memory controller	Hit/Miss (Open)	Open-Page	Hit/Miss (Close)	Close-Page
X	10 ns	empty	<b>50ns</b>	empty	<b>50ns</b> (Precharge to 70)
X	75 ns	rbh	<b>95ns</b>	empty	<b>115ns</b>
Y	100 ns	rbc	<b>160ns</b>	rbc	<b>175ns</b> (Precharge to 195)
X	190 ns	rbc	<b>250ns</b>	empty	<b>235ns</b> (Precharge to 255)
X	280 ns	rbh	<b>300ns</b>	empty	<b>320ns</b>
Y	290 ns	rbc	<b>360ns</b>	rbc	<b>380ns</b>