Full Name:
------------

Andrew ID (print clearly!):\_\_\_\_\_

# 740: Computer Architecture, Fall 2013

# **SOLUTIONS TO Midterm II**

November 25, 2013

#### Instructions:

- Make sure that your exam has 14 pages and is not missing any sheets, then write your full name and Andrew login ID on the front.
- This exam is closed book. You may not use **any** electronic devices. You may use one **single-sided** page of notes that you bring to the exam.
- Write your answers in the box provided below the problem. If you make a mess, clearly indicate your final answer.
- Be concise. You will be penalized for excessive verbosity. Use no more than 15 words per answer, unless otherwise stated.
- The exam lasts 1 hour 20 minutes.
- The problems are of varying difficulty. The point value of each problem is indicated. Do not spend too much time on one question. Good luck!

Problem	Your Score	Possible Points
1		55
2		28
3		30
4		30
5		42
6		20
Total		205

Please read the following sentence carefully and sign in the provided box:

"I promise not to discuss this exam with other students until Wednesday, November 27."

Signature:

# Problem 1: Potpourri (55 pts)

#### A) [10 pts] Thread prioritization

Suppose we are running a **multithreaded** application where threads are part of the same application on a multicore processor. The memory controller is shared between the cores.

**1**) Provide one reason why prioritizing a memory non-intensive thread over a memory-intensive one in the memory controller would improve performance. If this is not possible, write N/A and explain why.

Prioritizing latency-sensitive (memory non-intensive) threads can increase system throughput

**2**) Provide one reason why doing the same would degrade performance. If this is not possible, write N/A and explain why.

Can delay the critical/bottleneck thread which may not be memory non-intensive

#### B) [4 pts] Memory bandwidth

Under what conditions would an application's performance increase linearly as memory bandwidth is increased?

If memory bandwidth is the performance bottleneck

#### C) [4 pts] Fat trees

What problem does the fat tree interconnect solve that is present in the tree interconnect?

High link contention between root and subnodes – a fat tree increases the bandwidth of these links

#### D) [10 pts] Interconnect

You are observing a system with many processing elements connected through a network. There is currently no activity on the network (no messages are being sent). On cycle 10, one of the cores generates a message destined for a cache bank somewhere else on the network. You observe the network on cycle 20 and see that this message has not departed the source location. Assume that all components are enabled (not powered off) and operating at full speed. There are no other messages present in the system at this time. Why could this be?

The system is using circuit switching, and there is a large delay to set up all links between source and destination.

### E) [12 pts] Slack

As you recall, we have discussed the idea of slack based prioritization for on-chip interconnects in class. In fact, you reviewed a paper that introduced this concept. The key idea was to prioritize the packet that has the least slack over others in the router, where the slack of a packet (ideally) is defined as the number of cycles the packet can be delayed without hurting performance.

The concept of slack is actually more general. It can be applied to prioritization at any shared resource, assuming the slack of a "memory request" can be estimated well.

1) Suppose we have a mechanism that tries to estimate the exact slack of a memory request when the request is injected into the shared resources. Provide two reasons why estimating the exact slack of a packet might be difficult:

The exact latency of the request may not be known at the time of injection – the slack may change based on the state of the shared resources and the decisions made by them

How much the packet would affect performance may not be known at the time of injection – the overlap of latency of the packet may not be known at the time of injection

2) What performance issue can slack-based prioritization cause to other processors in the system? Why?

Can cause starvation to some threads

#### 3) How can you solve this problem?

Batching

#### F) [5 pts] Dataflow

What is the purpose of token tagging in dynamic dataflow architectures?

Supporting re-entrant code. Ensuring that tokens come from same context.

# G) [10 pts] Alpha 21264

The Alpha 21264 had a "Prefetch and evict next" instruction that "prefetched data into the L1 cache except that the block will be evicted from the L1 data cache on the next access to the same data cache set."

1) What access patterns could benefit from this instruction? Explain well.

Streaming or striding access pattern (no data reuse)

**2)** The Alpha 21264 processor employed a predictor that predicted whether a load would hit or miss in the cache before the load accessed the cache. What was the purpose of using this predictor? Explain concisely but with enough detail.

Allow speculative scheduling of consumers of the load

## Problem 2: Multithreading (28 pts)

Suppose your friend designed the following fine-grained multithreaded machine:

- The pipeline has 22 stages and is 1 instruction wide.
- Branches are resolved at the end of the 18th stage and there is a 1 cycle delay after that to communicate the branch target to the fetch stage.
- The data cache is accessed during stage 20. On a hit, the thread does not stall. On a miss, the thread stalls for 100 cycles, fixed. The cache is non-blocking and has space to accommodate 16 outstanding requests.
- The number of hardware contexts is 200.

Assuming that there are always enough threads present, answer the following questions:

**A) [7 pts]** Can the pipeline **always** be kept full and non-stalling? Why or why not? (*Hint: think about the worst case execution characteristics.*)

CIRCLE ONE: YES NO

NO - will stall when more than 16 outstanding misses in pipe

**B)** [7 pts] Can the pipeline always be kept full and non-stalling if all accesses hit in the cache? Why or why not?

CIRCLE ONE: YES NO

YES - switching between 200 threads is plenty to avoid stalls due to branch prediction delay

C) [7 pts] Assume that all accesses hit in the cache and your friend wants to keep the pipeline always full and non-stalling. How would you adjust the hardware resources (if necessary) to satisfy this while minimizing hardware cost? You cannot change the latencies provided above. Be comprehensive and specific with numerical answers. If nothing is necessary, justify why this is the case.

Reduce hardware thread contexts to 19, the minimum to keep pipe full/non-stalling

**D**) **[7 pts]** Assume that all accesses miss in the cache and your friend wants to keep the pipeline **always** full and non-stalling. How would you adjust the hardware resources (if necessary) to satisfy this while minimizing hardware cost? You cannot change the latencies provided above. Be comprehensive and specific with numerical answers. If nothing is necessary, justify why this is the case.

Reduce hardware thread contexts to 100, the minimum to keep pipe full/non-stalling. Increase capability to support 100 outstanding misses

#### Problem 3: Return of Tomasulo's Algorithm (30 pts)

The diagram below shows a snapshot at a particular point in time of various parts (reservation stations and register alias table) of the microarchitecture for an implementation supporting out-of-order execution in the spirit of Tomasulo's Algorithm. Note that there is an adder and a multiplier in this machine. The processor is supplied with a seven instruction program following reset. The state below was captured at some point in time during the execution of these seven instructions. Anything marked with a - is unknown and can't be relied upon for your answer. You should assume that the bottommost instruction in the reservation station arrived earliest and the topmost instruction in the reservation station arrived last.

SRC 1	
v	Tag

ID	V	Тад	Val	V	Тад	Val
-	-	-	-	-	-	-
E	0	С	-	0	А	-
F	0	А	-	1	-	20
С	0	А	-	0	А	-

SRC 2

		SRC 2				
ID	D V		V Tag		Val	V
-	-	-	-	-		
D	0 В -		-	0		
В	1	-	20	0		
А	1	-	20	1		

$\wedge$	$\land$
>	
$\sim$	$\bigvee$



Val

\_

\_

-

30

Tag -

Е

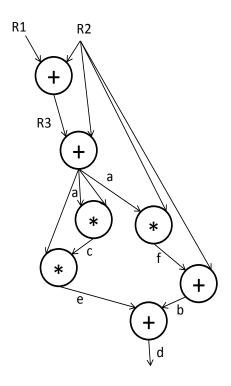
F

-

RAT

Reg	V	Tag	Val
RO	1	-	5
R1	0	А	10
R2	0	В	20
R3	1	-	30
R4	0	С	40
R5	0	D	50
R6	1	-	60
R7	1	-	70

A) [15 pts] Identify the instructions and draw the data flow graph for the seven instructions (use + for ADD and \* for MUL). Please label the edges of the data flow graph with the destination register tag if known. Label with register number if the tag is not known. Note that the first instruction is an ADD with destination register R3.



B) [15 pts] Fill in the instruction opcodes, source, and destination registers in the table below.

Instructions								
ОР	DEST	SRC1	SRC2					
ADD	R3	R1	R2					
ADD	R1	R2	R3					
MUL	R4	R1	R1					
MUL	R5	R2	R1					
ADD	R2	R2	R5					
MUL	R5	R4	R1					
ADD	R5	R2	R5					

#### Problem 4: Tiered-difficulty (30 pts)

Recall from your required reading on Tiered-Latency DRAM that there is a near and far segment, each containing some number of rows. Assume a very simplified memory model where there is just one bank and there are two rows in the near segment and four rows in the far segment. The time to activate and precharge a row is 25ns in the near segment and 50ns in the far segment. The time from start of activation to reading data is 10ns in the near segment and 15ns in the far segment. All other timings are negligible for this problem. Given the following memory request stream, determine the optimal assignment (minimize average latency of requests) of rows in the near and far segment (assume a fixed mapping where rows cannot migrate, a closed-row policy, and the far segment is inclusive).

```
time Ons: row 0 read
time 10ns: row 1 read
time 100ns: row 2 read
time 105ns: row 1 read
time 200ns: row 3 read
time 300ns: row 1 read
```

#### **Detailed solution**

If you were to map 0 and 2 (this is the answer) to near segment:

```
row 0: activated at time = 0
row 0: read at time = 10 (10ns latency)
row 1: activated at time = 25
row 1: read at time = 40 (30ns latency)
row 2: activated at time = 100
row 2: read at time = 110 (10ns latency)
row 1: activated at time = 125
row 1: read at time = 140 (35ns latency)
row 3: activated at time = 200
row 3: read at time = 215 (15ns latency)
row 1: activated at time = 300
row 1: read at time = 315 (15 ns latency)
```

#### total latency is 115ns

If you were to map 1 and 2 (an example incorrect answer) to near segment:

```
row 0: activated at time = 0
row 0: read at time = 15 (15ns latency)
row 1: activated at time = 50
row 1: read at time = 60 (50ns latency)
row 2: activated at time = 100
row 2: read at time = 110 (10ns latency)
row 1: activated at time = 125
row 1: read at time = 135 (30ns latency)
row 3: activated at time = 200
row 3: read at time = 215 (15ns latency)
row 1: activated at time = 300
row 1: read at time = 310 (10 ns latency)
```

#### total latency is 130ns

A) [6 pts] What rows would you place in near segment? Hint: draw a timeline.

rows 0 and 2. see above

**B**) **[6 pts]** What rows would you place in far segment?

rows 1 and 3 (also rows 0 and 2 since inclusive). see above

C) [6 pts] In 15 words or less, describe the insight in your mapping?

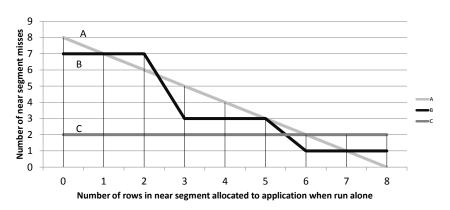
See TL-DRAM's WMC policy – the first access in near simultaneous requests causes the second to wait activation + precharge time. minimizing this wait by caching first row in near segment is better than caching second row in near segment (this decreases only time to read from start of activation), even if second row is accessed more frequently (see example above)

**D**) [6 pts] Assume now that the mapping is dynamic. What are the tradeoffs of an exclusive design vs. an inclusive design? Name one advantage and one disadvantage for each.

Exclusive requires swapping, but can use nearly full capacity of DRAM. Inclusive, the opposite.

Continued ...

**E)** [6 pts] Assume now that there are eight (8) rows in the near segment. Below is a plot showing the number of misses to the near segment for three applications (A, B, and C) when run alone with the specified number of rows allocated to the application in the near segment. This is similar to the plots you saw in your Utility-Based Cache Partitioning reading except for TL-DRAM instead of a cache. Determine the optimal static partitioning of the near segment when all three of these applications are run together on the system. In other words, how many rows would you allocate for each application? Hint: this should sum to eight. Optimal for this problem is defined as minimizing total misses across all applications.



1) How many near segment rows would you allocate to A?

	5	
2) How many near segme	ent rows would you allocate to B?	
	3	
3) How many near segme	ent rows would you allocate to C?	
	0	

#### Problem 5: GPUs (42 pts)

We define the **SIMD utilization** of a program running on a GPU as the fraction of SIMD lanes that are kept busy with active threads during the run of a program.

The following code segment is running on a GPU. Each thread executes a single iteration of the shown loop. Assume that the data values of the arrays A, B, and C are already in vector registers so there are no loads and stores in this program. (Hint: Notice that there are 5 instructions in each thread as labled below.) A warp in the GPU consists of 64 threads, and there are 64 SIMD lanes in the GPU.

A) [2 pts] How many warps does it take to execute this program?

16384/64 = 256

**B**) **[10 pts]** As shown below, assume array A has a repetitive pattern which has 32 ones followed by 96 zeros repetitively and array B has a different repetitive pattern which has 64 zeros followed by 64 ones repetitively. What is the SIMD utilization of this program?

A: 1	1	29 1s	1	0	0	93 Os	0	32 1s	96 0s	
B: 0	0	61 0s	0	1	1	61 1s	1	64 0s	64 1s	

When a warp is working on a segment of array A that has 64 0s, none of the threads in the warp will take the branch, which yields no branch divergence of the warp. Hence, the SIMD utilization of this particular input set is (64 + 64 + 32 \* 4)/(64 + 64 \* 5) = 66.7%

Continued ...

C) [10 pts] Is it possible for this program to yield a SIMD utilization of 25%?

CIRCLE ONE: YES NO

If YES, what should be true about arrays A and B for the SIMD utilization to be 25%? Be precise and show your work. If NO, explain why not.

Yes. For example, if only 4 elements in **every 64** elements of A are positive, we can have a SIMD utilization of (64 + 4 \* 4)/(64 \* 5) = 25%.

**D**) **[10 pts]** Is it possible for this program to yield a SIMD utilization of 20%?

CIRCLE ONE: YES NO

If YES, what should be true about arrays A and B for the SIMD utilization to be 20%? Be precise and show your work. If NO, explain why not.

No. The smallest SIMD utilization one can get is to have one and only one element in every 64 elements of A to be positive, which yields a minimal SIMD utilization of (64 + 1 \* 4)/(64 \* 5) = 21.25%, which is still greater than 20%.

**E**) **[10 pts]** During an execution with a particular input array A, which has exactly 24 positive elements in every 64 elements, Hongyi finds that the SIMD utilization of the program is 50%. Based on this observation, Hongyi claims that any input array that has an **average** of 24 out of 64 elements positive would yield a 50% SIMD utilization. Is Hongyi correct?

CIRCLE ONE: YES

If YES, show your work. If NO, provide a counterexample.

Hongyi is incorrect. If A has a repetitive pattern of 48 contiguous 1s followed by 80 contiguous 0s, in which case 37.5% of the elements are positive on average, then the SIMD utilization of the program will be 83.3% rather than 50%.

NO

### Problem 6: Hyperblock (20 pts)

As described in class, Hyperblock scheduling uses predication support to replace unbiased branches with predicates, which enables larger code blocks.

A) [2 pts] In one sentence, in terms of code optimizations, explain what benefit does larger scheduling code blocks provide?

Larger scheduling code blocks enable greater flexibility for instruction scheduling.

One optimization that can be applied to Hyperblock is **Instruction Promotion**. Instruction Promotion hoists the operation from a predicated instruction and replaces the original predicated instruction with a conditional move. With Instruction Promotion, operations can be scheduled and issued before their corresponding predicates are determined. Below shows an example of Instruction Promotion.

Before:

```
cmplt B6,B7-> B0
[B0] ld MEM1-> A5
[!B0] ld MEM2-> A5
nop 4
addi A5,8 -> A8
```

After Instruction Promotion:

```
ld MEM1-> A5
ld MEM2-> A6
cmplt B6,B7-> B0
[!B0] mv A6-> A5
nop 4
addi A5,8 -> A8
```

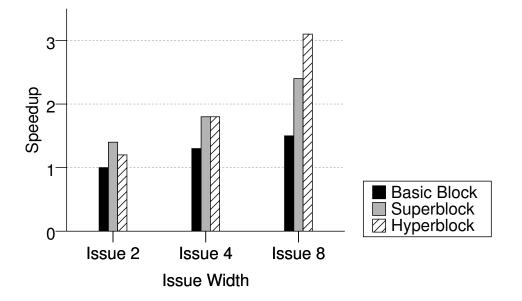
Assume we run this code on a processor that supports predication but can only issue a predicated instruction after its corresponding predicate has been resolved.

**B) [4 pts]** For the example above, can Instruction Promotion ever improve system performance? Why or why not?

Yes it can. With Instruction Promotion, the program can hide some of the load latency.

**C) [4 pts]** For the example above, can Instruction Promotion ever degrade system performance? Why or why not?

Yes it can. Instruction Promotion: (1) introduces extra instructions and (2) can increase register pressure. [Note that extra instructions may not always increase register pressure]



**D**) **[10 pts]** The graph above shows the performance comparison of a program optimized using Hyperblock and Superblock respectively with different issue widths. With all other factors being equal, as the figure shows, when the issue width is low, Superblock provides higher speedup than Hyperblock. However, when the issue width is high, Hyperblock provides higher speedup than Superblock. Explain why this can happen?

A wider issue width can tolerate the wasted instructions in a hyperblock, but does not benefit the superblock (all else being equal).

A more detailed explanation: Hyperblock uses predication which increases the total number of instructions to execute. When the issue width is low, executing extra predicated instructions requires extra work, which slows down the processor as all resources of the processor has already been fully utilized. When the issue width is high, however, Hyperblock provides a greater number of independent instructions from the multiple paths of control to fill the available processor resources. As Hyperblock also enables larger code blocks for better optimization for unbiased branches, Hyperblock provides better speedup.