CS 433 Midterm Exam Practice

Professor Sarita Adve

Time: 2 Hours

Please print your Name and NetID and circle your course section below.

Name:		
NetID:		
Section:	T3 (Undergraduate)	T4 (Graduate)

Instructions

- 1. No books, papers, notes, or any other typed or written materials are allowed. No calculators or other electronic materials are allowed.
- Please do not turn in loose scrap paper. Limit your answers to the space provided if possible. If this is not possible, please write on the back of the same sheet. You may use the back of each sheet for scratch work.
- 3. In all cases, show your work. No credit will be given if there is no indication of how the answer was derived. Partial credit will be given even if your final solution is incorrect, provided you show the intermediate steps in reaching the final solution.
- 4. If you believe a problem is incorrectly or incompletely specified, make a reasonable assumption and solve the problem. The assumption should not result in a trivial solution. In all cases, clearly state any assumptions that you make in your answers.
- 5. This exam has **5 problems** and **13 pages** (including this one). **All students** should solve **problems 1, 2A, and 3 through 5**. Only **graduate students should solve problem 2B and 2C**. Please budget your time appropriately. Good luck!

Problem	1	2	3	4	5	Total
Points T3	4 4	16 26	9 9	18 18	9 9	56 (undergrads) 66 (graduates)
Score						

Problem 1 [4 points]

Assume 90% of a sequential program's execution time can be parallelized.

Part A [2 points]

What speedup (from parallelism) is required on the parallelizable section to achieve an overall speedup of 4X for the full program? You may express your answer in terms of an equation with all variables explicitly substituted. You are not required to perform numerical calculations.

Part B [2 points]

What is the maximum possible speedup achievable on the above program through parallelization?

Problem 2 [16 points for undergraduates, 26 points for graduates]

This problem concerns Tomasulo's algorithm. Consider running the following loop on an architecture specified below.

1.	LP:	L.D	F0, 0(R1)
2.		ADD.D	F0, F0, F6
3.		DIV.D	F2, F2, F0
4.		L.D	F0, 8(R1)
5.		DIV.D	F4, F0, F8
6.		S.D	F4, 16(R1)
7.		DADDI	R1, R1, #-24
8.		BNEZ	R1, LP

Functional Unit Type	Cycles in EX	Number of Functional Units
Integer	1	1
FP Adder	5	1
FP Divider	15	1

- 1) Assume that you have unlimited reservation stations.
- 2) Memory accesses use the integer functional unit to perform effective address calculation during the EX stage. For stores, memory is accessed during the EX stage (Tomasulo's algorithm without speculation) or commit stage (Tomasulo's algorithm with speculation). All loads access memory during the EX stage. Loads and Stores stay in EX for 1 cycle.
- 3) Functional units are not pipelined.
- 4) If an instruction moves to its WB stage in cycle x, then an instruction that is waiting on the same functional unit (due to a structural hazard) can start executing in cycle x.
- 5) An instruction waiting for data on the CDB can move to its EX stage in the cycle after the CDB broadcast.
- 6) Only one instruction can write to the CDB in one clock cycle. Branches and stores do not need the CDB.
- 7) Whenever there is a conflict for a functional unit or the CDB, assume that the oldest (by program order) of the conflicting instructions gets access, while others are stalled.
- 8) Assume that the result from the integer functional unit is also broadcast on the CDB and forwarded to dependent instructions through the CDB (just like any floating point instruction).
- 9) Assume that the BNEZ occupies the integer functional unit for its computation and spends one cycle in EX.

Part A [16 points]

Complete the following table using Tomasulo's algorithm but without assuming any hardware speculation on branches. That is, an instruction after a branch cannot issue until the cycle after the branch completes its EX. Assume a *single-issue* machine. Fill in the cycle numbers in each pipeline stage for the first several instructions of the loop as given below, assuming the branch is always taken. The entries for the first instruction are filled in for you. Explain the reasons for any stalls.

Instruction	IS	EX	WB	Reason for Stalls
Iteration 1				
1. L.D F0, 0(R1)	1	2	3	
2. ADD.D F0, F0, F6				
3. DIV.D F2, F2, F0				
4. L.D F0, 8(R1)				
5. DIV.D F4, F0, F8				
6. S.D F4, 16(R1)				
7. DADDI R1, R1, #-24				
8. BNEZ R1, LP				
Iteration 2				
9. L.D F0, 0(R1)				

ONLY GRADUATE STUDENTS SHOULD SOLVE THE FOLLOWING PARTS B and C FOR PROBLEM 2.

Part B [4 points]

To improve the performance of the above loop execution, you are tasked with adding support for speculative execution. However, instead of using a reorder buffer (ROB) to store uncommitted (potentially speculative) values, your colleague proposes to build on the concept of renaming used in Tomasulo's algorithm and leverage the available large set of physical registers. The idea is to hold potentially speculative, uncommitted values that usually reside in the ROB in an extended set of physical registers. During instruction issue, a renaming process maps the names of the logical (i.e., architectural) registers (R0 to R31, F0 to F31) to this extended physical register set (P0 to P255), allocating a new unused register for the destination from a queue of free P-registers.

Using this idea, work out the register renaming necessary to enable the speculative execution of (part of) the second loop iteration in the table below. Assume that at the time the sequence in the table begins, registers R0 to R31 are mapped to P0 to P31 respectively and registers F0 to F31 are mapped to P32 to P63 respectively. Assume P64 to P255 are all on the free list and allocation of free registers is done in order from P64 to P255 (the lowest numbered one first), followed by any others that may be freed up in the execution below.

Instruct	Instruction		ion with renamed registers	Changes to the rename map
BNEZ	R1, LP	BNEZ	P1, LP	
Iteration	ı 2			
L.D	F0, 0(R1)	L.D	P64, 0(P1)	$F0 \rightarrow P64$
ADD.D	F0, F0, F6			
DIV.D	F2, F2, F0			
L.D	F0, 8(R1)			
DIV.D	F4, F0, F8			

ONLY GRADUATE STUDENTS SHOULD SOLVE THE FOLLOWING PART C FOR PROBLEM 2.

Part C [6 points]

Suppose that the first load instruction - L.D F0, 0(R1) - in the second loop iteration raises an exception. Assume the exception is raised after all the instructions in the table have arrived in the reorder buffer. Also assume the hardware provides precise exceptions.

When should the above exception be handled? For precise exceptions, the instructions after the load must be squashed before the exception is handled. Explain what steps the hardware needs to take to restore the correct architectural state after these instructions are squashed. If there is any additional information that the hardware needs to track to enable this, state the information clearly. Illustrate your answer by applying it to at least one of the squashed instructions.

Problem 3 [9 points]

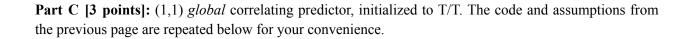
Consider the following piece of code:

```
DADDI R1, R0, #100
L1: DADDI R1, R1, #-1
BEQZ R1, END -- Branch 1
DADDI R12, R0, #2
L2: DADDI R12, R12, #-1
BNEZ R12, L2 -- Branch 2
J L1
END: ...
```

Assume R0 stores 0. Branch 1 is executed 100 times and branch 2 is executed a total of 198 times. For each branch, how many correct predictions will occur if we use the following prediction schemes? Assume at the beginning of execution, the last branch was *not taken*. Please explain your answers.

Part A [3 points]: 1-bit predictor initialized to T (taken)

Part B [3 points]: 2-bit saturating counter predictor initialized to 10 (taken)



The code and assumptions for Problem 3 are repeated below from the previous page for your convenience.

```
DADDI R1, R0, #100
L1: DADDI R1, R1, #-1
BEQZ R1, END -- Branch 1
DADDI R12, R0, #2
L2: DADDI R12, R12, #-1
BNEZ R12, L2 -- Branch 2
J L1
END: ...
```

Assume R0 stores 0. Branch 1 is executed 100 times and branch 2 is executed a total of 198 times. Assume at the beginning of execution, the last branch was *not taken*.

Problem 4 [18 points]

Consider the following code fragment:

Loop:	LD.D	F1, 0(R1)
	LD.D	F2, 0(R2)
	MUL.D	F3, F1, F10
	ADD.D	F4, F3, F2
	SD.D	F4, 0(R1)
	DADDUI	R1, R1, #8
	DADDUI	R2, R2, #8
	BNEZ	R1, R3, Loop

Consider a pipeline with the following latencies: 3 cycles between an FP multiply and its consumer, 1 cycle between an FP add and its consumer, and 0 cycles between all other pairs. Thus, there should be three stall cycles between the multiply and addition in the above code for correct operation. Assume that all functional units are pipelined. Assume the machine does NOT support delayed branches.

Unroll the above loop 4 times and write the resulting code to the left of the table on the next page (the above loop is repeated on the next page for your convenience). You have access to temporary registers T0...T63. Assume that the total number of iterations for the original loop is a multiple of 4. Schedule the unrolled loop for best performance on a VLIW machine where each VLIW instruction can contain one memory reference, one FP operation, and one integer operation. Write the scheduled instructions in the table on the next page to minimize the number of stalls. You may use L for L.D, M for MUL.D, etc.

Loop: LD.D F1, 0(R1)	Mem	FP ALU	Integer ALU
LD.D F2, 0(R2) MUL.D F3, F1, F10			
ADD.D F4, F3, F2 SD.D F4, 0(R1)			
DADDUI R1, R1, #8			
DADDUI R2, R2, #8 BNEZ R1, R3, Loop			
Please write the unrolled loop below			

Problem 5 [9 points]

Consider the following code fragment from an if-then-else statement of the form

if
$$(A == 0) A = B$$
;
else $A = A + 4$;

where A is at 0(R3) and B is at 0(R2):

1.		LD	R1, 0(R3)	; load A
2.		BNEZ	R1, L1	; test A
3.		LD	R1, 0(R2)	; then clause
4.		J	L2	; skip else
5.	L1:	DADDI	R1, R1, #4	; else clause
6.	L2:	SD	R1, 0(R3)	; store A

The machine on which the code will run does not have hardware speculation but does support compiler speculation where speculative instructions are marked with a speculation bit and poison bits are provided to deal with exceptions on speculative instructions. You are told that if two loads appear one after another in program order and happen to miss in the cache, then the machine can pipeline them in the memory system with significant performance boost. (We will learn about this type of load optimization soon in class. It is not necessary to know how it works to solve this problem.)

Given the above, you modify the above code as below to exploit compiler speculation and move the two loads next to each other for a possible performance boost. The second load is speculative (denoted by (s)).

```
LD
                     R1, 0(R3)
                                   ; load A
                     R14, 0(R2)
       (s)LD
                                   ; speculative load B
       BEQZ
                     R1, L3
                                   ; other branch of the if
       DADDI
                     R14, R1, #4
                                   ; else clause
L3:
       SD
                     R14, 0(R3)
                                   ; store A
```

Part A [4 points]

Assume A = 0 at the beginning of the execution and that the speculative load above incurs an exception. Fill the following table for the execution of your modified code above, showing the instructions executed (in program order), the speculative bit for the instruction (0 or 1), and the state of the poison bits for the different registers after the instruction is executed (0 or 1). State if and when (i.e., at which instruction) the exception incurred by the speculative load is handled and why. The first entry is filled for you.

	Speculative bit		Poiso	n bits	
Instruction		R1	R2	R3	R14
LD R1, 0(R3)	0	0	0	0	0

Part B [5 points]

Now assume that A!=0 at the start of the execution of your modified code and redo the following table with all of the other assumptions and instructions of Part A:

	Speculative bit	Poison bits			
Instruction		R1	R2	R3	R14
LD R1, 0(R3)	0	0	0	0	0