

CS 433 Midterm Exam: Oct 18, 2021

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.
2. 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 **14 pages** (including this one). **All students** should solve **problems 1 through 4**. Only **graduate students** should solve **problem 5**. Please budget your time appropriately. Good luck!

Problem	1	2	3	4	Graduate students 5	Total
Points	6	17	26	14	8	63 (undergrads) 71 (graduates)
Score						

Problem 1 [6 Points]

Suppose an important program you run has the following characteristics.

Instruction Type	% of Execution Time
Load from memory	12%
Store to memory	6%
FP Multiplication	17%
Other	65%

Part A [3 points]

You are considering upgrading your machine to one of two possible configurations. M1 reduces the contribution of loads to the execution time by $3\times$ and that of stores by $2\times$. M2 reduces the contribution of (only) floating-point multiplications to the execution time by $10\times$. What is the overall speedup of the program for each upgrade? 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 [3 points]

What is the maximum possible speedup if all memory accesses could be sped up infinitely? What is the maximum speedup if all (and only) floating-point multiplications could be sped up infinitely? You may express your answer in terms of an equation with all variables explicitly substituted. You are not required to perform numerical calculations.

Problem 2 [17 Points]

Consider running the following code on a machine with the assumptions below:

1. *loop*: L.D F1, 0(R5)
2. S.D F0, 0(R5)
3. ADD.D F0, F0, F1
4. DADDIU R4, R4, #-1
5. DADDIU R5, R5, #8
6. BNEZ R4, *loop*

- The machine has a **dual-issue**, out-of-order processor with hardware speculation, a **reorder buffer with 8 entries**, and the following functional units:

Functional Unit Type	Number of Functional Units	Cycles in EX
Integer ALU	1	1
FP Adder	1	4 (not pipelined)

- 2 instructions can be issued and committed each cycle.
- Loads use the Integer ALU for effective address calculation and for memory access during the EX stage. It takes a single cycle for the ALU to do both the address calculation and memory access.
- Stores use the integer ALU for effective address calculation in the EX stage and access memory in the CM stage (1 cycle for each).
- Branches use the Integer ALU for all computation in the EX stage. Branch direction and target are predicted perfectly.
- Instructions following a branch cannot issue in the same cycle as a branch.
- Branches do not have a branch delay slot.
- Only one instruction can write to the CDB in each clock cycle.
- 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 .
- If an instruction cannot move from EX to its WB stage in cycle x , it continues occupying the functional unit it was using.
- An instruction waiting for data from the CDB can move to its EX stage in the cycle after the CDB broadcasts the data.
- Branches and stores do not need the CDB.
- 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.
- Assume that the result from the Integer ALU is also broadcast on the CDB and forwarded to dependent instructions through the CDB, just like any floating point instruction.

The following table shows the first 11 instructions of the above code. Fill in the cycles each instruction is in each stage. Note **all reasons for all stalls**, including the type of hazard; the functional unit or register it is dependent on; and the instruction it is dependent on. **Note any stalls due to commit ordering and/or limited commit width**. Write “None” if there are no stalls. Some entries are already written for you.

#	Instruction	IS	EX	WB	CM	Reasons for Stalls
1	L.D F1, 0(R5)	1	2	3	4	None
2	S.D F0, 0(R5)	1	3	—	4	Structural INT from 1
3	ADD.D F0, F0, F1	2	4-7	8	9	RAW F1 from 1
4	DADDIU R4, R4, #-1	2				
5	DADDIU R5, R5, #8	3				
6	BNEZ R4, <i>loop</i>	3				
7	L.D F1, 0(R5)	4				
8	S.D F0, 0(R5)	4				
9	ADD.D F0, F0, F1	5				
10	DADDIU R4, R4, #-1	5				
11	DADDIU R5, R5, #8					

Problem 3 [26 Points]

Consider a single-issue, in-order five stage pipeline similar to those studied in class, but with the following specification:

Functional Unit	Cycles in EX	Number of Functional Units	Pipelined
Integer	1	1	Yes
FP Add/Subtract	3	1	Yes
FP/Integer Multiplier	8	1	Yes
FP/Integer Divider	24	1	No

- The integer functional unit performs integer addition (including effective address calculation for loads/stores), subtraction, and logic operations.
- There is full forwarding and bypassing, including forwarding from the end of a functional unit to the MEM stage for stores.
- Loads and stores complete in one cycle. That is, they spend one cycle in the MEM stage after the effective address calculation.
- There are as many registers, both FP and integer, as you need.
- Branches are resolved in ID and there is one branch delay slot.
- While the hardware has full forwarding and bypassing, it is the responsibility of the compiler to schedule such that the operands of each instruction are available when needed by each instruction.
- If multiple instructions finish their EX stages in the same cycle, then we will assume they can all proceed to the MEM stage together. Similarly, if multiple instructions finish their MEM stages in the same cycle, then we will assume they can all proceed to the WB stage together. In other words, for the purpose of this problem, you are to ignore structural hazards on the MEM and WB stages.

This problem explores the ability of the compiler to schedule code as efficiently as possible for such a pipeline. Consider the following code (also repeated on the next pages for reference):

```
loop:    L.D      F4, 0(R1)
        MUL.D   F8, F4, F0
        L.D      F6, 0(R2)
        ADD.D   F10, F6, F2
        ADD.D   F12, F8, F10
        S.D     F12, 0(R3)
        DADDIU  R1, R1, #8
        DADDIU  R2, R2, #8
        DADDIU  R3, R3, #8
        DSUB   R5, R4, R1
        BNEZ   R5, loop
```

Part A [6 Points]

Rewrite the above loop (repeated below for reference), but let every row take a cycle (each row can be an instruction or a stall). If an instruction can't be issued on a given cycle (because the current instruction has a dependency that will not be resolved in time), write "stall" instead, and move on to the next cycle (row) to see if it can be issued then. Assume that a NOP is scheduled in the branch delay slot (effectively stalling 1 cycle after the branch). *Explain the cause of all stalls*, but don't reorder instructions. How many cycles elapse before the second iteration begins?

```
loop: L.D      F4,    0(R1)
      MUL.D   F8,    F4,   F0
      L.D      F6,    0(R2)
      ADD.D   F10,   F6,   F2
      ADD.D   F12,   F8,   F10
      S.D     F12,   0(R3)
      DADDIU  R1,    R1,   #8
      DADDIU  R2,    R2,   #8
      DADDIU  R3,    R3,   #8
      DSUB    R5,    R4,   R1
      BNEZ    R5,    loop
```

Part B [6 Points]

Now reschedule the loop to compute the same results as quickly as possible. You can change immediate values and memory offsets and reorder instructions, but don't change anything else. Show any stalls that remain. How many cycles elapse before the second iteration begins? Show your work.

```
loop: L.D      F4,    0(R1)
      MUL.D   F8,    F4,    F0
      L.D      F6,    0(R2)
      ADD.D   F10,   F6,    F2
      ADD.D   F12,   F8,    F10
      S.D     F12,   0(R3)
      DADDIU  R1,    R1,    #8
      DADDIU  R2,    R2,    #8
      DADDIU  R3,    R3,    #8
      DSUB   R5,    R4,    R1
      BNEZ   R5,    loop
```

Part C [6 Points]

Now unroll the loop the minimum number of times needed to eliminate all stalls (with rescheduling). Show the unrolled and rescheduled loop. You can, and should, remove redundant instructions. How many original iterations of the loop are in an iteration of your new unrolled loop? How many cycles elapse before the next iteration of the unrolled loop begins? Don't worry about start-up or clean-up code outside the unrolled loop. Assume a very large number of iterations for the original loop. Show your work.

```
loop: L.D      F4,    0(R1)
      MUL.D   F8,    F4,    F0
      L.D     F6,    0(R2)
      ADD.D   F10,   F6,    F2
      ADD.D   F12,   F8,    F10
      S.D     F12,   0(R3)
      DADDIU  R1,    R1,    #8
      DADDIU  R2,    R2,    #8
      DADDIU  R3,    R3,    #8
      DSUB    R5,    R4,    R1
      BNEZ   R5,    loop
```

Part D [8 Points]

Consider a VLIW processor in which one instruction can support two memory operations (load or store), one integer operation (addition, subtraction, comparison, or branch), one floating point add or subtract, and one floating point multiply or divide. There is no branch delay slot. Now unroll the original loop four times (i.e., four original iterations in one new iteration), and schedule it for this VLIW processor to take as few stall cycles as possible. How many cycles do the four iterations take to complete? Use the table template on the next page to show your work.

#	MEMORY 1	MEMORY 2	INTEGER	FP ADD/SUB	FP MUL/DIV
1					
2					
3					
4					
5					
6					
7					
8					
9					
10					
11					
12					
13					
14					
15					
16					
17					
18					
19					
20					

Problem 4 [14 Points]

Consider a piece of code with three static branch instructions, B1, B2, and B3. During an execution of this code, the global history for these branches (i.e., their execution sequence and direction) is as follows:

Branch	B1	B2	B1	B2	B3	B1	B1	B2	B1
Direction	T	N	N	T	T	T	T	N	T

T stands for Taken; N stands for Not-taken. Thus, the execution starts with branch B1 being taken, then B2 is not taken, etc. Before the above execution, the outcome of all previous branches and the initial state of all predictors is not taken (N). Assume a (2, 1) correlating predictor is used for branch B1. Assume the state of the predictor is recorded in the form W/X/Y/Z where:

- W = state for when the last branch is TAKEN and the branch before the last is TAKEN
- X = state for when the last branch is TAKEN and the branch before the last is NOT TAKEN
- Y = state for when the last branch is NOT TAKEN and the branch before the last is TAKEN
- Z = state for when the last branch and the branch before the last are both NOT TAKEN

Assume the history at a branch is recorded in the form $H_2 H_1$ where H_1 is the last branch and H_2 is the branch before the last one.

Fill the entries in the two tables below assuming the (2,1) correlating predictor for branch B1 uses **local** branch history in Table 1 and **global** branch history in Table 2. (Recall that by local history for B1, we mean the history for only branch B1.)

Table 1: Assume the predictor uses LOCAL branch history

Branch B1 invocation #	History used for prediction	Prediction for B1	Actual direction of B1	New Predictor State
1	N N		T	
2			N	
3			T	
4			T	
5			T	

Table 2: Assume the predictor uses GLOBAL branch history

Branch B1 invocation #	History used for prediction	Prediction for B1	Actual direction of B1	New Predictor State
1	N N		T	
2			N	
3			T	
4			T	
5			T	

ONLY GRADUATE STUDENTS SHOULD SOLVE PROBLEM 5

Problem 5 [8 Points]

In this problem, we try to understand the implications of the Reorder Buffer (ROB) size on performance. Consider a processor implementing the ROB scheme described in class. Recall each instruction goes through issue (IS), writeback (WB), execute (EX), and commit (CM).

- Assume IS, WB, and CM each take one cycle (once all the conditions for these stages are met), as discussed in class.
- Assume our machine can fetch and commit 4 instructions each cycle.
- Assume a branch misprediction is handled when the branch instruction reaches the head of the ROB. It involves flushing that ROB entry and all entries following that entry.
- For now, assume there are no memory accesses (this will change in part C).

Part A [2 Points]

Suppose we have a perfect branch predictor and there is no data dependency between instructions. We have infinite execution units of each type and infinite reservation stations. All instructions take one cycle in the EX stage. What is the maximum achievable IPC? What is the minimum ROB size required to guarantee that IPC?

Part B [2 Points]

Suppose different FUs have different latencies in the EX stage, as given by the following table. Everything else is the same as in the previous part. What is the minimum size of the ROB required now to avoid any issue stalls due to a full ROB?

Functional Unit	Cycles
Integer ALU	1
FP Adder	5
FP Multiplier	10

Part C [2 Points]

In addition to the latencies above, now every 10th instruction is a load instruction. Assume the address calculation and cache/memory access parts of the load both happen in the EX stage. The hit rate in the data cache is 95% and the misses are uniformly spaced through the instruction stream. A hit takes 1 cycle in the EX stage. However, upon a miss, the data has to be fetched from the memory and this results in 100 cycles in the EX stage. What is the ROB size required now to avoid any issue stalls?

Part D [2 Points]

Now additionally assume we don't have perfect branch prediction anymore. Instead, we have a predictor with an accuracy of 95%. Assume every 8th instruction is a branch and the mispredictions are uniformly spaced through the instruction stream.

After a misprediction, how many instructions are issued before the next misprediction is encountered? In light of this result, do you think we need a ROB of the size you derived in Part C? Why/why not? If not, what do you think would be a good ROB size to have?