

**CS 152 Computer Architecture and Engineering
CS 252 Graduate Computer Architecture**

**Midterm #2
April 15, 2020
Professor Krste Asanović**

Name: _____

SID: _____

**I am taking CS152 / CS252
(circle one)**

80 Minutes, 17 pages.

Notes:

- Not all questions are of equal difficulty, so look over the entire exam!
- Please carefully state any assumptions you make.
- Please write your name on every page in the exam.
- Do not discuss the exam with other students who haven't taken the exam.
- If you have inadvertently been exposed to an exam prior to taking it, you must tell the instructor or TA.
- You will receive no credit for selecting multiple-choice answers without giving explanations if the instructions ask you to explain your choice.

Question	CS152 Point Value	CS252 Point Value
1	20	20
2	20	20
3	20	20
4	20	20
TOTAL	80	80

Problem 1: Vectors (20 points)

Problem 1.A (12 points)

State whether each of the following loops can be successfully vectorized and explain your reasoning. If it can only be vectorized under certain circumstances, give requirements on the inputs and any other necessary conditions.

In all cases, assume that A, B, and C are non-overlapping arrays in memory.

```
x = 0;
for (i = 0; i < N; i++) {
    x = x + (A[i] * B[i]);
}
```

```
for (i = 0; i < N; i++) {
    if (C[i] == 0) {
        if (A[i] > B[i])
            C[i] = A[i];
        else
            C[i] = B[i];
    }
}
```

For simplicity, assume N is evenly divisible by 2.

```
for (i = 0; i < N; i++) {  
    A[i] = B[i % (N/2)];  
}
```

```
for (i = 0; i < N; i++) {  
    C[A[i]] = C[B[i]];  
}
```


Problem 2: VLIW (20 points)

In this problem, we will optimize the following code sequence, which implements a prefix sum interleaved with a vector multiply, for a VLIW architecture.

```
for (i = 0; i < N; i++) {
    x = A[i] + x;
    C[i] = x * B[i];
}

# f1 contains initial value of x
addi x1, x0, N*8      # initialize loop boundary
addi x2, x0, 0        # initialize array index
loop:
    fld f2, A(x2)     # load A[i]
    fld f3, B(x2)     # load B[i]
    addi x2, x2, 8    # bump index
    fadd.d f1, f1, f2 # compute x
    fmul.d f4, f1, f3 # compute C[i]
    fsd f4, (C-8)(x2) # store C[i]
    bltu x2, x1, loop
```

- “A”, “B”, and “C” are immediates generated by the compiler that encode the base addresses of arrays A, B, and C, respectively.
- Arrays A, B, and C do not overlap in memory.
- N is a large number that is statically known.
- N is evenly divisible as needed for loop unrolling and software pipelining.
- Register f1 holds the initial value of x.
- Assume that no exceptions arise during execution.

The code is executed on an in-order VLIW machine with five execution units. All execution units are fully pipelined and latch their operands at issue.

- One integer ALU, 1-cycle latency, also used for branches
- One load unit, 2-cycle latency
- One store unit (ignore the latency of memory-memory dependencies for this problem)
- One floating-point adder, 3-cycle latency
- One floating-point multiplier, 3-cycle latency

Instructions are statically scheduled with no interlocks; all latencies are exposed in the ISA. All register operands are read before any writes from the same instruction take effect (i.e., no WAR hazards between operations within a single VLIW instruction).

Execution units write to the register file at the end of their last pipeline stage, and the results become visible at the beginning of the following cycle. There is no bypassing. **Old values can be read from registers until they have been overwritten.** (You may leverage this to more efficiently schedule VLIW code.)

The unoptimized scheduling of the above assembly code is shown in the following table.

Label	ALU	LOAD	STORE	FADD	FMUL
	addi x1, x0, N*8				
	addi x2, x0, 0				
loop:		f1d f2, A(x2)			
	addi x2, x2, 8	f1d f3, B(x2)			
				fadd.d f1, f1, f2	
					fmul.d f4, f1, f3
	bltu x2, x1, loop		fsd f4, (C-8)(x2)		

Problem 2.C: Loop Unrolling Performance (1 point)

What is the throughput of the unrolled loop (Part 2.A) in floating-point operations per cycle (FLOPS/cycle)? Only consider the steady-state behavior of the loop. Do not count memory operations.

Problem 2.D: Software Pipelining Performance (1 point)

What is the throughput of the software-pipelined loop (Part 2.B) in floating-point operations per cycle (FLOPS/cycle)? Only consider the steady-state behavior of the loop. Do not count memory operations.

Problem 3: Multithreading (20 points)

Consider the following code, which performs an in-place slide operation that moves non-zero elements forward in array A by a displacement M. To parallelize the loop on a multithreaded processor, suppose that we split the loop so that each thread executes every iteration for which $(i \% T) == \text{TID}$, where T is the total number of threads and TID is a thread identifier from 0 to T-1 inclusive that is uniquely assigned to each thread.

```
for (i = 0; i < N; i++) {
    if (A[i+M] != 0)
        A[i] = A[i+M];
}
```

N and M are arbitrary integers, and $N > M$ and $N > T$.

The code is executed on a multithreaded in-order core with no data cache, perfect branch prediction with no penalty for both taken and not-taken branches, and no threading overhead.

- Main memory latency is 60 cycles.
- After the processor issues a load, it can continue executing instructions until it reaches an instruction that is dependent on the load value.
- Integer arithmetic operations have a 1-cycle latency.

Problem 3.A (5 points)

For the loop to be safely parallelized this way, what constraint(s) must there be on T, the number of threads? Explain your reasoning.

Assume there is no synchronization among threads while executing the loop.

Problem 3.B (5 points)

Write the assembly code that is executed by each thread. Treat the elements of A as 32-bit integers. You may use any available register, such as t1-t6.

```
# A is passed in a0
# M is passed in a1
# T is passed in a2
# TID is passed in a3
addi t0, a0, N*4    # Initialize loop boundary (t0 = A + N)
slli a1, a1, 2      # Scale M to bytes
slli a2, a2, 2      # Scale T to bytes
slli a3, a3, 2      # Scale TID to bytes
add a0, a0, a3      # Offset pointer by thread ID
loop:
```

Problem 3.C (5 points)

Suppose that threads are switched every cycle using a fixed round-robin schedule. If the thread is not ready to run on its turn, a bubble is inserted into the pipeline.

What is the minimum number of threads needed to always fully utilize the pipeline while maintaining correct execution? Show your work.

Assume that $M=90$ and that N is arbitrarily large.

Problem 3.D (5 points)

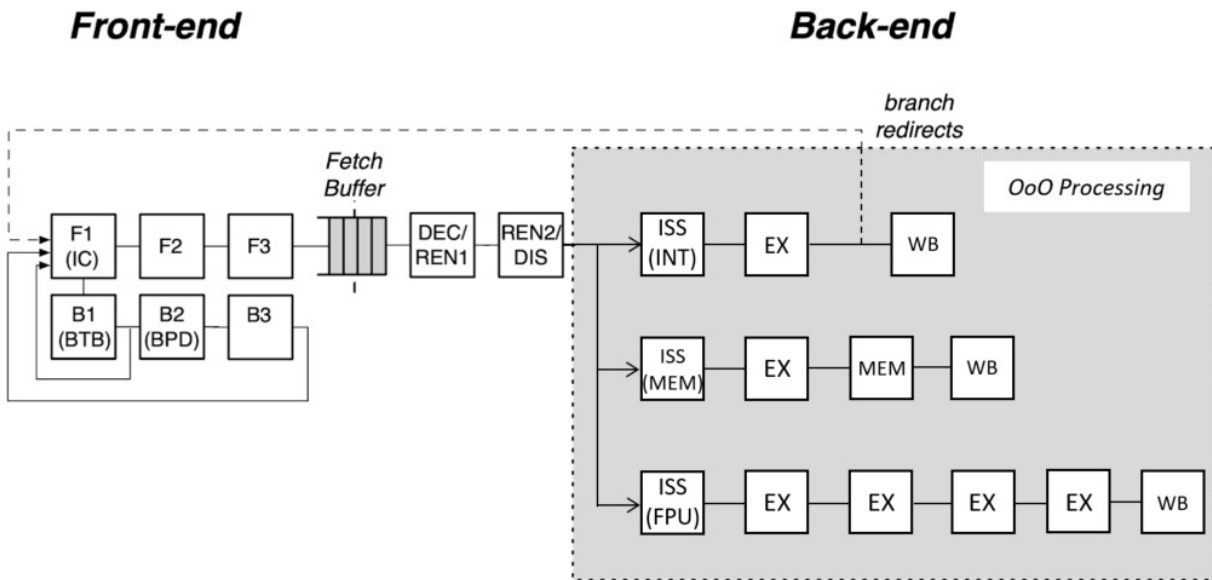
Now consider a dynamic scheduling policy that switches threads only when the next instruction would stall due to a data dependency.

What is the minimum number of threads needed to always fully utilize the pipeline while maintaining correct execution? Show your work.

Assume that $M=90$ and that N is arbitrarily large.

Problem 4: Out-of-Order Execution (20 points)

For this problem, we consider the following dual-issue out-of-order superscalar processor with a unified physical register file.



Dispatch

- **Up to two instructions** can be renamed and dispatched per cycle.
- Register renaming follows the **unified physical register file** scheme.
- Instructions are written into the ROB at the end of the REN2/DIS stage.

Issue

- **There are three issue windows separate from the ROB:**
 - ALU operations and branches (INT)
 - Memory operations (MEM)
 - Float-point operations (FP)
- Instructions are written into the issue window at the end of the REN2/DIS stage.
- **Up to one instruction** can be issued to execution per cycle from each issue window.
- Assume that an age-based scheduler always selects the oldest instruction to issue.
- An instruction may issue in the same cycle when the last operand that it depends on is in the writeback stage.
- Register operands are read during the ISS stage.
- The physical register file has 6 read and 3 write ports. (Assume no structural hazards.)

Execution

- All functional units are fully pipelined with the following latencies:
 - ALU operations: 1 cycle
 - Loads and stores: 2 cycles latency (Assume that all accesses hit in the data cache.)
 - Floating-point operations: 4 cycles
- Writeback occurs in a separate WB stage.
- **Mispredicted branches redirect the frontend and trigger a pipeline flush in the cycle after they are resolved in the INT EX stage.**

Commit

- **Up to two instructions** can be committed per cycle.
- Commit is handled by a decoupled unit that looks at the ROB entries.
- The earliest that an instruction can commit is in the cycle following writeback.

Problem 4.A (10 points)

The following instruction sequence is executed on the out-of-order core described above.

<code>fmul f1, f3, f2</code>	
<code>add x1, x1, x2</code>	
<code>fsw f1, 4(x1)</code>	
<code>lw x2, 0(x3)</code>	
<code>bnez x2, done</code>	Misprediction
<code>flw f1, 0(x1)</code>	Page fault
<code>fadd f3, f1, f2</code>	

Fill out the table with the cycles at which instructions enter the ROB, issue to the functional units, complete and write back to the physical register file, and commit. If an instruction is killed before issuing, completing, or committing, mark the corresponding entries with “-”.

- The ROB is initially empty and contains enough entries for the instructions shown.
- All instructions are present already in the fetch buffer.
- `bnez` is initially predicted to be not taken but later resolves as taken. The “done” branch target points to unrelated code elsewhere.
- A page fault is detected for `flw`.

The first instruction has been done for you.

	Dispatch	Issue	Completion	Commit
<code>fmul</code>	0	1	6	7
<code>add</code>	0			
<code>fsw</code>				
<code>lw</code>				
<code>bnez</code>				
<code>flw</code>				
<code>fadd</code>				

Problem 4.B (10 points)

For the same code as Part 4.A (reproduced below), show the state of the ROB, issue windows, rename table, and free list **in the cycle after recovering from all mispredicts and exceptions** – i.e., immediately after precise architectural state has been restored and the processor has been redirected to the correct branch target. Assume that mispredictions and exceptions use the same rollback procedure, which happens instantaneously.

The first instruction has been done for you.

fmul f1, f3, f2	
add x1, x1, x2	
fsw f1, 4(x1)	
lw x2, 0(x3)	
bnez x2, done	Misprediction
flw f1, 0(x1)	Page fault
fadd f3, f1, f2	

The same assumptions as Part 4.A apply:

- The ROB and issue windows are initially empty.
- bnez is initially predicted to be not taken but later resolves as taken. The “done” branch target points to unrelated code elsewhere.
- A page fault is detected for flw.

For each entry in the rename table, show all changes in sequence. Unused architectural registers are omitted from the rename table for clarity.

The free list is treated as a FIFO, and entries are dequeued from the top and appended to the bottom. Cross out (or mark with “x”) the entries from the free list that have been dequeued.

Rename Table	
x1	P10
x2	P7
x3	P2
f1	P5 → P9
f2	P8
f3	P11

Free List	
P9	x
P4	
P6	
P3	
P1	
P12	

