

## Problem 1: Virtual Memory

**Multiple choice: Check one unless otherwise noted**

A) **(2 Point)** In a virtually indexed, physically tagged cache, what part of the virtual address is used to select the cache set? **Note: a cache may use part of the VPN if aliasing is corrected via another mechanism. However, this is a *pick one* question, so “Page offset” is the best answer. Those who marked both “VPN” and “Page offset” received partial credit.**

☐ VPN      ☐ PPN      ☒ Page offset      ☐ Tag

B) **(2 Point)** A page-table walker (PTW) handles what kind of events?

☐ SegFaults      ☐ Demand paging      ☐ Page faults.      ☒ TLB misses

C) **(3 Points)** Which of the following are advantages of page-based virtual memory over a base-and-bound scheme? **Check all that apply.**

☐ Protecting one process from other processes on the system  
☐ Translating addresses to virtualize the resource of physical memory  
☒ Reduced external fragmentation  
☐ Reduced address translation cost

D) **(3 Points)** Explain why TLBs are critical for good performance in a paged virtual memory system.

In the base case, address translation requires multiple memory accesses to find a physical page number by traversing the page table hierarchy. Not only are these accesses added to every load and store, but this must be performed on every instruction fetch to translate the virtual address of the instruction. A TLB avoids these accesses by caching VPN-to-PPN translations.

For full credit, you must explicitly mention:

- The base case of translation requires *extra memory accesses*
- TLBs avoid this via *caching*

E) **(4 points)** Consider a system with 4096B pages and page-table entries that are 32 bits at all levels of the page table hierarchy. How many levels of page tables are needed to support a 32-bit virtual address space?

Each level of the page table hierarchy must fit in a page. Therefore, there are  $4096 / 4 = 1024$  page table entries per level of the hierarchy. The page offset is  $\log_2(4096) = 12$  bits, and each level of the page table hierarchy uses 10 bits of the VPN as an index.

$$N_{\text{levels}} = (32b - 12b) / 10b = 2$$

F) **(9 points)** Consider a system with the following specifications:

- 8192B page sizes
- 64-bit virtual address space
- 48-bit physical address space
- 32KiB, 4-way set-associative L1 cache with 64B block size
- 256KiB, 4-way set-associative L2 cache with 64B block size
- TLB entries contain 8 bits of metadata

Fill in the table. Show organized work outside the table so partial credit may be assigned.

	# of bits
Physical page number	35
Virtual page number	51
Page offset	13
TLB entry	94
Cache block offset	6
L1 cache tag	35
L1 cache index	7
L2 cache tag	32
L2 cache index	10

- Page offset bits =  $\log_2(8192) = 13$  bits
- PPN bits = (PA bits) – (page offset bits) =  $48 - 13 = 35$  bits
- VPN bits = (VA bits) – (page offset bits) =  $64 - 13 = 51$  bits
- TLB entry bits = (VPN bits) + (PPN bits) + (metadata bits) =  $(51 + 35 + 8) = 94$  bits
- Cache block offset =  $\log_2(64) = 6$  bits
- L1 cache tag = (PA bits) –  $\log_2(32\text{Ki} / 4) = 48 - 13 = 35$  bits
- L1 cache index =  $\log_2(32\text{KiB} / (64\text{B} * 4)) = \log_2(128) = 7$  bits
- L2 cache tag = (PA bits) –  $\log_2(256\text{Ki} / 4) = 48 - 16 = 32$  bits
- L2 cache index =  $\log_2(256\text{KiB} / (64\text{B} * 4)) = \log_2(1024) = 10$  bits

G) **(5 points)** In the system from (F), what would be a reasonable technique(s) to apply to avoid aliasing in the L1 cache while minimizing the overhead of TLB lookups? The L2 cache? Justify your response.

A reasonable technique to avoid aliasing in the L1 cache while minimizing the overhead of TLB lookups would be to make it virtually indexed and physically tagged, and to perform the TLB lookup in parallel with set indexing. This avoids serializing the delay of the TLB check with the full delay of the cache, but will avoid aliasing, since the L1 cache offset and index bits comprise 13 total bits, all of which are taken from the 13-bit page offset.

A reasonable technique to avoid aliasing in the L2 cache while minimizing the overhead of TLB lookups would be to make it physically indexed and physically tagged. Since the physical address is needed by the end of an L1 access to determine whether there has been an L1 miss, it does not add meaningful overhead to provide the PPN bits at the start of any L2 access. With this fully physical addressing scheme, aliasing is impossible.

## Problem 2: Uniprocessor Caches and Memory Hierarchy

- A) (3 Points) How does total capacity, access latency, and bandwidth scale as we move between layers of a uniprocessor's memory hierarchy? Indicate the general trend each with an arrow pointing in the direction of an **increase** in that parameter.

Memory	Example Parameter	Capacity	Access Latency	Bandwidth
Register File	↑	↓	↓	↑
On-chip Caches				
Off-Chip Memory				

- B) (3 Points) Suppose we build a processor with a 1-cycle L1 hit latency, a 10-cycle L1 miss penalty, and a 20-cycle L2 miss penalty. Assuming a L1 hit rate of 90% and a L2 local hit rate of 80%, what is the average memory access time seen by this processor?

$$\begin{aligned}
 \text{AMAT} &= T_{L1, \text{hit}} + (1 - P_{L1, \text{hit}}) * (\text{MissPenalty}_{L1} + (1 - P_{L2, \text{hit}}) * \text{MissPenalty}_{L2}) \\
 \text{AMAT} &= 1 + 0.1 * 10 + 0.1 * 0.2 * 20 \\
 \text{AMAT} &= 2.4
 \end{aligned}$$

2.4 cycles

- C) (4 Points) What is the difference between an inclusive and exclusive multi-level cache? Give one advantage of each approach.

In an inclusive cache hierarchy, an inner cache may only cache lines also cached by adjacent outer level of the cache hierarchy, whereas in an exclusive cache hierarchy lines cached by an inner layer are not cached by the outer layer.

Advantage of exclusion:

- Increased on-chip capacity (@ constant area)

Advantage of inclusion:

- Easier to implement cache coherence

We accepted other well justified answers.

The remainder of this problem evaluates cache performance for different loop orderings. Consider the following two loops, written in C, which calculates the sum of all elements of a 16 by 512 matrix of 32-bit integers:

<i><b>Loop A</b></i>	<i><b>Loop B</b></i>
<pre>int sum = 0; for (i = 0; i &lt; 16; i++)     for (j = 0; j &lt; 512; j++)         sum += A[i][j];</pre>	<pre>int sum = 0; for (j = 0; j &lt; 512; j++)     for (i = 0; i &lt; 16; i++)         sum += A[i][j];</pre>

The matrix A is stored contiguously in memory in row-major order. Row-major order means that elements in the same row of the matrix are adjacent in memory. You may assume A starts at 0x0, thus  $A[i][j]$  resides in memory location  $[4 * (512 * i + j)]$ .

Assume:

- caches are initially empty.
- only accesses to matrix A cause memory references and all other necessary variables are stored in registers.
- instructions are in a separate instruction cache.

**D) (7 points)** Consider a 4KiB direct-mapped data cache with 64-byte cache lines. Calculate the number of cache misses that will occur when running Loop A. Calculate the number of cache misses that will occur when running Loop B. You must show your work for full credit!

Notes:

- to index into a set of this cache we use bits [6:11] of the address.
- $64 / 4 = 16$  words per cache line
- there are  $2^4 * 2^9 = 8192$  total elements  $\rightarrow 2^9$  total cache lines to store the matrix
- each element is accessed only once

For loop A:

- Consecutive elements within a row are 4 bytes apart
- Every 16 elements we experience a cache miss, but then hit on all remaining 15 elements in that line  $\rightarrow$  one miss per cache line  $\rightarrow$  512 misses

For loop B:

- Consecutive elements within a column are  $512 * 4 = 2^{11}$  bytes apart  $\rightarrow$  the MSB of the set address toggles every access; and the tag changes every pair of accesses
- For any given column, all accesses strike two alternating sets:  $\{0,1\}XXXXX \rightarrow$  every access misses  $\rightarrow$  8192 misses

The number of cache misses for Loop A:

512

The number of cache misses for Loop B:

8192

**E) (7 points)** Consider a 4KiB fully-associative data cache with 64-byte cache lines. This data cache uses a first-in/first-out (FIFO) replacement policy. Calculate the number of cache misses that will occur when running Loop A, and when running Loop B. You must show your work for full credit!

Loop A's solution remains unchanged: we fully consume each 64B cache line as we stride along the row.

Loop B:

- This cache has  $2^{12} / 2^6 = 64$  lines, like the last cache
- The array has 16 rows  $< 64$  lines
- Accessing the columns  $j \% 16 == 0$ , will miss on every access, but every accessed line remains in cache (lines brought in 64-49 misses ago will be evicted)
- Columns  $j \% 16 != 0$ , hit on every access  $\rightarrow$  1 miss per cache line  $\rightarrow$  512 misses like loop

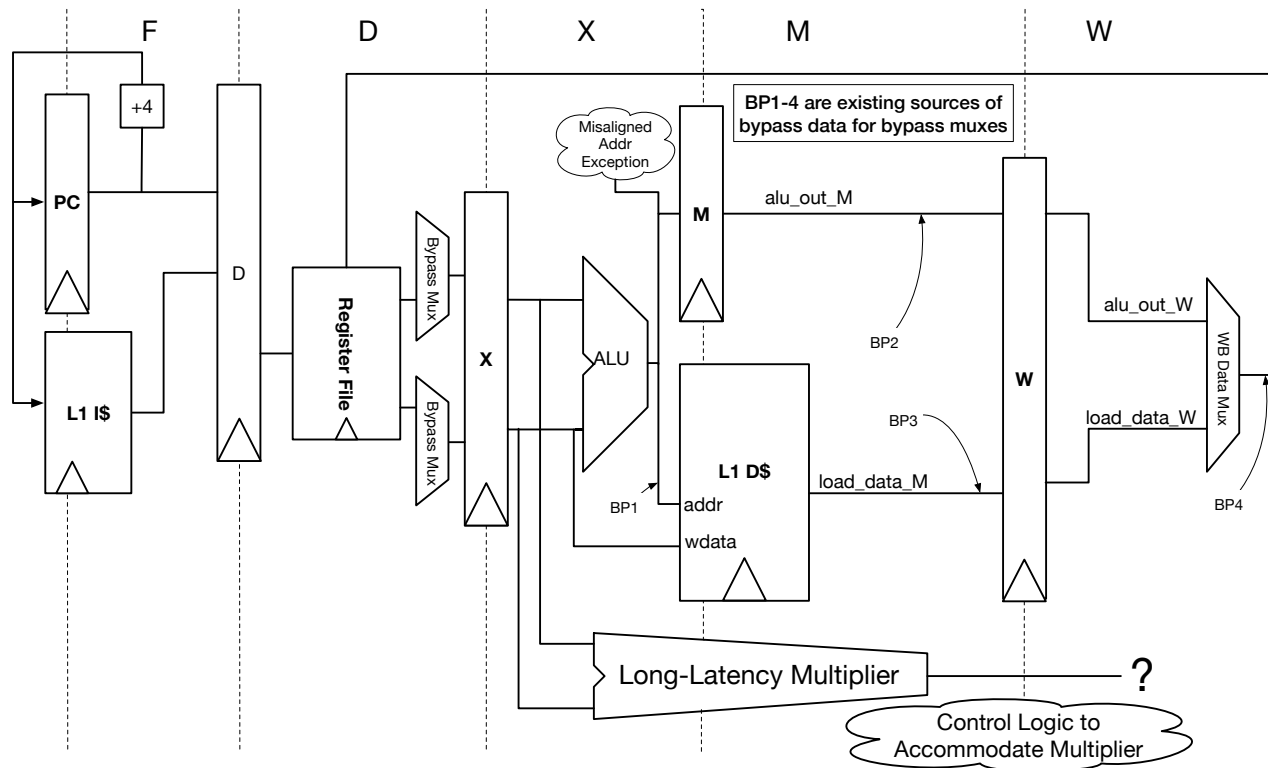
The number of cache misses for Loop A:	512
The number of cache misses for Loop B:	512

### Problem 3: Pipelining and The Iron Law

#### Iron Law (15 points total)

	Instructions / Program	Cycles / Instruction	Seconds / Cycle	Execution Time
Pipelining a single-cycle implementation	<p>Negligible effect / no change</p> <p>This is a purely microarchitectural change</p>	<p>Increase</p> <p>A fully-pipelined implementation will have CPI &gt; 1, since hazards may exist among instructions in different stages.</p>	<p>Decrease</p> <p>Pipelining inserts registers in the middle of the datapath, decreasing the number of serial levels of logic that data must propagate through in one cycle.</p>	<p>Almost certainly decrease</p> <p>Almost any design can benefit from some degree of pipelining. In other words, one pipeline stage is rarely the optimal number for performance.</p>
Pipelining an unpipelined multi-cycle implementation, while keeping the <i>latency</i> of each instruction constant	<p>Negligible effect / no change</p> <p>This is a purely microarchitectural change</p>	<p>Decrease</p> <p>An unpipelined implementation cannot exploit ILP and waits for each operation to complete before beginning the next.</p>	<p>Negligible effect / no change</p> <p>A pipelined implementation may be a bit more complex, but well-designed stalls or replay logic will avoid deep serial paths.</p>	<p>Decrease</p> <p>Pipelining will generally help the performance of multi-cycle processing pipelines, at the expense of increased complexity and area.</p>
Reducing the number of stages in a pipeline	<p>Negligible effect / no change</p> <p>This is a purely microarchitectural change</p>	<p>Decrease</p> <p>Shallower pipelines will be exposed to fewer hazards, as they have fewer overlapped instructions that may have dependencies.</p>	<p>Increase</p> <p>Merging or lengthening stages will cause more work to be done in a single clock cycle, and the merged logic in the datapath will often be serial.</p>	<p>Ambiguous</p> <p>This is highly dependent on both the baseline design and workload. The CPI decrease may be relatively smaller or larger than the clock period increase.</p>

## Modifying the 5-stage Pipeline: Long Latencies



**A) (2 Point)** If we allow non-dependent operations to proceed while a multiply is outstanding, which parts of execution are proceeding “out-of-order” in the processor? Assume that multiply operations cannot generate exceptions based on their input or output values. **Check all that apply.**

☐ Issue      ☒ Completion      ☐ Commit      ☐ None of these

**B) (3 Points)** Suppose we want to add a multiplier with an 8-cycle latency and an 8-cycle occupancy to a 5-stage pipeline, as shown above. We must keep around some information about registers whose values are pending an outstanding multiply operation. What type of microarchitectural structure is commonly used for this purpose?

Scoreboard



## Problem 4: Cache Coherence

A) **(2 Point)** What is the fewest number of metadata bits (i.e., not including the line's tag and data) a writeback cache must track per cache line to implement an MSI coherence protocol (assume no transient states).

☐ 1 bit      ☒ 2 bits      ☐ 3 bits      ☐ 4 bits

B) **(3 Point)** Which sequence of memory operations would produce less coherency traffic under MESI vs MSI. Assume all memory operations act on the same cache line, and that neither processor P1 nor P2 have the line in cache initially. **Check all that apply.**

☐ P1.read, P2.read, P1.write  
☒ P1.read, P1.write, P2.read  
☐ P1.write, P2.read, P1.read

C) **(3 Points)** In general, which of the following are advantages of snoop cache coherence over directory-based cache coherence. **Check all that apply.**

☒ Simpler to implement  
☐ Scalable to more cores  
☒ Lower latency on cache misses  
☐ Uses less interconnect bandwidth

D) **(4 Points)** In the table below, indicate which memory operations experience a hit, true sharing miss, or false sharing miss under an MSI protocol. Assume x1, x2 are in the same cache line and the line is initially cached in a shared state by both processors (P1 and P2). The first row is filled out for you.

Time	P1	P2	Hit	True Sharing Miss	False Sharing Miss
1		Write x2		<b>X</b>	
2	Read x1				<b>X</b>
3		Read x1	<b>X</b>		
4	Write x2			<b>X</b>	
5		Write x1		<b>X</b>	

In the remainder of this question, we'll study a simplified version of the directory-based cache coherence scheme from HW5 (**detach Appendix B, attached to the reverse of the exam**). Our system consists of 16 cores, each with write-back, write-allocate private caches. All caches are connected to a single directory. Specifically, we're interested the series of coherence events that transpire to service a CPU's load or store under different initial conditions. For example, consider a load-miss in CPU0 to a line is currently not cached by any CPU in the system:

Time	Agent	Current State	Event/Message Received	Next State	Message(s) Sent
0	CPU 0	C-nothing	Load	C-pending	ShReq(0)
1	Directory	R( $\epsilon$ ) ( <i>empty</i> )	ShReq(0)	R(0)	ShResp(0)
2	CPU 0	C-pending	ShResp(0)	C-shared	N/A
Total Messages Sent:					2

Simplifications:

- assume agents can send and receive multiple messages simultaneously
- assume messages take one time step to reach their destinations
- **if a set of events *may* happen in parallel, indicate this by setting the “time” field to the same value.**
- include only the ID field of messages (i.e., neglect data, address fields)

E) **(6 Points)** In the table below, indicate the series of events that occur to service CPU0, load-miss when CPU 4 has the line cached in the C-exclusive state.

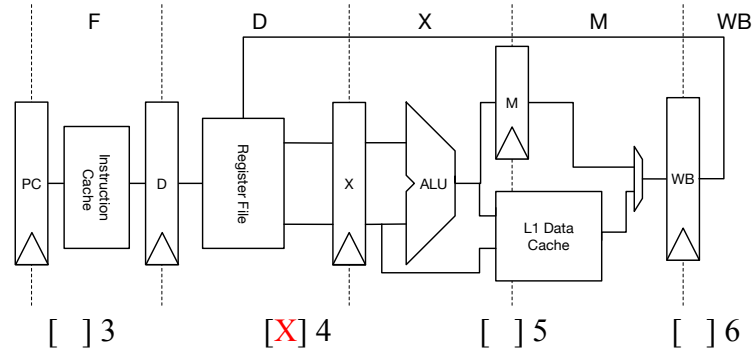
Time	Agent	Current State	Event/Message Received	Next State	Message(s) Sent
0	CPU 0	C-nothing	Load	C-pending	ShReq(0)
1	Directory	W(4)	ShReq(0)	$T_w(0,4)$	WbReq(4)
2	CPU4	C-exclusive	WbReq(4)	C-shared	WbResp(4)
3	Directory	$T_w(0,4)$	WbResp(4)	R({0,4})	ShResp(0)
4	CPU 0	C-pending	ShResp(0)	C-shared	N/A
Total Messages Sent:					4

F) **(10 Points)** In the table below, indicate the series of events that occur to service CPU0's store-miss when CPUs 0, 4 and 8 have the line cached in a shared state.

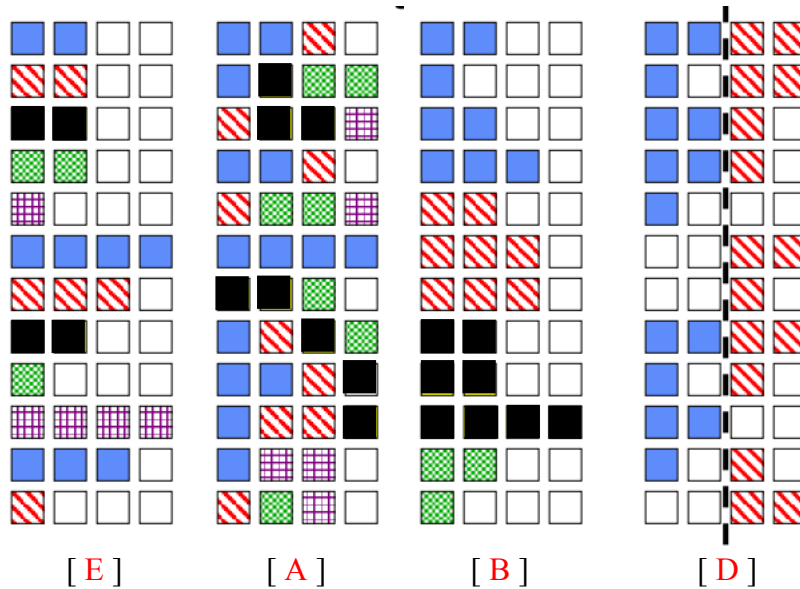
Time	Agent	Current State	Event/Message(s) Received	Next State	Message(s) Sent
0	CPU 0	C-Shared	Store	C-pending	ExReq(0)
1	Directory	$R(\{0,4,8\})$	ExReq(0)	$T_R(0, \{4,8\})$	InvReq(4);InvReq(8)
2	CPU4	C-Shared	InvReq(4)	C-nothing	InvResp(4)
2	CPU8	C-Shared	InvReq(8)	C-nothing	InvResp(8)
3	Directory	$T_R(0, \{4,8\})$	InvResp(4); InvResp(8)	W(0)	ExResp(0)
4	CPU0	C-pending	ExResp(0)	C-exclusive	N/A
Total Messages Sent:					6

## Question 5: Multithreading

- A) (2 Points) If we multithread a classic 5-stage RISC pipeline with no bypassing (the result of an instruction can be read in decode one cycle after writeback) using a fixed-interleave policy, how many threads are required to avoid interlocks (assume no other hazards)?



- B) (4 Points) The following diagrams indicate how the slots of a four execution units are being utilized, with each row corresponding to a different cycle and each column corresponding to a functional unit slot of the machine. Instructions belonging to different threads are indicated with a different color and texture. A white square indicates an empty slot. Label each illustration with letter (A – G) of the corresponding multithreading scheme.



### Labels

- A: Simultaneous MT
- B: Coarse-grained MT
- C: Vertical MT
- D: Multiprocessing
- E: Fine-grained MT
- F: Horizontal MT
- G: Parallel MT

- C) (3 Points) Which of the following **must** be duplicated per-thread in a multithreaded processor. **Check all that apply.**

- ☒ GPRs
- ☐ TLBs
- ☒ PCs
- ☐ Branch Predictors
- ☐ Reorder Buffers

In remainder of this question, we will consider the execution of the following C kernel:

```
void kernel(float * A, float * B, float * C, int N) {
    for (int i = 0; i < N; i++)
        C[i] += A[i] * B[i];
}
```

The code above can be translated into the following RISC-V assembly code:

```
# a1, a2, a3 hold A, B, and C respectively
# a4 holds N
# t0 is initially 0
loop:
    flw f1, 0(a1)
    flw f2, 0(a2)
    flw f3, 0(a3)
    fmul f4, f1, f2
    fadd f5, f4, f3
    fst f5, 0(a3)
    addi a1, a1, 4
    addi a2, a2, 4
    addi a3, a3, 4
    addi t0, t0, 1
    bne t0, a4, loop
```

Each cycle, the processor can fetch and issue one instruction that performs any of the following operations:

- load/store, 20-cycle latency (fully pipelined)
- integer add, 1-cycle latency
- floating-point add, 1-cycle latency
- floating-point multiply, 5-cycles latency (fully pipelined)
- branch, no delay slots, 1-cycle latency

The processor does not have a cache. Each memory operation directly accesses main memory. If an instruction cannot be issued due to a data dependency, the processor stalls. We also assume that the processor has a perfect branch predictor with no penalty for both taken and not-taken branches. Assume N is very large.

- D) (5 points) Consider a single-issue in-order, multithreaded pipeline, where 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. Each thread executes the above algorithm, and is calculating its own independent piece of the A array (i.e., there is no communication required between threads). In steady state, how many cycles does the machine take to execute each loop iteration for a very large value of N, without rescheduling (changing) the assembly and **assuming only a single thread**. Explain.

#	Instruction	Issue Cycle	Complete Cycle
0	flw f2, 0(a1)	0	20
1	flw f2, 0(a2)	1	21
2	flw f3, 0(a3)	2	22
3	fmul f4, f1, f2	21	26
4	fadd f5, f4, f3	26	27
5	fst f5, 0(a3)	27	47
6	addi a1, a1, 4	28	29
7	addi a2, a2, 4	29	30
8	addi a3, a3, 4	30	31
9	addi t0, t0, 1	31	32
10	bne t0, a4, loop	32	33
	flw f2, 0(a1)	33	34

33 cycles

- E) (4 points) Without rescheduling (changing) the assembly, how many threads are required to saturate the machine from part D? Explain.

Memory operations are the longest latency operations at 20 cycles: we'll need to hide that latency to fully utilize the machine. The fmul (3) and fadd (4) are the only two instructions that depend on loads, each has an intermediate instruction (2) and (3) between them respectively. As such we'll need N threads, where:

$$N = \lceil \text{latency} / (\# \text{ of intermediate instructions} + 1) \rceil = \lceil 20 / (2) \rceil = 10$$

You can also check the fadd dependency on fmul isn't problematic with the same expression.

$$N = \lceil \text{latency} / (\# \text{ of intermediate instructions} + 1) \rceil = \lceil 5 / 1 \rceil = 5 (5 < 10)$$

10 threads

**F) (12 points)** What is the **minimum** number of threads we need to achieve peak performance on the machine from part D, assuming you **may** reschedule the code as necessary without loop unrolling. Explain, giving your final schedule.

**You may perform the following optimizations:**

- **Reordering instructions**
- **Renaming / re-allocating registers**
- **Changing load/store immediate offsets**

Without tricks like software pipelining or loop unrolling which we (tried to) forbid, the objective here is to stuff as many existing intermediate operations between long latency operations and their dependent instructions as possible. We're looking for a schedule that minimizes the largest  $N_{i,j}$ , where  $i$  and  $j$  are indexes of dependent instructions in program order ( $i$  depends on  $j$ ):

$$N_{i,j} = \text{latency}_j / (i - j)$$

```
loop:
0 flw f1, 0(a1)
1 flw f2, 0(a2)
2 flw f3, 0(a3)
3 addi a1, a1, 4
4 addi a2, a2, 4
5 addi a3, a3, 4
6 fmul f4, f1, f2
7 addi t0, t0, 1 # many students forgot to put an insn here
8 fadd f5, f4, f3
9 fst f5, -4(a3)
10 bne t0, a4, loop
```

There are four long-latency ( $> 1$  cycle) operations we need to consider:

fmul  $\rightarrow$  flw (0):  $N_{6,0} = \lceil 20 / (6 - 0) \rceil = 4$

fmul  $\rightarrow$  flw (1):  $N_{6,1} = \lceil 20 / (6 - 1) \rceil = 4$

fadd  $\rightarrow$  flw (2):  $N_{8,2} = \lceil 20 / (8 - 2) \rceil = 4$

fadd  $\rightarrow$  fmul:  $N_{8,6} = \lceil 5 / (8 - 6) \rceil = 3$

This schedule requires 4 threads. Some crafty students tried software pipelining: a schedule requiring only three threads received full credit (an extra instruction would have made it possible to do it in two.)

<b>4</b> threads
------------------

## Problem 6 7: Memory Consistency Models (20 points)

- A) (3 Points) Is it possible to translate code that assumes sequential consistency to the RISC-V weak memory consistency model? Explain.

Yes, by inserting memory fences where violations of SC may be visible to another thread.  
Existence proof: for example, one foolproof way to achieve this is to insert a fence<sub>rw,rw</sub> between every memory operation in each thread.

- B) (3 Points) Given the instruction sequences below, check all possible combinations of P1.r2, P2.r2 after both threads have executed, assuming an ISA that is sequentially consistent. X and Y are non-overlapping, and initially X = 0, Y = 0,

**P1:**

```
li r1, 1
st r1, X
lw r2, Y
```

**P2:**

```
li r1, 2
st r1, Y
ld r2, X
```

- [ ] P1.r2 = 0; P2.r2 = 0
- [X] P1.r2 = 0; P2.r2 = 1
- [X] P1.r2 = 2; P2.r2 = 0
- [X] P1.r2 = 2; P2.r2 = 1

- C) (2 Points) Given the same instruction sequences and initial conditions from part B, which of the following combinations are possible under an ISA with TSO memory consistency model.

- [X] P1.r2 = 0; P2.r2 = 0
- [X] P1.r2 = 0; P2.r2 = 1
- [X] P1.r2 = 2; P2.r2 = 0
- [X] P1.r2 = 2; P2.r2 = 1

- D) (4 Points) In general, what is the difference between a weak and a strong memory consistency model? Give one advantage of each.

Weak memory models relax some combination of the program-order and/or write-atomicity constraints imposed by stronger memory models.

Advantage of a strong memory model:

- More intuitive memory semantics make it easier to write and debug concurrent code
- Similarly, easier to write correct compilers for high-level languages

Advantage of a weak memory model:

- Easier to build simple implementations (more design flexibility, many simple optimizations would violate a strong MCM without considerably more complexity)
- Easier to build high-performance implementations, as the implementation can aggressively reorder memory ops without needing to speculate on the MCM



**(10 points)** The following RISC-V assembly encodes a request-response relationship between two threads. The requestor thread puts work in a memory location `request` and then sets `go = 1`. It then spins waiting on the responder to produce the response. The responder thread spins waiting for work from the master, by checking if `go` has been set. Once available, it reads the request data, computes the response, and writes the response data to a memory location `response` before setting `done = 1`.

Requestor thread:

```
li a0, 1
sw data, request
sw a0, go
spin:
lw a1, done
beqz a1, spin
lw a2, response
sw zero, done
```

Responder thread:

```
spin:
lw a1, go
beq a1, spin
lw a2, request
sw zero, go
... a3 = process request
sw a3, response
li a0, 1
sw a0, done
```

Under a fully relaxed memory consistency model, insert fences where necessary to ensure this code functions correctly? Use the least restrictive fences for full credit. **Assume each thread executes its code only once.**

```
li a0, 1

sw data, request

fencew,w

sw a0, go

spin:

lw a1, done

beqz a1, spin

fencer,r

lw a2, response

sw zero, done
```

```
spin:

lw a1, go

beq a1, spin

fencer,r

lw a2, request

sw zero, go

... a3 = process request

sw a3, response

fencew,w

li a0, 1

sw a0, done
```

## Problem 8 7: Vector Machines

A) (2 Point) What is the minimum number of lanes a machine must have to exploit data-level parallelism on vector instructions when VL is set to 4?

☒ 1      ☐ 2      ☐ 4      ☐ 8

Even single-lane vector machines may exploit DLP by pipelining vector operations in a manner that takes advantage of the lack of inter-element dependencies.

B) (16 points) Vectorize the following **32-bit integer** C code using the RISC-V vector specification described in lab 4. See appendix A for the vector instruction set listing.

```
for (i = 0; i < N; i++) {  
    B[i] = (A[i] < 0) ? -A[i] : A[i]; // A and B do not overlap  
}
```

```
# v0, v2-v8: configured to hold 32-bit integer values  
# v1: configured to hold 8-bit integers, used as mask register  
# a0 and a1: hold pointers A and B, respectively  
# a2: holds N.  
# t0-t3 may be used as scalar temporaries
```

```
# Your code begins:
```

```
zero:  
    setvl    t0, a2  
    vsub     v0, v0, v0    # v0[i] = 0
```

```
stripmine_loop:  
    setvl    t0, a2  
    slli     t1, t0, 0x2  
    vld      v2, 0(a0)  
    vslti    v1, v2, v0  
    vsub     v2, v0, v2, v1t  
    vst      v2, 0(a1)  
    add      a0, a0, t1  
    add      a1, a1, t1  
    sub      a2, a2, t0  
    bne      a2, r0, stripmine_loop
```

```
done:  
    ret
```

## Appendix A: Vector Architecture for Question 4 7

This instruction listing is identical to that in Lab 4, but with a `setvl` instruction that has identical semantics to the preprocessor macro provided in Lab 4. This instruction first sets VL to  $\min(\text{maximum vector length}, rs1)$ ; and then returns the new VL.

- The two vector mask arguments are `vl t` (perform op only if `vl[i]` is true) and `vl f` (perform op only if `vl[i]` is false). Omitting the final vector mask (`vm`) argument to all instructions is legal, and treats all elements  $i < VL$  as active.
- Unlike shortening VL, lengthening VL causes the elements extending past the previous vector length to have undefined values.

Instruction	Operation
<code>setvl rd, rs1</code>	$VL := \min(MVL, rs1); (rd) := VL$
<code>vld vd, offset(rs1), vm</code>	$vd[i] := \text{mem}[(rs1) + \text{offset} + i]$
<code>vst vs3, offset(rs1), vm</code>	$\text{mem}[(rs1) + \text{offset} + i] := vs3[i]$
<code>vlds vd, offset(rs1), rs2, vm</code>	$vd[i] := \text{mem}[(rs1) + \text{offset} + i * rs2]$
<code>vsts vs3, offset(rs1), rs2, vm</code>	$\text{mem}[rs1 + \text{offset} + i * (rs2)] := vs3[i]$
<code>vldx vd, offset(rs1), vs2, vm</code>	$vd[i] := \text{mem}[(rs1) + \text{offset} + vs2[i]]$
<code>vstx vs3, offset(rs1), vs2, vm</code>	$\text{mem}[(rs1) + \text{offset} + vs2[i]] := vs3[i]$
<code>vadd vd, vs1, vs2, vm</code>	$vd[i] := vs1[i] + vs2[i]$
<code>vsub vd, vs1, vs2, vm</code>	$vd[i] := vs1[i] - vs2[i]$
<code>vmul vd, vs1, vs2, vm</code>	$vd[i] := vs1[i] * vs2[i]$
<code>vdiv vd, vs1, vs2, vm</code>	$vd[i] := vs1[i] / vs2[i]$
<code>vrem vd, vs1, vs2, vm</code>	$vd[i] := vs1[i] \% vs2[i]$
<code>vmax vd, vs1, vs2, vm</code>	$vd[i] := \max(vs1[i], vs2[i])$
<code>vmin vd, vs1, vs2, vm</code>	$vd[i] := \min(vs1[i], vs2[i])$
<code>vsll vd, vs1, vs2, vm</code>	$vd[i] := vs1[i] \ll vs2[i]$
<code>vsrl vd, vs1, vs2, vm</code>	$vd[i] := vs1[i] \gg vs2[i]$
<code>vseq vd, vs1, vs2, vm</code>	$vd[i] := vs1[i] == vs2[i] ? 1 : 0$
<code>vsne vd, vs1, vs2, vm</code>	$vd[i] := vs1[i] != vs2[i] ? 1 : 0$
<code>vsllt vd, vs1, vs2, vm</code>	$vd[i] := vs1[i] < vs2[i] ? 1 : 0$
<code>vsge vd, vs1, vs2, vm</code>	$vd[i] := vs1[i] \geq vs2[i] ? 1 : 0$
<code>vaddi vd, vs1, imm, vm</code>	$vd[i] := vs1[i] + \text{imm}$
<code>vslli vd, vs1, imm, vm</code>	$vd[i] := vs1[i] \ll \text{imm}$
<code>vsrli vd, vs1, imm, vm</code>	$vd[i] := vs1[i] \gg \text{imm}$
<code>vmadd vd, vs1, vs2, vs3, vm</code>	$vd[i] := vs1[i] * vs2[i] + vs3[i]$
<code>vmsub vd, vs1, vs2, vs3, vm</code>	$vd[i] := vs1[i] * vs2[i] - vs3[i]$
<code>vnmadd vd, vs1, vs2, vs3, vm</code>	$vd[i] := -(vs1[i] * vs2[i] + vs3[i])$
<code>vnmsub vd, vs1, vs2, vs3, vm</code>	$vd[i] := -(vs1[i] * vs2[i] - vs3[i])$
<code>vslide vd, vs1, rs2, vm</code>	$vd[i] := 0 \leq (rs2) + i < VL ? vs1[(rs2) + i] : 0$
<code>vinserl vd, vs1, rs2, vm</code>	$vd[(rs2)] := (rs1)$
<code>vextract rd, vs1, rs2, vm</code>	$(rd) := vs1[(rs2)]$
<code>vmfirst rd, vs1</code>	$(rd) := ([i \text{ for } i \text{ in range}(0, VL) \text{ if } \text{LSB}(vs1[i]) == 1] + [-1])[0]$
<code>vmpop rd, vs1</code>	$(rd) := \text{len}([i \text{ for } i \text{ in range}(0, VL) \text{ if } \text{LSB}(vs1[i]) == 1])$
<code>vselect vd, vs1, vs2, vm</code>	$vd[i] := vs2[i] < VL ? vs1[vs2[i]] : 0$
<code>vmerge vd, vs1, vs2, vm</code>	$vd[i] := \text{LSB}(vm[i]) ? vs2[i] : vs1[i]$

## Appendix B: Directory-based Cache Coherence Protocol (Abridged Handout 6)

### Changes from the handout:

- Rep renamed Resp (Rep will still be accepted)
- **Tr and Tw now include the site id ( $id_{req}$ ) that initialized the request**
- Removed unnecessary messages, columns, and rows for the exam
- Removed home sites – there is only one directory site

**Cache states:** For each cache line, there are 4 possible states:

- C-nothing: The accessed data is not resident in the cache.
- C-shared: The accessed data is resident in the cache, and possibly also cached at other sites. The data in memory is valid.
- C-exclusive: The accessed data is exclusively resident in this cache and has been modified.
- C-pending: The accessed data is in a *transient* state.

**Directory states:** For each memory block, there are 4 possible states:

- $R(dir)$ : The memory block is shared by the sites specified in  $dir$  ( $dir$  is a set of sites). The data in memory is valid in this state. If  $dir$  is empty (i.e.,  $dir = \epsilon$ ), the memory block is not cached by any site.
- $W(id)$ : The memory block is exclusively cached at site  $id$ , and has been modified at that site. Memory does not have the most up-to-date data.
- $T_R(id_{req}, dir)$ : The memory block is in a transient state waiting for the acknowledgements to the invalidation requests that the directory has issued, before giving exclusive access to the site  $id_{req}$  that initiated the request.
- $T_W(id_{req}, id)$ : The memory block is in a transient state waiting for a block exclusively cached at site  $id$  (i.e., in C-modified state) to make the memory block at the directory up-to-date, before servicing initial request made by site  $id_{req}$ .

### Protocol messages:

Category	Messages
Cache to Memory Requests	ShReq, ExReq
Memory to Cache Requests	WbReq, InvReq
Cache to Memory Responses	WbResp( $v$ ), InvResp
Memory to Cache Responses	ShResp( $v$ ), ExResp( $v$ )

Current State	Handling Message	Next State	Action
C-nothing	Load	C-pending	ShReq(id)
C-nothing	Store	C-pending	ExReq(id)
C-nothing	ShResp (a)	C-shared	updates cache with prefetch data
C-nothing	ExResp (a)	C-exclusive	updates cache with data
C-shared	Store	C-pending	ExReq(id)
C-shared	InvReq(a)	C-nothing	InvResp(id)
C-exclusive	WbReq(a)	C-shared	WbResp(id, data(a))
C-pending	ShResp(a)	C-shared	updates cache with data
C-pending	ExResp(a)	C-exclusive	update cache with data

Table B-1: Cache State Transitions

Current State	Message Received	Next State	Action
R(dir) & (dir = $\epsilon$ )	ShReq(a)	R({id})	ShResp(id, data(a))
	ExReq(a)	W(id)	ExResp(id, data(a))
R(dir) & (id $\notin$ dir) & (dir $\neq \epsilon$ )	ShReq(a)	R(dir + {id})	ShResp(id, data(a))
	ExReq(a)	Tr(id, dir)	InvReq(dir, a) // Note: here id is also = id <sub>req</sub>
R(dir) & (dir = {id})	ExReq(a)	W(id)	ExResp(id, data(a))
R(dir) & (id $\in$ dir) & (dir $\neq$ {id})	ExReq(a)	Tr(id, dir - {id})	InvReq(dir - {id}, a) // Note: here id is also = id <sub>req</sub>
W(id') (id' $\neq$ id)	ShReq(a)	Tw(id, id')	WbReq(id', a) // Note: here id is also = id <sub>req</sub>
Tr(id <sub>req</sub> , dir) & (id $\in$ dir) & dir $\neq$ {id}	InvResp(a)	Tr(id <sub>req</sub> , dir - {id})	None
Tr(id <sub>req</sub> , dir) & (dir = {id})	InvResp(a)	W(id <sub>req</sub> )	ExResp(id <sub>req</sub> , data(a))
Tw(id <sub>req</sub> , id)	WbResp(a)	R({id <sub>req</sub> , id})	data-> memory; ShResp(id <sub>req</sub> )

Table B-2: Directory State Transitions, messages sent from site id.