

Problem M11.1: Synchronization Primitives

One of the common instruction sequences used for synchronizing several processors are the LOAD RESERVE/STORE CONDITIONAL pair (from now on referred to as LdR/StC pair). The LdR instruction reads a value from the specified address and sets a local reservation for the address. The StC attempts to write to the specified address provided the local reservation for the address is still held. If the reservation has been cleared the StC fails and informs the CPU.

Problem M11.1.A

Describe under what events the local reservation for an address is cleared.

Problem M11.1.B

Is it possible to implement LdR/StC pair in such a way that the memory bus is not affected, i.e., unaware of the addition of these new instructions? Explain

Problem M11.1.C

Give two reasons why the LdR/StC pair of instructions is preferable over atomic read-test-modify instructions such as the TEST&SET instruction.

Problem M11.1.D

LdR/StC pair of instructions were conceived in the context of snoopy busses. Do these instructions make sense in our directory-based system in Handout #13? Do they still offer an advantage over atomic read-test-modify instructions in a directory-based system? Please explain.

Problem M11.2: Implementing Directories

Ben Bitdiddle is implementing a directory-based cache coherence invalidate protocol for a 64-processor system. He first builds a smaller prototype with only 4 processors to test out the cache coherence protocol described in Handout #13. To implement the list of sharers, **S**, kept by **home**, he maintains a bit vector per cache block to keep track of all the sharers. The bit vector has one bit corresponding to each processor in the system. The bit is set to one if the processor is caching a shared copy of the block, and zero if the processor does not have a copy of the block. For example, if Processors 0 and 3 are caching a shared copy of some data, the corresponding bit vector would be 1001.

Problem M11.2.A

The bit vector worked well for the 4-processor prototype, but when building the actual 64-processor system, Ben discovered that he did not have enough hardware resources. Assume each **cache block is 32 bytes**. What is the overhead of maintaining the sharing bit vector for a 4-processor system, as a **fraction of data storage bits**? What is the overhead for a 64-processor system, as a **fraction of data storage bits**?

Overhead for a 4-processor system: _____

Overhead for a 64-processor system: _____

Problem M11.2.B

Since Ben does not have the resources to keep track of all potential sharers in the 64-processor system, he decides to limit **S** to keep track of only 1 processor using its 6-bit ID as shown in Figure M11.2-A (**single-sharer scheme**). When there is a load $[C2P_Req(a) S]$ request for a shared cache block, Ben invalidates the existing sharer to make room for the new sharer (home sends an invalidate request $[P2C_Req(a) I]$ to the existing sharer, the existing sharer sends an invalidate response $[C2P_Rep(a) I]$ to home, home replaces the exiting sharer's ID with the new sharer's ID and sends the load response $[P2C_Rep(a) I S]$ to the new sharer).

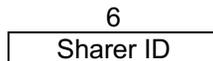


Figure M11.2-A

Consider a 64-processor system. To determine the efficiency of the bit-vector scheme and single-sharer scheme, **fill in the number of invalidate-requests** that are generated by the protocols for each step in the following two sequences of events. Assume cache block **B** is uncached initially for both sequences.

Sequence 1	bit-vector scheme # of invalidate-requests	single-sharer scheme # of invalidate-requests
Processor #0 reads B	0	0
Processor #1 reads B		
Processor #0 reads B		

Sequence 2	bit-vector scheme # of invalidate-requests	single-sharer scheme # of invalidate-requests
Processor #0 reads B	0	0
Processor #1 reads B		
Processor #2 writes B		

Problem M11.2.C

Ben thinks that he can improve his original scheme by adding an extra “**global bit**” to **S** as shown in Figure M11.2-B (**global-bit scheme**). The global bit is set when there is more than 1 processor sharing the data, and zero otherwise.

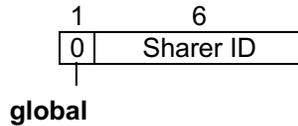


Figure M11.2-B

When the global bit is set, home stops keeping track of a specific sharer and assumes that all processors are potential sharers.

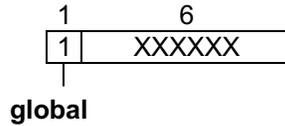


Figure M11.2-C

Consider a 64-processor system. To determine the efficiency of the global-bit scheme, **fill in the number of invalidate-requests** that are generated for each step in the following two sequences of events. Assume cache block **B** is uncached initially for both sequences.

Sequence 1	global-bit scheme # of invalidate-requests
Processor #0 reads B	0
Processor #1 reads B	
Processor #0 reads B	

Sequence 2	global-bit scheme # of invalidate-requests
Processor #0 reads B	0
Processor #1 reads B	
Processor #2 writes B	

Problem M11.3: Tracing the Directory-based Protocol

For the problem we will be using the following sequences of instructions. These are small programs, each executed on a different processor, each with its own cache and register set. In the following **R** is a register and **X** is a memory location. Each instruction has been named (e.g., B3) to make it easy to write answers.

Assume data in location X is initially 0.

Processor A	Processor B	Processor C
A1: ST X, 1	B1: R := LD X	C1: ST X, 6
A2: R := LD X	B2: R := ADD R, 1	C2: R := LD X
A3: R := ADD R, R	B3: ST X, R	C3: R := ADD R, R
A4: ST X, R	B4: R:= LD X	C4: ST X, R
	B5: R := ADD R, R	
	B6: ST X, R	

These questions relate to the directory-based protocol in Handout #13 (as well as Lecture 15). Unless specified otherwise, assume all caches are initially empty and *no voluntary responses are sent* (i.e. responses are sent only on receiving a request).

Problem M11.3.A

Suppose we execute Program A, followed by Program B, followed by Program C and all caches are initially empty. Write down the sequence of messages that will be generated. We have omitted ADD instructions because they cannot generate any messages. EO indicates the global execution order.

Processor A			Processor B			Processor C		
Ins	EO	Messages	Ins	EO	Messages	Ins	EO	Messages
A1	1	<M,A,Req,x,M> <A,M,Rep,x,I,M,0>	B1	4		C1	8	
A2	2		B3	5		C2	9	
A4	3		B4	6		C4	10	
			B6	7				

How many messages are generated? _____

Problem M11.3.B

Is there an execution sequence that will generate even fewer messages? Fill in the EO columns to indicate the global execution order. Also, fill in the messages.

Processor A			Processor B			Processor C		
Ins	EO	Messages	Ins	EO	Messages	Ins	EO	Messages
A1			B1			C1		
A2			B3			C2		
A4			B4			C4		
			B6					

How many messages are generated? _____

Problem M11.3.C

Can the number of messages in Problem M11.3.B be decreased *by using voluntary responses*? Explain.

Problem M11.3.D

What is the execution sequence that generates the most messages *without any voluntary responses*? Fill in the global execution order (EO) and the messages generated. Partial credit will be given for identifying a bad, but not necessarily the worst sequence.

Processor A			Processor B			Processor C		
Ins	EO	Messages	Ins	EO	Messages	Ins	EO	Messages
A1			B1			C1		
A2			B3			C2		
A4			B4			C4		
			B6					

How many messages are generated? _____

Problem M11.4: Snoopy Cache Coherent Shared Memory

In this problem, we investigate the operation of the snoopy cache coherence protocol in Handout #14.

The following questions are to help you check your understanding of the coherence protocol.

- Explain the differences between **CR**, **CI**, and **CRI** in terms of their purpose, usage, and the actions that must be taken by memory and by the different caches involved.
- Explain why **WR** is not snooped on the bus.
- Explain the I/O coherence problem that **CWI** helps avoid.

Problem M11.4.A

Where in the Memory System is the Current Value

In Table M11.4-1, M11.4-2, and M11.4-3, column 1 indicates the initial state of a certain address X in a cache. Column 2 indicates whether address X is currently cached in any other cache. (The “cached” information is known to the cache controller only immediately following a bus transaction. Thus, the action taken by the cache controller must be independent of this signal, but state transition could depend on this knowledge.) Column 3 enumerates all the available operations on address X, either issued by the CPU (read, write), snooped on the bus (**CR**, **CRI**, **CI**, etc), or initiated by the cache itself (replacement). Some state-operation combinations are impossible; you should mark them as such. (See the first table for examples). In columns 6, 7, and 8 (corresponding to this cache, other caches and memory, respectively), **check all possible locations where up-to-date copies of this data block could exist after the operation in column 3 has taken place** and ignore column 4 and 5 for now. Table M11.4-1 has been completed for you. Make sure the answers in this table make sense to you.

Problem M11.4.B

Mbus Cache Block State Transition Table

In this problem, we ask you to fill out the state transitions in Column 4 and 5. In column 5, fill in the resulting state after the operation in column 3 has taken place. In column 4, list the necessary Mbus transactions that are issued by the cache as part of the transition. Remember, the protocol should be optimized such that data is supplied using **CCI** whenever possible, and only the cache that *owns* a line should issue **CCI**.

Problem M11.4.C

Adding atomic memory operations to MBus

We have discussed the importance of atomic memory operations for processor synchronization. In this problem you will be looking at adding support for an atomic fetch-and-increment to the MBus protocol.

Imagine a dual processor machine with CPUs A and B. Explain the difficulty of CPU A performing fetch-and-increment(x) when the most recent copy of x is cleanExclusive in CPU B's cache. You may wish to illustrate the problem with a short sequence of events at processor A and B.

Fill in the rest of the table below as before, indicating state, next state, where the block in question may reside, and the CPU A and MBus transactions that would need to occur atomically to implement a fetch-and-increment on processor A.

State	other cached	ops	actions by this cache	next state	this cache	other caches	mem
Invalid	yes	read					
		write					

initial state	other cached	ops	actions by this cache	final state	this cache	other caches	mem	
Invalid	no	none	none	I			√	
		CPU read	CR	CE	√		√	
		CPU write	CRI	OE	√			
		replace	none	<i>Impossible</i>				
		CR	none	I		√	√	
		CRI	none	I		√		
		CI	none	<i>Impossible</i>				
		WR	none	<i>Impossible</i>				
		CWI	none	I				√
Invalid	yes	none	same as above	I		√	√	
		CPU read		CS	√	√	√	
		CPU write		OE	√			
		replace		<i>Impossible</i>				
		CR		I		√	√	
		CRI		I		√		
		CI		I		√		
		WR		I		√	√	
		CWI		I				√

initial state	other cached	ops	actions by this cache	final state	this cache	other caches	mem
cleanExclusive	no	none	none	CE			
		CPU read					
		CPU write					
		replace					
		CR		CS			
		CRI					
		CI					
		WR					
CWI							

Table M11.4-1

initial state	other cached	ops	actions by this cache	final state	this cache	other caches	mem
ownedExclusive	no	none	none	OE			
		CPU read					
		CPU write					
		replace					
		CR		OS			
		CRI					
		CI					
		WR					
CWI							

initial state	other cached	ops	actions by this cache	final state	this cache	other caches	mem
cleanShared	no	none	none	CS			
		CPU read					
		CPU write					
		replace					
		CR					
		CRI					
		CI					
		WR					
CWI							
cleanShared	yes	none	same as above				
		CPU read					
		CPU write					
		replace					
		CR					
		CRI					
		CI					
		WR					
CWI							

Table M11.4-2

initial state	other cached	ops	actions by this cache	final state	this cache	other caches	mem
ownedShared	no	none	none	OS			
		CPU read					
		CPU write					
		replace					
		CR					
		CRI					
		CI					
		WR					
ownedShared	yes	none	same as above				
		CPU read					
		CPU write					
		replace					
		CR					
		CRI					
		CI					
		WR					
CWI							

Table M11.4-3

Problem M11.5: Snoopy Cache Coherent Shared Memory

This problem improves the snoopy cache coherence protocol presented in Handout #14. As a review of that protocol:

When multiple shared copies of a *modified* data block exist, one of the caches *owns* the current copy of the data block instead of the memory (the owner has the data block in the OS state). When another cache tries to retrieve the data block from memory, the owner uses *cache to cache intervention* (CCI) to supply the data block. CCI provides a faster response relative to memory and reduces the memory bandwidth demands. However, when multiple shared copies of a *clean* data block exist, there is no owner and CCI is *not* used when another cache tries to retrieve the data block from memory.

To enable the use of CCI when multiple shared copies of a *clean* data block exist, we introduce a new cache data block state: *Clean owned shared* (COS). This state can only be entered from the clean exclusive (CE) state. The state transition from CE to COS is summarized as follows:

initial state	other cached	ops	actions by this cache	final state
cleanExclusive (CE)	no	CR	CCI	COS

There is no change in cache bus transactions but a slight modification of cache data block states. Here is a summary of the possible cache data block states (**differences from problem set highlighted in bold**):

- Invalid (**I**): Block is not present in the cache.
- Clean exclusive (**CE**): The cached data is consistent with memory, and no other cache has it. **This cache is responsible for supplying this data instead of memory when other caches request copies of this data.**
- Owned exclusive (**OE**): The cached data is different from memory, and no other cache has it. This cache is responsible for supplying this data instead of memory when other caches request copies of this data.
- Clean shared (**CS**): The data has not been modified by the corresponding CPU since cached. Multiple **CS** copies and at most one **OS** copy of the same data could exist.
- Owned shared (**OS**): The data is different from memory. Other **CS** copies of the same data could exist. This cache is responsible for supplying this data instead of memory when other caches request copies of this data. (Note, this state can only be entered from the **OE** state.)
- **Clean owned shared (COS): The cached data is consistent with memory. Other CS copies of the same data could exist. This cache is responsible for supplying this data instead of memory when other caches request copies of this data. (Note, this state can only be entered from the CE state.)**

Problem M11.5.A

Fill out the state transition table for the new COS state:

initial state	other cached	ops	actions by this cache	final state
COS	yes	none	none	COS
		CPU read		
		CPU write		
		replace		
		CR		
		CRI		
		CI		
		WR		
		CWI		

Problem M11.5.B

The COS protocol is not ideal. Complete the following table to show an example sequence of events in which multiple shared copies of a clean data block (*block B*) exist, but CCI is *not* used when another cache (*cache 4*) tries to retrieve the data block from memory.

cache transaction	source for data	state for data block B			
		cache 1	cache 2	cache 3	cache 4
0. <i>initial state</i>	—	I	I	I	I
1. cache 1 reads data block B	memory	CE	I	I	I
2. cache 2 reads data block B	CCI	COS	CS	I	I
3. cache 3 reads data block B	CCI	COS	CS	CS	I
4.					
5.					

Problem M11.5.C

As an alternative protocol, we could eliminate the CE state entirely, and transition directly from I to COS when the CPU does a read and the data block is not in any other cache. This modified protocol would provide the same CCI benefits as the original COS protocol, but its performance would be worse. **Explain the advantage of having the CE state.** You should not need more than one sentence.

Problem M11.6: Snoopy Caches

This part explores multi-level caches in the context of the bus-based snoopy protocol discussed in Lecture 14 (2017). Real systems usually have at least two levels of cache, smaller, faster L1 cache near the CPU, and the larger but slower L2. The two caches are usually inclusive, that is, any address in L1 is required to be present in L2. L2 is able to answer every snoop inquiry immediately but usually operates at 1/2 to 1/4th the speed of CPU-L1 interface. For performance reasons it is important that snoopers steal as little bandwidth as possible from L1, and does not increase the latency of L2 responses.

Problem M11.6.A

Consider a situation when the L2 cache has a cache line marked Sh, and an ExReq comes on the bus for this cache line. The snoopers asks both L1 and L2 caches to invalidate their copies but responds OK to the request, even before the invalidations are complete. Suppose the CPU ends up reading this value in L1 before it is truly discarded. What must the cache and snoopers system do to ensure that sequential consistency is not violated here?

Hint: Consider how much processing can be performed safely on the following sequences after an invalidation request for x has been received

Ld x; Ld y; Ld x

Ld x; St y; Ld x

Problem M11.6.B

Consider a situation when L2 has a cache line marked Ex and a ShReq comes on the bus for this cache line. What should the snoopers do in this case, and why?

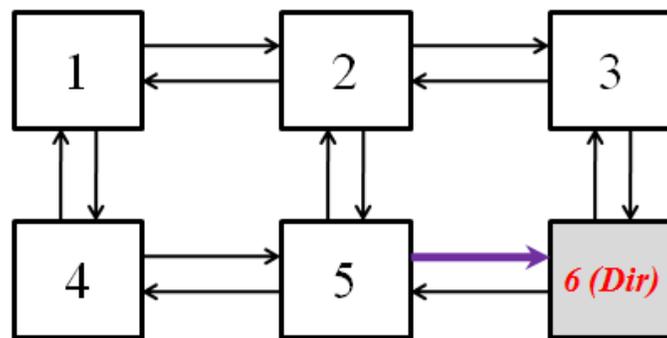
Problem M11.6.C

When an ExReq message is seen by the snoopers and there is a Wb message in the C2M queue waiting to be sent, the snoopers replies *retry*. If the cache line is about to be modified by another processor, why is it important to first write back the already modified cache line? Does your answer change if cache lines are restricted to be one word? Explain.

Problem M11.7: Directory-based Protocol

Problem M11.7.A

The following questions deal with the directory-based protocol discussed in class. Assume XY routing, and message passing is FIFO. (**XY routing algorithm** first routes packets horizontally, and then vertically towards their Y coordinates.) Protocol messages with the same source and destination sites are always received in the same order as that in which they were sent. **For this question, assume that the cache coherence protocol is free from deadlock, livelock and starvation.**



Assume the node 6 serves as the home directory, where the states for memory blocks are stored. Assume all caches are initially empty and no responses are sent voluntarily (i.e. every response is caused by a request)

Processor 1	Processor 4	Processor 5
I1.1: ST X, 10	I4.1: LD R1, X	I5.1: ST X, 20

Suppose the global execution order is as follows:

I4.1 => I5.1 => I1.1

Assume that the next instruction will start its execution only when the previous instruction has completed. For each instruction, list all protocol messages that are sent over the link 5 -> 6 (the purple link in the above figure).

I4.1:

I5.1:

I1.1:

Problem M11.7.B

For the directory protocol, we assume the message passing to be FIFO, meaning protocol messages with the same source and destination are always received in the same order as that in which they were sent. Now suppose messages can be delivered out-of-order for the same source and destination pairs. Describe one scenario that the cache coherence protocol will break due to this out-of-order delivery.

Problem M11.7.C

Under the 6823 directory-based protocol, a cache will receive a writeback request from the directory $\langle M2C_Req, a, S \rangle$ for address “a” when it is in state M and another cache wants a shared copy. Is it possible for a cache in the S state to receive $\langle M2C_Req, a, S \rangle$? Describe how this scenario can occur using the messages passed between the cache and the memory, and the state transitions.

Problem M11.8: Synchronicity (Spring 2014 Quiz 4, Part B)

You are writing a queue to be used in a multi-producer/single-consumer application. (Producer threads write messages that are read by one consumer.) We assume here a queue with infinite space. The basic code is shown below.

TST *rs*, Imm(*rt*) is the test-and-set instruction, which *atomically* loads the value at Imm(*rt*) into *rs*, and if the value is zero, updates the memory location at Imm(*rt*) to 1. This atomic instruction is useful for implementing locks: a value of 1 at the memory location indicates that someone holds the lock, and a value of 0 means the lock is free.

Producer pushes a message onto queue: (memory operations in bold)

```
void push(int** tail_ptr, int* tail_write_lock, int message) {
    while (lock_try(tail_write_lock) == false);
    **tail_ptr = message;
    *tail_ptr++;
    lock_release(tail_write_lock);
}

# R1 - contains address of data to enqueue
# R2 - contains the address of the tail pointer of queue
# R3 - address of tail pointer write lock
P1 SpinLock: TST R4, 0(R3)      # try to acquire tail write lock
P2          BNEZ R4, R4, SpinLock
P3          LD R4, 0(R2)       # get tail pointer
P4          ST R1, 0(R4)       # write message to tail
P5          ADD R4, R4, 4      # update tail pointer
P6          ST R4, 0(R2)       # release lock
P7          ST R0, 0(R3)
```

Consumer pops a message off queue: (memory operations in bold)

```
int pop(int** head_ptr, int** tail_ptr) {
    while (*head_ptr == *tail_ptr);
    int message = **head_ptr;
    *head_ptr++;
    return message;
}

# R1 - will receive address contained in message
# R2 - contains the address of the head pointer of queue
# R3 - contains the address of the tail pointer of the queue
C1 Retry:   LD R4, 0(R2)       # get head pointer
C2          LD R5, 0(R3)       # get tail pointer
C3          SUB R5, R4, R5     # is there a message?
C4          BNEZ R5, Pop
C5          JMP Retry
C6 Pop:     LD R1, 0(R4)       # read message from queue
C7          ADD R4, R4, 4      # update head pointer
C8          ST R4, 0(R2)
```

Problem M11.8.A

You are trying to port this code to an architecture that does not have the TST instruction (but, happily, the rest of the ISA is unchanged). Instead the new architecture has load-reserve/store-conditional instructions. Implement `TST rs, 0(rt)` using load-reserve/store-conditional:

```
LR rs, Imm(rt):
    rs ← Memory[(rt) + Imm]
    Track address (rt) + Imm

SC rs, Imm(rt):
    If (rt) + Imm modified:
        rs ← 0                                # Fail
    Else:
        Memory[(rt) + Imm] = (rs)           # Succeed
        rs ← 1
```

Problem M11.8.B

This new architecture is also *not* sequentially consistent. Give an example of memory orderings between the producer and consumer that would result in incorrect behavior. *Explain your answer fully or you will not receive credit.*

Your answer should look something like:

P1, P3, P4, C1, C2, P6, P7, C1, C2, C6, C8

(Except that this is a sequentially consistent ordering, so it is not a correct answer.)

Problem M11.8.C

Show where memory fences should be added to the producer and consumer code to ensure correctness with a weak consistency model. Explain your answer fully.

P1 SpinLock:TST R4, 0(R3) # try to acquire tail write lock

P2 BNEZ R4, R4, SpinLock

P3 LD R4, 0(R2) # get tail pointer

P4 ST R1, 0(R4) # write message to tail

P5 ADD R4, R4, 4 # update tail pointer

P6 ST R4, 0(R2)

P7 ST R0, 0(R3) # release lock

C1 Retry: LD R4, 0(R2) # get head pointer

C2 LD R5, 0(R3) # get tail pointer

C3 SUB R5, R4, R5 # is there a message?

C4 BNEZ R5, Pop

C5 JMP Retry

C6 Pop: LD R1, 0(R4) # read message from queue

C7 ADD R4, R4, 4 # update head pointer

C8 ST R4, 0(R2)

Problem M11.8.D

Let's next consider performance with a single producer thread and consumer thread. The following happens repeatedly:

1. The producer executes all instructions to push a message on the queue.
2. The consumer executes all instructions to pop a message off the queue.

Assume data, head, and tail pointers all lie in different, non-conflicting cache blocks.

First, after a few messages have been sent through the queue, will the consumer ever miss reading the head pointer? Will the producer ever miss reading the tail write lock, or fail to acquire the tail write lock? Explain in one or two sentences.

Problem M11.8.E

We'll now focus on the tail pointer only. Assuming a MSI invalidate coherence protocol, show the state of the tail pointer in the producer and consumer cache after each operation in the sequence below. Show any data or permissions transfers, e.g. "Memory→C" or "C invalidates P".

Operation	Producer tail pointer state	Consumer tail pointer state	Transfers
P1 TST try lock			
P3 LD tail_ptr			
P4 ST message			
P6 ST new_tail			
P7 ST release lock			
C1 LD head_ptr			
C2 LD tail_ptr			
C6 LD message			
C7 ST new_head			
P1 TST try lock			
P3 LD tail_ptr			
P4 ST message			
P6 ST new_tail			
P7 ST release lock			
C1 LD head_ptr			
C2 LD tail_ptr			
C6 LD message			
C7 ST new_head			

How many state transitions occur per message in the steady state?

Problem M11.8.F

Stay focused on the tail pointer only. Assume an update coherence protocol where the state of each line is either valid (V) or invalid (I). Show the state of the tail pointer in the producer and consumer cache after each operation in the sequence below in the steady state. Show any data or permissions transfers, e.g. “Memory→C” or “C invalidates P”.

Operation	Producer tail pointer state	Consumer tail pointer state	Transfers
P1 TST try lock			
P3 LD tail_ptr			
P4 ST message			
P6 ST new_tail			
P7 ST release lock			
C1 LD head_ptr			
C2 LD tail_ptr			
C6 LD message			
C7 ST new_head			
P1 TST try lock			
P3 LD tail_ptr			
P4 ST message			
P6 ST new_tail			
P7 ST release lock			
C1 LD head_ptr			
C2 LD tail_ptr			
C6 LD message			
C7 ST new_head			

How many state transitions occur per message in the steady state?

Problem M11.8.G

Your new architecture supports “remote access” for cached lines. This lets you assign a “home cache” for lines so that all memory operations will be sent *over the network* to operate remotely on the line *without allocating it in the requesting cache*.

For example, if line 0x100 is homed to processor A, and processor B writes 0x100, then *processor A’s cache will be updated* and processor B’s will be unchanged.

Assume the tail pointer is mapped to the producer’s cache, and the cache uses an MSI invalidate protocol (similar to Question 5). Once again, show the state of the tail pointer for the sequence of operations in the steady state and data/permission transfers:

Operation	Producer tail pointer state	Consumer tail pointer state	Transfers
P1 TST try lock			
P3 LD tail_ptr			
P4 ST message			
P6 ST new_tail			
P7 ST release lock			
C1 LD head_ptr			
C2 LD tail_ptr			
C6 LD message			
C7 ST new_head			
P1 TST try lock			
P3 LD tail_ptr			
P4 ST message			
P6 ST new_tail			
P7 ST release lock			
C1 LD head_ptr			
C2 LD tail_ptr			
C6 LD message			
C7 ST new_head			

How many state transitions occur per message in the steady state?

Problem M11.9: Cache Coherence (Spring 2015 Quiz 3, Part B)

Ben Bitdiddle is designing a snoopy-based, write-invalidate MSI protocol for write-back caches. Under the standard MSI protocol, when a cache observes a Bus Read Exclusive message (BusRdX), it has to invalidate its own copy of the cache block. Ben instead proposes an optimization, called delayed invalidation, to potentially reduce the number of read misses. The optimization works as follows:

Delayed invalidation: When a cache observes a Bus Read Exclusive message (BusRdX) and it has a copy of the block in the Shared (S) state, the cache delays the invalidation of the block until before a cache miss happens. In other words, the cache will treat any subsequent requests from its own processor as if the BusRdX had not happened, until one of those requests causes a miss. At that point, all pending invalidations are performed before processing the miss.

Problem M11.9.A

Suppose processors P1 and P2 have private, snoopy caches. Both caches are initially empty. Consider the following sequence of accesses:

I0	P2: read	A
I1	P1: write	A
I2	P2: read	A
I3	P1: write	A
I4	P2: read	A
I5	P2: read	B
I6	P2: read	A

Assume blocks A and B do not conflict in the cache. Compare Ben's delayed invalidation optimization with the standard MSI protocol by filling the states (on the next page) for each cache block after each operation is done and calculate the number of misses in both cases.

Assume we use the standard MSI protocol. Fill in the following table.

Standard MSI Protocol				
	Processor P1's Cache		Processor P2's Cache	
Initial State	A: I	B: I	A: I	B: I
After P2 reads A	A: I	B: I	A: S	B: I
After P1 writes A	A:	B:	A:	B:
After P2 reads A	A:	B:	A:	B:
After P1 writes A	A:	B:	A:	B:
After P2 reads A	A:	B:	A:	B:
After P2 reads B	A:	B:	A:	B:
After P2 reads A	A:	B:	A:	B:

How many misses occur in the two caches?

Assume we adopt Ben's delayed invalidation optimization. Fill in the following table. If there is a delayed invalidation, write it in the invalidation queue (the "Inv Queue" column). For example, "Inv L" means there is a delayed invalidation on block L.

MSI Protocol with Delayed Invalidation						
	Processor P1's Cache			Processor P2's Cache		
	MSI state		Inv Queue	MSI state		Inv Queue
Initial State	A: I	B: I		A: I	B: I	
After P2 reads A	A: I	B: I		A: S	B: I	
After P1 writes A	A:	B:		A:	B:	
After P2 reads A	A:	B:		A:	B:	
After P1 writes A	A:	B:		A:	B:	
After P2 reads A	A:	B:		A:	B:	
After P2 reads B	A:	B:		A:	B:	
After P2 reads A	A:	B:		A:	B:	

How many misses occur in the two caches?

Problem M11.9.B

Does Ben's delayed invalidation optimization violate cache coherence rules? Please explain your answer in one or two sentences.

Problem M11.9.C

Suppose the original system guarantees sequential consistency. Does adding the delayed invalidation optimization break sequential consistency? Please explain your answer in one or two sentences. If your answer is yes, please provide a sequence of load/store operations that violates sequential consistency.

Problem M11.9.D

Ben only applies delayed invalidation on cache blocks that are in the S state. When a cache observes a Bus Read Exclusive message (BusRdX) and the associated cache block is in the Modified (M) state, it sends out the data in response to a BusRdX message and changes the cache state to Invalid (I).

Is it possible to delay invalidation when the cache block is in the Modified (M) state? If it is not, please explain why. If it is possible, please describe how to make delayed invalidations work when the block is in the M state. In other words, please describe the actions the cache needs to take when the cache observes a BusRdX message, how to handle subsequent read and write accesses if the invalidation is delayed, and when the invalidation needs to be processed.

Problem M11.10: Cache Coherence (Spring 2015 Quiz 3, Part C)

Please use Handout #15 to answer the questions in this part.

Problem M11.10.A

Ben designs an architecture that does not have the atomic compare-and-swap (CAS) instruction but has load-reserve (LR) and store-conditional (SC) instructions.

Help Ben implement a Boolean compare-and-swap instruction BCAS *old*, *new*, Imm(*base*) using load-reserve and store-conditional instructions:

```
LR rs, Imm(rt):
    <flag, addr> ← <1, rt + Imm>
    rs ← Memory[rt + Imm]

SC rs, Imm(rt):
    If <flag, addr> == <1, rt + Imm>:
        Memory[rt + Imm] ← rs
        rs ← 1                # Succeed
    Else:
        rs ← 0                # Fail
```

BCAS is a simplified CAS instruction that only deals with values 0 and 1. You can use temporary registers (*tmp1*, *tmp2*, *tmp3*...) and any algorithmic, logical, memory, and branch instructions in the MIPS instruction set.

Problem M11.10.B

Suppose the hardware where the shared-memory queue from Handout #15 is executed has a weak consistency model that relaxes all the orderings of reads and writes. Give an example of memory orderings between the producer and consumer that would result in incorrect behavior.

Please fully explain your answer to get full credit.

Your memory ordering example should look something like:

P1, C2, P2, C4, P4, C5, C7, C9, C10

Problem M11.10.C

Please add the minimum number of memory fences ($FENCE_{WR}$, $FENCE_{RW}$, $FENCE_{WW}$, or $FENCE_{RR}$) to the producer and consumer codes to ensure correctness with a weak consistency model. Please explain your answer fully.

Code for producer to enqueue a message:

```
P1:  LD  R3, 0(R2)  # get tail pointer

P2:  ST  R1, 0(R3)  # write message to tail

P3:  ADD R3, R3, 4  # update tail pointer

P4:  ST  R3, 0(R2)
```

Code for consumer to dequeue a message:

```
C1: SpinLock: MOV  R6, R0          # set R6 to 0

C2:          CAS  R6, R5, 0(R4) # try to acquire lock

C3:          BNEZ R6, SpinLock

C4:          LD   R7, 0(R2)      # get head pointer

C5: Retry:  LD   R8, 0(R3)      # get tail pointer

C6:          BEQ  R7, R8, Retry # is there a message?

C7:          LD   R1, 0(R7)      # read message from queue

C8:          ADD  R7, R7, 4      # update head pointer

C9:          ST   R7, 0(R2)

C10:         ST   R0, 0(R4)      # release lock
```