

CS4223 Assignment 2: Design and Implementation of Tracing Simulator for MESI, MOESI and Dragon Cache Coherency Protocols

JULIUS PUTRA TANU SETIAJI (A0149787E), KYLE TIMOTHY NG CHU (A0155181J)

1 CHOICE OF PROGRAMMING LANGUAGE

Initially, we tried to implement the trace simulator using Python. However, after a while, we realised that it gets unwieldy, unreadable, and difficult to extend given the complexity of implementing a trace simulator. Particularly, while Python has enumerations, they are not as powerful as Algebraic Data Types (ADT) in other languages, especially with regards to associated values. Python's lack of pattern matching also makes implementing a finite state machine (FSM) a hassle. While Python's dynamic typing allows more rapid prototyping, a better compile-time typecheck would be really useful to prevent bugs.

Haskell seems to fit the bill really well here. However, since we are essentially simulating a mutable, eagerly-evaluating von Neumann architecture machine, Haskell's characteristic of being immutable and lazy can become an impediment. We need a language that has some of Haskell's features such as being strongly and statically typed, pattern matching and ADT, while still allowing us to have mutability. There are some candidates, such as Scala, Rust and OCaml. For these reasons as well as familiarity, completeness of the standard library and availabilities of third party libraries we chose to implement our tracing simulator using **Scala**.

2 IMPLEMENTATION

Our code can be divided into 2 large parts: common code for all the protocols, and the specific code for the trace simulator of each cache coherency protocols.

Overall, our code uses the delegation pattern, where each device can choose to become a delegate of another device. This way, devices can have a two-way communication of some sort. The devices themselves are designed as Finite State Machines (FSM) to allow for more robust code.

There are some potential issues that may prevent correctness in this concurrent system, namely deadlock, livelock, starvation and race conditions [1]. As we explain the design of the devices in our simulator, we will mention how we try to prevent these problems from happening.

2.1 `Main.scala` and `Protocol.scala`

This is the main driver code of the simulator. It takes in arguments passed by user and parses them into the appropriate types, such as `Int`, `String` and `Protocol`. `Protocol` is actually defined inside `Protocol.scala`, which is an enumeration of supported cache coherency protocols. In our case we have `MESI` and `Dragon`. If all the arguments are parsed successfully, the appropriate tracing simulator is then invoked.

2.1.1 `Debug.scala`. This file provides code for debug printing that can be elided for final production code.

2.1.2 `Address.scala`. This file provides the case class for addressing in an n-way associative cache through tag and set index.

2.2 coherence.bus

2.2.1 `Bus`. Across the different protocols, we use the same bus implementation. The design that we chose for our bus is **an atomic bus**, where one cache is granted exclusive access to the bus to perform one transaction (a transaction is defined as the cache granted exclusive access placing command on the bus, and then receiving response by other caches).

When a device wants to place command on the bus, it has to first request access from the bus, and the device has to be a `BusDelegate`. This ensures that a cache will have exclusive access to the bus, with the bus arbitrator acting like a mutex lock. The bus arbitrator uses a first come first serve arbitration policy to *prevent starvation*.

When the bus arbitrator decides that it is the turn for one `BusDelegate` to obtain exclusive access to the bus, the bus will invoke `busAccessGranted()` on the `BusDelegate` to obtain the command that the device wants to place on the bus. This is to *prevent race condition*, because the transition of the states of the cache lines are not actually atomic even though the bus is atomic. Consider the case in MESI protocol when Processors P1 and P2 write to cache line A simultaneously (both need to issue `BusUpgr` to move the line from the S state to the M state). Suppose P1 wins bus access and places `BusUpgr` on the bus. P2 is still waiting for bus access to place its own `BusUpgr`, but it cannot proceed because P1 has exclusive access to the bus. P2 then receives `BusUpgr` from P1, thus it must invalidate line A according to the MESI protocol. As such, P2 must also change its pending `BusUpgr` request to a `BusRdX`. With our design, we will only ask the cache what command to place on the bus when they are granted access, thus allowing the caches to change the command that they send between the moments access is requested and access is granted.

After a device places a command on the bus, `onBusCompleteMessage` method on all other devices as `BusDelegate` will get invoked. This allows them to make internal state transitions in response to the command on the bus, as well as place their response on the bus. The bus will arbitrarily grant one device to place their response on the bus. Once the device having exclusive access to the bus decided that it has received enough responses, it will indicate to the bus that it wishes to relinquish its exclusive access. This fact is propagated to all other devices through invocation of the `onBusTransactionEnd` of the `BusDelegate`. The bus will then pick another pending access request, and perform the same cycle all over again. Note that to *avoid livelock*, the device that obtains exclusive access must be allowed to complete before exclusive access is relinquished, thus our design of letting the device decide when to relinquish exclusive access of the bus. Consider two processors P1 and P2 writing to cache block B using MESI cache coherency protocol. First, P1 acquired access to the bus and issues `BusRdX`, resulting in P2 invalidating its cache line. However, without exclusive ownership, this scenario can happen: before P1 can perform the write, P2 acquires bus, issues `BusRdX`, resulting in P1 invalidating its cache line, and so on so forth, which results in a livelock.

Note that in our bus design, each device connected to the bus is not blocked on requesting access to the bus. They are still able to service incoming transactions while waiting to be granted exclusive access to the bus. This is done to *prevent protocol-level deadlock*, sometimes called *fetch-deadlock* [1].

2.2.2 BusDelegate. Specifies the hooks in a device that wants to hear back from the Bus regarding events on the Bus. Of interest is the `hasCopy` method, which is invoked in order to check the “SHARED” OR line, required by the protocols to check whether a particular cache line already exists in the other caches.

2.2.3 BusStatistics. This is a “template” for each protocol to provide a class that takes note of statistics in the Bus

2.2.4 MessageMetadata. A case class that stores the message that a cache wants to place on the bus, as well as the address affected by that message. The size of each message here is assumed to be 1 word long.

2.2.5 ReplyMetadata. A case class that stores the reply from other caches, as well as the size of the reply. The size here will be used in the calculation of how long it takes to transmit the reply and associated data on the Bus.

2.3 coherence.cache

This package contains classes necessary to build a cache with a Least Recently Used (LRU) cache eviction policy and where each cache line stores a state.

2.3.1 LRUCache. This is a simple cache with LRU cache eviction policy that can be used for each set in an n-way associative cache. It is implemented on top of Java’s built-in `LinkedHashMap` which uses a hash table and linked list to keep track of the order of when an item was last used.

2.3.2 CacheLine. This is a simple case class representing a cache line. It is generic in the type parameter `State`, and stores only the state of this cache line.

2.4 coherence.devices

This package contains mostly the “template” for each type of specific devices. This helps to keep our code DRY¹ as the common code among all the protocol are put here.

2.4.1 Device. A simple “template” for all devices in the simulator. It only requires that every device has a method `cycle()` which is invoked to notify the device to advance one cycle and perform what they need to do in a cycle.

2.4.2 Memory. This is the “template” for a memory in each protocol. The memory will always return false as it is the identity value in the OR operation on the shared line, thus the memory will not affect the “SHARED” OR line.

¹Don’t Repeat Yourself, one important Software Engineering Principle

2.4.3 Processor. This is a concrete class that represents a processor in the simulator. Each processor owns a cache, where the processor is the cache's CacheDelegate. It is an FSM with 4 states, with transitions as shown in Figure 1:

- Ready – the processor is ready to execute the next instruction.
- Operation – the processor is currently executing a non-memory instruction
- Cache – the processor is waiting for the result of a memory instruction from its cache
- Finished – the processor has no more instruction to execute

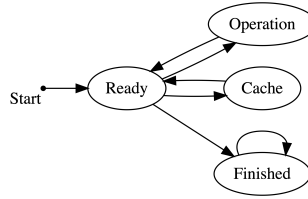


Fig. 1. The state diagram for the processor

2.4.4 ProcessorOp. This trait is effectively an ADT representing the three possible instructions: the memory instructions Load and Store which has the address as the associated value as well as the non-memory instruction Other with the number of cycles required for execution as the associated value. The associated singleton object of the ProcessorOp trait provides a function to parse each line in the input trace into a ProcessorOp.

2.4.5 Cache. This is the “template” for the protocol-specific implementation of the cache.

2.4.6 CacheDelegate. This is the “template” for a cache delegate, which requires only one method, requestCompleted which is invoked when the operation is complete.

2.4.7 CacheOp. An ADT representing the possible cache operations.

2.4.8 CacheStatistics. This is a “template” for each protocol to provide a class that takes note of statistics in the Cache

2.5 MESI Protocol (coherence.mesi)

The MESI protocol consists of four states: Modified (M), Exclusive (E), Shares (S), and Invalid (I). The state transition diagram can be found in Figure 2. For any given pair of caches, the permitted states of a given cache line is given in Table 1

2.5.1 Message. There are 4 possible messages to be placed by a cache that gained exclusive access to the bus:

- (1) BusRd – a read request to a cache block
- (2) BusRdX – a read exclusive request to a cache block, typically when a cache wants to write to a cache block it does not already have.

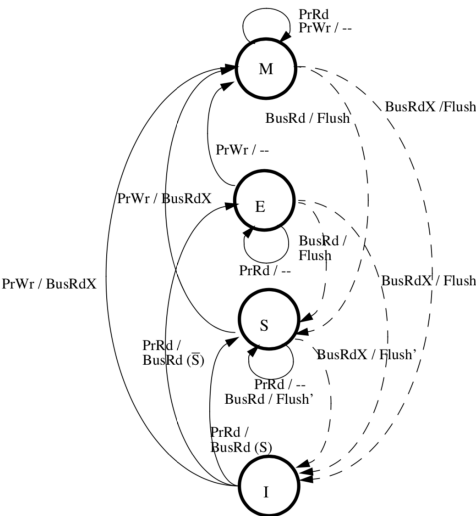


Fig. 2. State Transition Diagram of the MESI Cache Coherency Protocol

	M	E	S	I
M	No	No	No	Yes
E	No	No	No	Yes
S	No	No	Yes	Yes
I	Yes	Yes	Yes	Yes

Table 1. Permitted States for a Pair of Caches in MESI Protocol

- (3) BusUpgr – a request to gain read exclusive access to a shared cache block that a cache already has.
- (4) Flush – a request to write back the cache block to the main memory

2.5.2 Reply. There are 4 possible replies to be placed by a cache or memory to a Message:

- (1) Flush – when a cache block must be written back to the main memory
- (2) FlushOpt – when a cache block is posted on the bus to supply it to another cache
- (3) MemoryRead – when the memory posts a cache block on the bus
- (4) WritebackOk – acknowledgement from the memory that a cache block has been written successfully

2.5.3 Memory. The memory in our MESI Protocol simulator snoops on the bus and will do one of these three things:

- (1) Do nothing if the command is BusUpgr
- (2) Read from the address in the command placed by a cache if the command is BusRd or BusRdX
- (3) Write back to memory if the command is Flush

If the memory detects that another cache is providing the cache line instead, e.g. if the state is S, then the memory will cancel the current read request. If the memory detects that a cache replied with OK, it will cancel the read request, and start the writeback process instead.

2.5.4 Cache. The cache is an FSM with 8 states, whereby each state will inform the cache about what to do with regards to the bus:

- Ready – the cache is ready to accept commands from the processor
- WaitingForBus – the cache has requested access to the bus but is waiting for that access to be granted
- WaitingForReplies – the cache has exclusive access to the bus and has placed a message on the bus, and is waiting for replies from the other caches.
- WaitingForWriteback – the cache has exclusive access to the bus and has snooped a Flush request, and is currently waiting for the memory to acknowledge that the writeback succeeded.
- WaitingForResult – the cache has the cache line requested, and that cache line is valid, and the cache is currently experiencing its hit latency.
- WaitingForBusUpgrPropagation – the cache sent a BusUpgr message and is waiting for the message to reach the other caches such that the other caches will invalidate their cache lines.
- EvictWaitingForBus – the cache needs to evict a cache line, has requested access to the bus and is currently waiting for that access to be granted.
- EvictWaitingForWriteback – the cache currently has exclusive access to the bus and has placed the flush command on the bus and is waiting for the memory to acknowledge the writeback.

The state transitions are shown in Figure 3

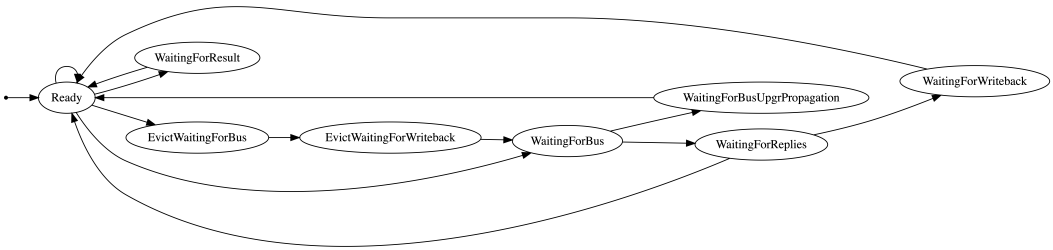


Fig. 3. The state diagram for the cache

2.6 MOESI Protocol (coherence.moesi): Optimisation to Basic MESI Protocol

Some processor architectures such as UltraSPARC [3] and AMD Opteron [2] use an optimised version of the MESI protocol by adding one more state: O for Owned, representing a cache line that is both modified and shared. The idea is to allow cache having a modified cache line to supply that data to another cache without performing a slow writeback

to memory. The state transition diagram can be found in Figure 4. For any given pair of caches, the permitted states of a given cache line is given in Table 2.

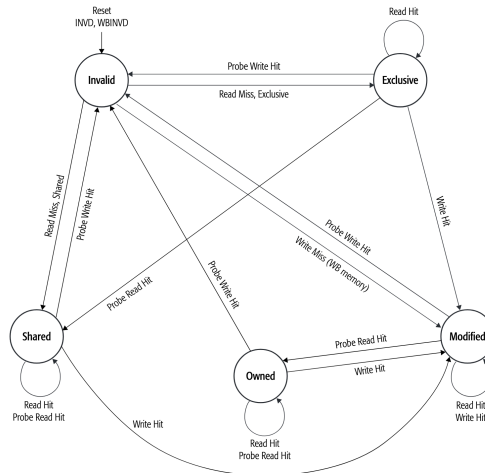


Fig. 4. State Transition Diagram of the MOESI Cache Coherency Protocol

	M	O	E	S	I
M	No	No	No	No	Yes
O	No	No	No	Yes	Yes
E	No	No	No	No	Yes
S	No	Yes	No	Yes	Yes
I	Yes	Yes	Yes	Yes	Yes

Table 2. Permitted States for a Pair of Caches in MESI Protocol

2.6.1 Differences Compared to MESI.

2.7 Dragon Protocol

	E	Sc	Sm	M
E	No	No	No	No
Sc	No	Yes	Yes	No
Sm	No	Yes	No	No
M	No	No	No	No

Table 3. Permitted States for a Pair of Caches in Dragon Protocol

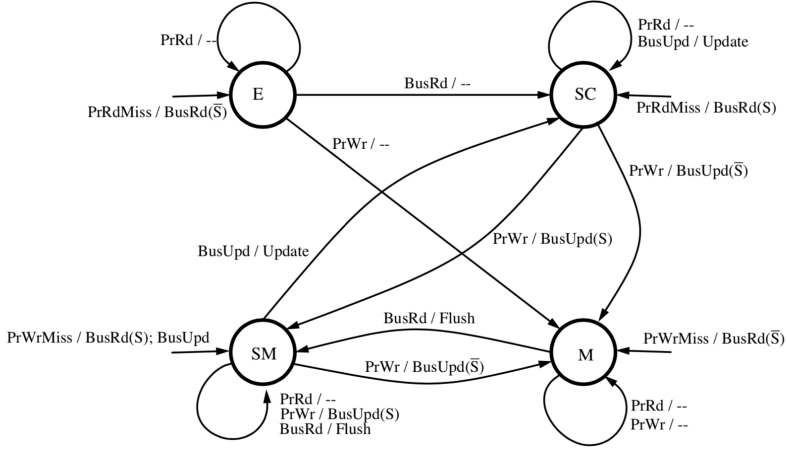


Fig. 5. State Diagram for the Dragon Protocol

3 QUANTITATIVE ANALYSIS

REFERENCES

- [1] David Culler, Jaswinder Pal Singh, and Anoop Gupta. 1998. *Parallel Computer Architecture: A Hardware/Software Approach*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.
- [2] Bruce Jacob, Spencer Ng, and David Wang. 2007. *Memory Systems: Cache, DRAM, Disk*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.
- [3] Sun Microsystems. 1997. UltraSPARC™ User's Manual. *STP1030-UG, Sun Microsystems (January 1996) (1997)*.