# Memory Barriers: A Hardware view for Software Hackers

Paul E. McKenney
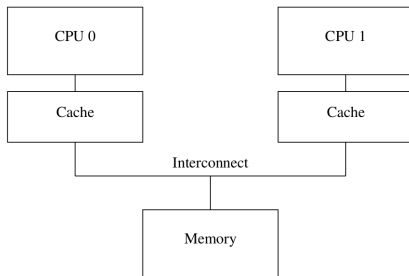
Presented by
Akshay Gopalakrishnan

December 6, 2021

# Introduction

- This paper represents the inner working of hardware that results in several non-sequential behaviors of our concurrent programs

- The paper is rife with examples as well as showcasing the reasons for having such hardware features which in turn help in our programs performing better.

- Along with the positives the author also carefully cautions why such rampant changes for performance might result in highly non-trivial behaviors being showcased by the hardware running our programs.

- The paper concludes by discussing the then versions of several concurrent hardware that exhibit different non-sequential behaviors.

## Cache structures



- An extra chunk of memory local to a given cpu (or multiple cpus in bigger systems).
- Data read from memory will also be saved in the cache if there is space.
- Any subsequent accesses to the same memory will be first looked at in the cache and if not present, from the main memory.

## The usefullness of Caches

- Used for faster access of memory.
- Useful when a single shared memory location is being accessed several times but not changed.
- Need to read from main memory which is at least 10x times slower than accessing caches.
- Overall performance of program is significantly improved.

In modern systems there are multiple cache levels that exist, each of which is based on the scope. For instance, L1 cache is local to just one core. Whereas L2 is to multiple cores in the same processor (could also be others). All this layering is done for performance, part of the reason why we have such highly non-trivial behaviors of our concurrent programs.

- Multiple caches need to be in synchronization to ensure no stale data in caches exist.
- Caches can communicate with each other via the interconnect network.
- We require a protocol to ensure that caches lines are updated accordingly.

An example of such a protocol is MESI (Modified Exclusive Shared and Invalid).

- Modified - cache line has the upto date data which resides in memory and the data has been stored by the corresponding CPU.
- Exclusive - cache line is not updated with the recent memory store done by the corresponding CPU.
- Shared - cache line is in read-only state (CPU needs to consult with other CPU caches before being able to write to it)
- Invalid - "empty" cache line and new data can be put here.

THe MESI protocol

- Multiple caches need to be in synchronization to ensure no stale data in caches exist.
- Caches can communicate with each other via the interconnect network.
- We require a protocol to ensure that caches lines are updated accordingly.

An example of such a protocol is MESI (Modified Exclusive Shared and Invalid).

- Modified - cache line has the upto date data which resides in memory and the data has been stored by the corresponding CPU.
- Exclusive - cache line is not updated with the recent memory store done by the corresponding CPU.
- Shared - cache line is in read-only state (CPU needs to consult with other CPU caches before being able to write to it)
- Invalid - "empty" cache line and new data can be put here.

Modified state is when my CPU is writiing to memory some data as well as duly updated the cache line. Now this line has the only copy of the latest data in memory. Exlcusive is one step behind modified state, wherein the CPU has updated memory, but not just its cache line. Shared state is when more than one cache line has the same data of memory. This only means, if I want to change one cache line, I must inform the others too. Invalid state can contain any stale data that is never going to be read by the CPU (as it is stale).

## MESI Protocol Messages

The following messages are passed by a cache line to other caches in the system.

- Read - Contains the physical address of the cache line to be read.
- Read Response - Contains the data requested by a previous Read message.
- Invalidate - Contains the physical address of the cache line to be invalidated.
- Invalidate Acknowledge - CPU receiving an Invalidate message must respond with this message once the specfied cache line is invalidated.
- Read Invalidate - Does the action of both Read and Invalidate in one message.
- Writeback - Contains both the address and the data tobe written back to memory (could also implicitly update other cache lines with this).

# Example of Cache Communication

The following table represents a sequence of actions done by CPU and the different states (MESI) of the caches after the action has been done.

| Sequence # | CPU # | Operation | CPU Cache | | | | Memory | |
|---|---|---|---|---|---|---|---|---|
| | | | 0 | 1 | 2 | 3 | 0 | 8 |
| 0 | | Initial State | -/I | -/I | -/I | -/I | V | V |
| 1 | 0 | Load | 0/S | -/I | -/I | -/I | V | V |
| 2 | 3 | Load | 0/S | -/I | -/I | 0/S | V | V |
| 3 | 0 | Writeback | 8/S | -/I | -/I | 0/S | V | V |
| 4 | 2 | RMW | 8/S | -/I | 0/E | -/I | V | V |
| 5 | 2 | Store | 8/S | -/I | 0/M | -/I | I | V |
| 6 | 1 | Atomic Inc | 8/S | 0/M | -/I | -/I | I | V |
| 7 | 1 | Writeback | 8/S | 8/S | -/I | -/I | V | V |

- The 0th sequence represents the default state of cache before being used. Each cache line is set to "Invalid" state.
- The 1st sequence represents a load done by CPU 0 to fetch data from address 0. The content in stored in CPU0 cache and the state changes to "Shared".
- The 2nd sequences represents a load done by CPU 0 to fethc data from address 0. The content is stored in CPU1 cache and the state changes to "Shared".

## Sequence 3 to 5

- The 3rd sequence represents a load done by CPU 0 to fetch data from address 8. This implicity means to Invalidate one's own cache line by sending an invalidate message to it. This gives space to store data from address 0, which now is in the "Shared" state.

- The 4th sequence represents an RMW done by CPU 2 on address 0. Before it can do this, it needs to invalidate other caches having data on this address by sending "Invalidate" message. Once that is done, it's own cache (CPU 2) is set to "Exclusive" state, while others which had data at address 0 to "Invalid".

- The 5th sequence represents the actual store done by CPU 2 (as part of RMW). This changes its own cache to state "Modified" and sets the memory of address 0 to "Invalid".

## Sequence 6 and 7

- The 6th sequence represents CPU 1 performing an atomic write to data at Address 0. Since CPU 2 has its cache in modified state, the increment needs to change that value, while also invalidating their cache. So a "Read Invalidate" message is sent to each cache line, get the updated store value and increment it by 1. Now it sets its own cache line (CPU 1) to modified state. While other caches having data at address 0 is set ot invalidate.

- The 7th sequence represents the actual commiting of the new data to memory. This can either be done by actually issuing a writeback or forcing the cache to make space for other address data (via a Load). Here, a Load is issued for address 0 by CPU 1, which forces a writeback to address 0. The value at memory address 0 is updated and the state becomes "Valid". Meanwhile the cache line of CPU 1 has a "Shared" state.

## Another example

Consider the example where CPU 0 wants to write to memory whose data is on cache line of CPU 1. This would require the following actions

- CPU 0 sends an Invalidate message to CPU 1.
- CPU 1 receives the message, sets the appropriate cache line to "Invalid" state.
- CPU 1 sends the Acknowledgement message along with the data on the cache line to CPU 0.
- CPU 0 receives the Acknowledgement and data.
- CPU 0 does the write to memory and updates its own cache.

Notice that in the example above, CPU 0 needs to stall itself until the Acknowledgement message is received from CPU 1. This is showcased in the figure below.



Figure 4: Writes See Unnecessary Stalls

CPU 0 need not stall as it will eventually do the write that it is supposed to do. Rather it can continue doing some other task.

# Solution to write stalling: Write Buffers

One solution to this is to have a store buffer for each CPU. Writes to be done will be stored in this buffer and the CPU can continue doing future tasks. When the Acknowledgement messages are received by other CPUs, the write can be committed to its cache line (and eventually to memory).



Figure 5: Caches With Store Buffers

## Added Complications

Adopting the above system does create its own problems.

- The main cause is that writes to memory committed to store buffers are not read by the same CPU from it.
- Instead, the same memory is read either from cache or main memory.
- This means, the CPU uses a stale value of memory which results in incorrect execution of programs.

```
1    a = 1;
2    b = a + 1;
3    assert(b == 2);
```

Here is the reason why the example above fails. Consider only CPU 0 and 1 exist.

- CPU 0 wants to do the write $a = 1$. So it sends a Read Invalidate message to CPU 1 and commits the store to its Store buffer.

- CPU 1 receives the message and sends the invalidate with the current value of $a$ from its cache.

- CPU 0 meanwhile wants to do $b = a+$, so starts reading $a$ from its own cache line. It receives the value 0.

- Now CPU 0 receives the message from CPU 1 and the store buffer flushes the write to CPU 0 cache line.

- CPU 0 does $b = a + 1$ having read $a = 0$ from its cache before.

- Now $b = 1$.

- The assertion fails.

## Solution seems simple: but?

The straightforward solution to this is to first let CPUs check their own store buffers during loads from memory and then if not found checking cache and/or memory. However, this also does not prevent all our problems. Consider the code below, where CPU 0 holds exclusive rights to cache line for *b*.

```
1 void foo(void)
2 {
3   a = 1;
4   b = 1;
5 }
6
7 void bar(void)
8 {
9   while (b == 0) continue;
10   assert(a == 1);
11 }
```

The problem is that CPU 0 (if running the first code snippet) can upadte value of *b* before committing the value of *a* to cache. This may result in the assertion to fail.

The reason why this would happen is the following steps

- CPU 0 wants to do $a = 1$, but $a$'s cache line is not exclusive/modified, so it sends the write to store buffer and sends Read Invalidate to CPU 1.

- CPU 1 wants to do *while*($b == 0$) so tries to read value of $b$. It does not exist in its cache line, so sends a read to CPU 0.

- CPU 0 now wants to do $b = 1$ and since it owns this cache line, commits the write immididately to cache.

- CPU 0 receives the read message from CPU 1 and sends $b = 1$ as a response to CPU 1.

Memory Barriers: A Hardware view for Software Hackers

2021-12-06

└─Solution seems simple: but?



- CPU 1 receives the read invalidate for *a* and sends Acknowledgement along with the value of *a* to CPU 0.

- CPU 1 then receives the value of *b*. Ends the loop (as $b = 1$) and moves to assertion.

- The assertion fails.

- CPU 0 receives the Acknowledgement from CPU 1 on *a*. So now it commits from store buffer to cache line.

The problem now is that the hardware does not recognize dependencies between memory accesses (meaning things like conditionals and assertions like the above example). At this point, the hardware itself cannot do much.

The solution to this was to have an instruction which specifically tells the hardware that some dependency exists. This instruction is called a Write Memory Barrier. The following code with a write barrier solves our issue.

```
1 void foo(void)
2 {
3    a = 1;
4    smp_mb();
5    b = 1;
6 }
7
```
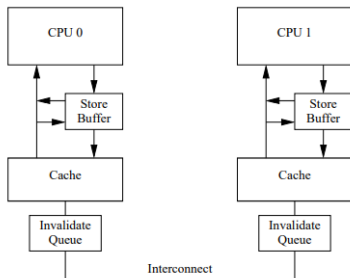
The reason inserting a write memory barrier there solves our problem is
that CPU 0 will never do the wrtie $b = 1$ until the store buffer holding
$a = 1$ is first committed to cache.

Store buffers are good for performance. But in hardware, these buffers are quite small. So if they get full or the CPU needs to process a write barrier instruction, it must wait until all the writes are flushed. This in turn depends on whether the CPU has received invalidate acknowledgement from all the CPUs in the hardware. This may take time as other CPUs may be busy accessing main memory or other cache lines at the moment, which take a long time. In addition, to send acknowledgement, the CPU must first access its cache and tag the appropriate cache line Invalid. All this takes longer and this stalling becomes a bottleneck for performance.

# Invalidate queues

One observation is that the CPU need not actually invalidate the cache line immediately before sending acknowledgement. The CPU only needs to do this if another message for the same cache line is requested.

For this, each CPU is equipped with its own invalidate queue. Invalidate requests are stored in this queue, promising the CPU requesting the acknowledgement that it would be done eventually. So the acknowledgement message is often sent immediately without modifying the cache line state.

## Added Complications

While doing the above is good, the CPUs still are not concerned with the state of the cache if they themselves receive a read/write access to be performed. In essence, on a read, the CPU might end up reading a stale value from cache line (with the invalidate message for that line still in queue).

The following example, run with store buffers and invalidate queues even with a write barrier will produce an incorrect result.

```
1 void foo(void)
2 {
3   a = 1;
4   smp_mb();
5   b = 1;
6 }
7
8 void bar(void)
9 {
10  while (b == 0) continue;
11  assert(a == 1);
12 }
```

The reason why the above example may lead to assertion failure is as follows. Suppose CPU0 runs foo and CPU1 bar. Also suppose $a$ is in shared state in cache while CPU0 has exclusive rights to cache line of $b$.

- CPU0 performs $a = 1$. But the cache line is "Shared". So it sends an invalidate message to CPU1.

- CPU1 receives the message, immidiately sends acknowledgement. But CPU1 does not yet change the cache line of $a$ to "Invalid".

- CPU1 processes $while(b == 0)$, by issuing a read request to CPU0 cache.

- CPU0 receives the acknowledgement. Now it updates its cache with the modified value of $a$.

- CPU0 processes the write barrier and proceeds to do $b = 1$. Since it has exclusive right to it, updates the cache line of $b$ right away.

- Now CPU0 receives the read request for $b$ and sends the value 1.

- CPU1 receives this value and ends the loop.

- Now for the assertion, CPU1 reads the value of $a$ from its own cache (without processing hte invalidate message for it as it thinks nobody else it bothering it with another cache message for that line).

- Hence the assertion fails.

The above example indicates that the memory barrier is sort of ignored. One solution for this example is to ensure that when the CPU comes across a memory barrier, it ensures that its invalidate queue is serviced (emptied) before proceeding. In literature one can say this as a read barrier. So the solution to the above problem is a barrier that will be inserted in bar too.

```
 8 void bar(void)
 9 {
10   while (b == 0) continue;
11   smp_mb();
12   assert(a == 1);
13 }
```

Here comes Read Memory Barriers

The above example indicates that the memory barrier is sort of
ignored. One solution for this example is to ensure that when the
CPU comes across a memory barrier, it ensures that its invalidate
queue is serviced (emptied) before proceeding. In literature one
can say this as a read barrier. So the solution to the above
problem is a barrier that will be inserted in bar too.

```
8  void bar(void)
9  {
10     while (b == 0) continue;
11     smp_mb();
12     assert(a == 1);
13 }
```

Currently, this is only viewed as a memory barrier instruction. In many
languages/hardwares, different variants of these barriers are provided. Of
these, we have looked at Read and Write barriers. However, we can further
dissect it based on what accesses should not occur before those accesses
which precede the barrier. Doing so, we get Read-Read(Read memory
barrier), Read-Write, Write-Read and Write-Write(write memory barrier)
variants of barriers.

While memory barriers are useful for ordering accesses, one has to note that they only order accesses for the CPU that performs these accesses. This means, CPUs might observe incoming writes from other CPUs in different orders. This in literature is known as non-multicopy-atomicity.

## Conclusion

- Caches and buffers definitely help in performance improvements.
- However, this comes at a cost of non-trivial program behaviors.
- These behaviors can be controlled using barriers, but only upto some point.
- The age-old tradeoff between performance and complexity in semantics to explain program execution can be seen firsthand here.
- Perhaps we should not crave for so much performance?
- Perhaps just improve the cache coherence protocol system?

Questions?