

Lecture 15 — Memory Consistency

Patrick Lam and Jeff Zarnett

2019-12-18

OpenMP Memory Model, Its Pitfalls, and How to Mitigate Them

OpenMP uses a **relaxed-consistency, shared-memory** model. This almost certainly doesn't do what you want. Here are its properties:

- All threads share a single store called *memory*—this store may not actually represent RAM.
- Each thread can have its own *temporary* view of memory.
- A thread's *temporary* view of memory is not required to be consistent with memory.

We'll talk more about memory models later. Now we're going to talk about the OpenMP model and why it's a problem.

Memory Model Pitfall. Consider this code.

```

a = b = 0
/* thread 1 */          /* thread 2 */

atomic(b = 1) // [1]     atomic(a = 1) // [3]
atomic(tmp = a) // [2]    atomic(tmp = b) // [4]
if (tmp == 0) then        if (tmp == 0) then
    // protected section    // protected section
end if                    end if

```

Does this code actually prevent simultaneous execution? Let's reason about possible states.

Order				t1 tmp	t2 tmp
1	2	3	4	0	1
1	3	2	4	1	1
1	3	4	2	1	1
3	4	1	2	1	0
3	1	2	4	1	1
3	1	4	2	1	1

Looks like it (at least intuitively).

Sorry! With OpenMP's memory model, no guarantees: the update from one thread may not be seen by the other.

Restoring Sanity with Flush. We do rely on shared memory working “properly”, but that's expensive. So OpenMP provides the **flush** directive.

```
#pragma omp flush [(list)]
```

This directive makes the thread's temporary view of memory consistent with main memory; it:

- enforces an order on the memory operations of the variables.

The variables in the list are called the *flush-set*. If you give no variables, the compiler will determine them for you.

Enforcing an order on the memory operations means:

- All read/write operations on the *flush-set* which happen before the **flush** complete before the flush executes.
- All read/write operations on the *flush-set* which happen after the **flush** complete after the flush executes.
- Flushes with overlapping *flush-sets* can not be reordered.

To show a consistent value for a variable between two threads, OpenMP must run statements in this order:

1. t_1 writes the value to v ;
2. t_1 flushes v ;
3. t_2 flushes v also;
4. t_2 reads the consistent value from v .

Let's revise the example again.

```

                a = b = 0
/* thread 1 */
atomic(b = 1)
flush(b)
flush(a)
atomic(tmp = a)
if (tmp == 0) then
    // protected section
end if

/* thread 2 */
atomic(a = 1)
flush(a)
flush(b)
atomic(tmp = b)
if (tmp == 0) then
    // protected section
end if
```

OK. Will this now prevent simultaneous access?

Well, no.

The compiler can reorder the `flush(b)` in thread 1 or `flush(a)` in thread 2. If `flush(b)` gets reordered to after the protected section, we will not get our intended operation.

Correct Example. We have to provide a list of variables to flush to prevent re-ordering:

```

                a = b = 0
/* thread 1 */
atomic(b = 1)
flush(a, b)
atomic(tmp = a)
if (tmp == 0) then
    // protected section
end if

/* thread 2 */
atomic(a = 1)
flush(a, b)
atomic(tmp = b)
if (tmp == 0) then
    // protected section
end if
```

Where There Is Implicit Flush:

- `omp barrier`
- at entry to, and exit from, **omp critical**;
- at exit from **omp parallel**;
- at exit from **omp for**;
- at exit from **omp sections**;
- at exit from **omp single**.

Where There's No Implicit Flush:

- at entry to **for**;
- at entry to, or exit from, **master**;
- at entry to **sections**;
- at entry to **single**;
- at exit from **for**, **single** or **sections** with a **nowait**
 - **nowait** removes implicit flush along with the implicit barrier

This is not true for OpenMP versions before 2.5, so be careful.

Final thoughts on flush. We've seen that it's very difficult to use flush properly. Really, you should be using mutexes or other synchronization instead of flush [Sue07], because you'll probably just get it wrong. But now you know what flush means.

Why Your Code is Slow

OpenMP code too slow? Avoid these pitfalls:

1. Unnecessary flush directives.
2. Using critical sections or locks instead of atomic.
3. Unnecessary concurrent-memory-writing protection:
 - No need to protect local thread variables.
 - No need to protect if only accessed in **single** or **master**.
4. Too much work in a critical section.
5. Too many entries into critical sections.

Example: Too Many Entries into Critical Sections.

```
#pragma omp parallel for
for (i = 0; i < N; ++i) {
    #pragma omp critical
    {
        if (arr[i] > max) max = arr[i];
    }
}
```

would be better as:

```
#pragma omp parallel for
for (i = 0 ; i < N; ++i) {
    #pragma omp flush(max)
    if (arr[i] > max) {
        #pragma omp critical
        {
            if (arr[i] > max) max = arr[i];
        }
    }
}
```

Memory Consistency, Memory Barriers, and Reordering

Now we'll talk a bit more about memory consistency, memory barriers and reordering in general. We'll start with instruction reordering by the CPU and move on to reordering initiated by the compiler. I'll also touch on some CPU instructions for atomic operations.

Memory Consistency. In a sequential program, you expect things to happen in the order that you wrote them. So, consider this code, where variables are initialized to 0:

```
T1: x = 1; r1 = y;  
T2: y = 1; r2 = x;
```

We would expect that we would always query the memory and get a state where some subset of these partially-ordered statements would have executed. This is the *sequentially consistent* memory model.

“... the result of any execution is the same as if the operations of all the processors were executed in some sequential order, and the operations of each individual processor appear in this sequence in the order specified by its program.” — Leslie Lamport

What are the possible values for the variables?

Another view of sequential consistency:

- each thread induces an *execution trace*.
- always, the program has executed some prefix of each thread's trace.

It turns out that sequential consistency is too expensive to implement. (Think how much coordination is needed to get a few people to agree on where to go for lunch; now try to get a group of people to agree on what order things happened in. Right. Now imagine it's a disagreement between threads so they don't have the ability to negotiate.) So most systems actually implement weaker memory models, such that both `r1` and `r2` might end up unchanged. Recall the **flush** example from last time.

Reordering. Compilers and processors may reorder non-interfering memory operations within a thread. For instance, the two statements in `T1` appear to be independent, so it's OK to execute them—or, equivalently, to publish their results to other threads—in either order. Reordering is one of the major tools that compilers use to speed up code.

When is reordering a problem? Spin locks, as we'll see below.

Memory Consistency Models

Here are some flavours of memory consistency models:

- Sequential consistency: no reordering of loads/stores.
- Sequential consistency for data-race-free programs: if your program has no data races, then sequential consistency.

- Relaxed consistency (only some types of reorderings):
 - Loads can be reordered after loads/stores; and
 - Stores can be reordered after loads/stores.
- Weak consistency: any reordering is possible.

In any case, **reorderings** only allowed if they look safe in current context (i.e. they reorder independent memory addresses). That can still be problematic, though.

Compilers and reordering. When it can prove that a reordering is safe with respect to the programming language semantics, the **compiler** may reorder instructions (so it's not just the hardware).

Reordering example. Say we want thread 1 to print a value set in thread 2.

```

                                f = 0

/* thread 1 */                  /* thread 2 */
while (f == 0) /* spin */;      x = 42;
printf("%d", x);                f = 1;
```

If thread 2 reorders its instructions, will we get our intended result? *No!*

Memory Barriers

We previously talked about OpenMP barriers: at a `#pragma omp barrier`, all threads pause, until all of the threads reach the barrier. Lots of OpenMP directives come with implicit barriers unless you add `nowait`.

A rather different type of barrier is a *memory barrier* or *fence*. This type of barrier prevents reordering, or, equivalently, ensures that memory operations become visible in the right order. A memory barrier ensures that no access occurring after the barrier becomes visible to the system, or takes effect, until after all accesses before the barrier become visible.

The x86 architecture defines the following types of memory barriers:

- `mfence`. All loads and stores before the barrier become visible before any loads and stores after the barrier become visible.
- `sfence`. All stores before the barrier become visible before all stores after the barrier become visible.
- `lfence`. All loads before the barrier become visible before all loads after the barrier become visible.

Note, however, that while an `sfence` makes the stores visible, another CPU will have to execute an `lfence` or `mfence` to read the stores in the right order.

Consider the example again:

```
                f = 0

/* thread 1 */      /* thread 2 */
while (f == 0) /* spin */;  x = 42;
// memory fence    // memory fence
printf("%d", x);      f = 1;
```

This now prevents reordering, and we get the expected result.

You can use the `mfence` instruction to implement *acquire barriers* and *release barriers*. An acquire barrier ensures that memory operations after a thread obtains the mutex doesn't become visible until after the thread actually obtains the mutex. The release barrier similarly ensures that accesses before the mutex release don't get reordered to after the mutex release. Note that it is safe to reorder accesses after the mutex release and put them before the release.

Preventing Memory Reordering in Programs: Compiler Barriers. First: Don't use `volatile` in C/C++ on variables [Inc08]. Remember that the `volatile` keyword is supposed to tell the compiler not to put the value in a register. That does NOT make it a synchronization construct. If you use the correct synchronization primitives, you will get the behaviour you want. However, you can prevent reordering using compiler-specific calls.

- Microsoft Visual Studio C++ Compiler:

```
_ReadWriteBarrier()
```

- Intel Compiler:

```
__memory_barrier()
```

- GNU Compiler:

```
__asm__ __volatile__ ("" ::: "memory");
```

The compiler also shouldn't reorder across e.g. Pthreads mutex calls.

Aside: gcc Inline Assembly. Just as an aside, here's gcc's inline assembly format

```
__asm__ ( assembler template
        : output operands          /* optional */
        : input operands           /* optional */
        : list of clobbered registers /* optional */
        );
```

Note that we've just seen `__volatile__` with `__asm__`. This isn't the same as the normal C `volatile`. It means:

- The compiler may not reorder this assembly code and put it somewhere else in the program.

Back to Memory Reordering in Programs. Fortunately, an OpenMP **flush** (or, better yet, mutexes) also preserve the order of variable accesses. That is, it prevents reordering from both the compiler and hardware. For GNU, flush is available as `__sync_synchronize()`;

`volatile`. This qualifier ensures that the code does an actual read from a variable every time it asks for one (i.e. the compiler can't optimize away the read). It does not prevent re-ordering nor does it protect against races.

Note: proper use of memory fences makes `volatile` not very useful (again, `volatile` is not meant to help with threading, and will have a different behaviour for threading on different compilers/hardware).

C/C++11 Memory Model

We have talked about memory models in the context of OpenMP. Let's talk about the core languages now—that is, C and C++ [Lov14, BA08]—when not using OpenMP.

What outputs are possible from this example?

Thread 1:	Thread 2:
<code>foo = 7;</code>	<code>printf("%d\n", foo);</code>
<code>bar = 42;</code>	<code>printf("%d\n", bar);</code>

You might think “undefined”, but actually it's worse than that. The C11 and C++11 language definitions don't even say what a thread is. Of course, there is Pthreads, but that is a library, not the language itself. So you can't even ask this question when talking about pre-C11/C++11 versions of C and C++. Well. Okay. You can ask, but the question makes no sense. Sort of like putting your hand up in class and saying “Where did you bury your oranges?”. It's a syntactically valid question and it describes something that's possible, but it still makes no sense.

The older standards made no reference to any kind of CPU, memory architecture, cache strategy, or anything like that. Which on the one hand is nice and general, but on the other hand, leads to problems. The “abstract machine” that the C and C++ standards refer to is inherently single threaded, making it actually impossible to write a portable multithreaded C or C++ program [Lov14]. The part that's impossible is that word “portable” – people write multithreaded C and C++ programs all the time (in this class, even) but they have system specific code and implementation-defined behaviour. Open up a pthreads library or some equivalent and sure enough, you will find something architecture specific in there.

C++11 (and C11) have improved the situation, though. There is actually a memory model (based on an abstract machine) and threading primitives such as mutexes, atomics, and memory barriers—the concepts that we have seen in this course. Now there are rules! Yes. We like rules. Okay. I [JZ], at least, like rules. C++11 defines how a compiler can generate code that accesses memory even when there is concurrency. There are also standard mutex operations and atomics and barriers and all those lovely things.

Now, we can ask the question about the behaviour of the above example. It does have undefined behaviour, since there is contended access to the variables `foo` and `bar`. How can we fix that?

Atomics. A good exam question: if `foo` is atomic, what are the possible outputs?

Thread 1:	Thread 2:
<code>foo.store(7);</code>	<code>printf("%d\n", foo.load());</code>
<code>bar.store(42);</code>	<code>printf("%d\n", bar.load());</code>

Alright, we have some defined behaviour now. Honestly, it depends how these things are scheduled, but the answer is one of the following set: {0/0, 7/42, 7/0, 0/42}. The answer depends on how they are interleaved. But at least we get some certainty that the output will be one of those four things and there's no chance of garbage because the print takes place during an assignment operation.

We probably still don't like this because we don't have mutual exclusion here and we can get several different answers, some of which are probably “wrong” (for whatever definition of a correct answer is), but at least our set of potential wrong answers is smaller. So that's a start. Compilers have to follow the new rules in generating code, so their output will behave as if the architecture followed the standard memory model. That's something.

References

[BA08] Hans J. Boehm and Sarita V. Adve. Foundations of the c++ concurrency memory model, 2008. Online; accessed 16-December-2015. URL: <http://rsim.cs.illinois.edu/Pubs/08PLDI.pdf>.

- [Inc08] Stack Exchange Inc. Using C/Pthreads: do shared variables need to be volatile?, 2008. Online; accessed 14-December-2015. URL: <http://stackoverflow.com/questions/78172/using-c-pthreads-do-shared-variables-need-to-be-volatile>.
- [Lov14] Robert Love. How are the threading and memory models different in c++ as compared to c?, 2014. Online; accessed 16-December-2015. URL: <http://www.quora.com/C++-programming-language/How-are-the-threading-and-memory-models-different-in-C++-as-compared-to-C>.
- [Sue07] Michael Suess. Please don't rely on memory barriers for synchronization, 2007. Online; accessed 12-December-2015. URL: <http://www.thinkingparallel.com/2007/02/19/please-dont-rely-on-memory-barriers-for-synchronization/>.