

Contents of Lecture 3

- The need for memory consistency models
- The uniprocessor model
- Sequential consistency
- Relaxed memory models
- Weak ordering
- Release consistency

The Need for Memory Consistency Models

- In this lecture we assume all threads are executing on different CPUs. This makes the presentation simpler so that we can use the words "processor" and "thread" interchangeably.
- A memory consistency model is a set of rules which specify when a written value by one thread can be read by another thread.
- Without these rules it's not possible to write a correct parallel program. We must reason about the consistency model to write correct programs.
- The memory consistency model also affects which programmer/compiler and hardware optimizations are legal.

Who Must Know About the Consistency Model?

- The memory consistency model affects
 - System implementation: hardware, OS, languages, compilers
 - Programming correctness
 - Performance
- There are several different consistency models.
- One lets us use normal variables for synchronization, while most don't.
- Instead we must use some form of lock/unlock calls to ensure that our threads will see the data most recently written by another thread.

Uniprocessor Memory Consistency

- To the programmer, the last written value to a memory location will always be the value read from that same location.
- This is the contract to the programmer.
- Most programming languages, including C99, are specified using this model for single threaded programs. In Lecture 6 we will look at the new C11 standard which specifies it's memory model with support for multiple threads and atomic variables.
- Why is the following program invalid?

```
a = 1;  
a = ++a + a--;  
printf("a = %d\n", a);
```

- In C/C++ there is a rule which says a variable may only be modified once between two sequence points (e.g. semicolons) so for C/C++ the uniprocessor model has additional constraints.

Effects of the Uniprocessor Model

- We can as programmers reason about which values are legal (provided we know about the additional rules in C/C++).
- Performance optimizations are possible both for optimizing compilers and processor designers.
- Compilers can put a variable in a register as long as it wishes and can produce code which behaves **as if** the variable was accessed from memory.
- Compilers can thus swap the order of two memory accesses if it can determine that it will not affect program correctness under the uniprocessor model.
- A superscalar processor can execute two memory access instructions in any order as long as they don't modify the same data.
- A processor may reorder memory accesses e.g. due to lockup-free caches, multiple issue, or write buffer bypassing.

Lockup-Free Caches

- Suppose a cache miss occurs.
- In a lockup-free cache other subsequent cache accesses can be serviced while a previous cache miss is waiting for data to arrive.
- **A lockup-free cache can reorder memory accesses.**
- This is useful in superscalar processors since execution can proceed, until the next instruction to execute actually needs the data from the cache miss.
- Another name for lockup-free caches is non-blocking caches.
- Data prefetching fetches data before it is needed and loads it into the cache (but not into a processor register — but the compiler may do so if it wishes anyway).
- Data prefetching requires the cache to be lockup-free — otherwise the processor would be stalled and the prefetching would be rather pointless.

Write Buffer Bypassing

We want our reads to be serviced as quickly as possible

- Between the L1 cache and L2 cache is a buffer, called the first level write buffer.
- A second level write buffer is located between the L2 cache and the bus interface.
- By letting read misses bypass writes in the buffer to other addresses, the reads can be serviced faster.
- **Such bypasses reorders the memory accesses.**
- A superscalar processor has queues in the load-store functional unit where accesses also can be reordered.

Sequential Consistency

- Sequential consistency (SC) was published by Leslie Lamport in 1979 and is the simplest consistency model.
- Neither Java, Pthreads, nor C11/C++ require it. They work on relaxed memory models.
- Sequential consistency can be seen from the programmer as if the multiprocessor has no cache memories and all memory accesses go to memory, which serves one memory request at a time.
- This means that program order from one processor is maintained and that all memory accesses made by all processors can be regarded as atomic (i.e. not overlapping).

Definition of SC

- Lamport's definition: *A multiprocessor system is sequentially consistent if 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.*
- Consider the program execution:

```
int A = B = C = 0;
```

```
T1:  
A = 1;
```

```
T2:  
  
if (A)  
    B = 1;
```

```
T3:  
  
  
if (B)  
    C = A;
```

- Since all memory accesses are atomic, writes are seen in the same order so T3 must read the value 1 when reading A.

Dekker's Algorithm

```
bool    flag[2] = { false, false };
int     turn = 0;

void work(int i)
{
    for (;;) {
        flag[i] = true;
        while (flag[!i]) {
            if (turn != i) {
                flag[i] = false;
                while (turn != i)
                    ;
                flag[i] = true;
            }
        }

        /* critical section */

        /* leave critical section */

        turn = !i;
        flag[i] = false;
    }
}
```

- SC ensures that Dekker's algorithm works.

Implementing SC in a System Without Caches

- We will examine the effects of each of the following hardware optimizations:
 - Write buffer with read bypass
 - Overlapping writes
 - Non-blocking reads
- As we will see, none of these can be used even if there are no caches.
- This reduces performance considerably.

Write Buffer with Read Bypass

- Assume a bus-based multiprocessor.
- Since there are no caches, the write buffer, a FIFO queue, sits between the CPU and the bus interface.
- With read bypass, it is thus meant that a read skips the queue in the buffer and goes first to the bus before any write (to a different address).
- In Dekker's algorithm both CPUs can set their `flag[i]` to true and put that write into its write buffer.
- Then the reading of the other thread's flag will bypass the write in the write buffer.
- When bypassing the old values of the `flag[!i]` can be returned (e.g. if there were other writes before in the buffers) and both can enter the critical section!

The read bypass destroys the atomicity and hence the sequential order.

Overlapping Writes

- In a bus-based system, the FIFO write buffer queue ensures that all writes from the CPU are ordered.
- In a general topology, however, different nodes typically are located at different distances and writes easily can arrive in an order different from the program order.
- In the example with variables A, B, and C, the new value of B may reach T3 before A does which violates SC.

T2 should not be allowed to start its write to b before T3 has become aware of the write to a.

Non-Blocking Reads

- Consider the following example from the first lecture again.
- With speculative execution and non-blocking reads, T2 can read the value of `a` before it leaves the while-loop, which violates SC.

```
int a, f;
```

```
// called by T1
```

```
void v(void)
```

```
{
```

```
    a = u();
```

```
    f = 1;
```

```
}
```

```
// called by T2
```

```
void w(void)
```

```
{
```

```
    while (!f)
```

```
        ;
```

```
    printf("a = %d\n", a);
```

```
}
```

Implementing SC in a System With Caches

- The three issues in systems without caches can violate SC also with caches.
- For example a read cache hit must not be allowed to precede a previous read miss.
- In addition, since there can be multiple copies of a variable, there must be a mechanism which controls e.g. whether a cached copy is allowed to be read — it is not if another processor just has modified it, for example.
- This mechanism is called a *cache coherence protocol*.
- A cache coherence protocol has three main tasks, as we will see next.

Cache Coherence Protocols

- ① At a write, the cache coherence protocol should either remove all other copies, including the memory copy, or send the newly written data to update each copy.
 - A protocol which uses the former technique is called a **write invalidate protocol** while the latter is called a **write update protocol**.
 - Which is best depends on the sharing behaviour of the application but write invalidate is almost always better. More details follow!
- ② Detecting when a write has completed so that the processor can perform the next memory access.
- ③ Maintaining the illusion of atomicity — with memory in multiple nodes the accesses cannot be atomic but a SC machine must behave as if they are.

Detecting Write Completion

- Consider a write to a memory location which is replicated in some caches.
- How can the writing processor know when it's safe to proceed?
- The write request is sent to the memory where the data is located.
- The memory knows which caches have a copy (recall from the previous lecture this information is stored in a directory, e.g. as a bit vector).
- The memory then sends either updates or invalidations to the other caches.
- The receiving caches then must acknowledge they have received the invalidation message from memory.
- The acknowledgement is typically sent to the memory and then when all acknowledgements have been received, a message is sent to the writing processor (actually, its cache) to tell it the write has completed.
- After that, the processor can proceed.

Write Atomicity 1(2)

- There are two different problems:
 - (1) Write atomicity for a particular memory location, i.e. ensuring all CPUs see the same sequence of writes to a variable.
 - (2) Write atomicity and reading the modified value.
- For (1) it is sufficient to ensure that writes to the same memory location are serialized, but not for (2). See next page.
- The memory controller can easily enforce serialization.
- Assume writes from two different caches are sent to it.
- One of them must arrive first. The other can simply be replied to with a negative acknowledgement of "try again later!"
- When the cache receives that message it will simply try again and after a while it will be its turn.

Write Atomicity 2(2)

- Let us now consider (2): reading the modified value.
- A write has completed when all CPUs with a copy have been notified.
- However, if one CPU is allowed to read the written value before the write has completed, SC can be violated.

```
int A = B = C = 0;
```

```
T1:  
A = 1;
```

```
T2:  
if (A)  
    B = 1;
```

```
T3:  
  
if (B)  
    C = A;
```

- Assume all variables are cached by all threads, and T2 reads A before T3 knows about the write.
- Then T2 can write to B which might be so close to T3 that T3 can read A from its cache before the invalidation of A reaches T3.
- The solution is to disallow any CPU from reading A before the write is complete, which can be implemented in write invalidate as for case (1).

Write Atomicity in Write Update Protocols

- Recall that in a write update protocol, instead of invalidating other copies, new values are sent to the caches which replicate the value.
- So A is sent to T2 and to T3.
- While it's tempting to read A for T2 it's not permitted to.
- In write update, the updates are done in two phases. First is the data sent, and all CPUs acknowledge they have received it. Second each CPU is sent a message that it may read the new data.
- There are other problems with write update as well, for instance updates may be sent to CPUs which no longer are interested in the variable, thus wasting network traffic.

Optimizing Compilers and Explicitly Parallel Codes

- We have now seen the restrictions on hardware so that it does not reorder memory accesses and thus violate SC.
- The same restrictions must of course be put on optimizing compilers.
- The compiler must preserve "source code order" of all memory accesses to variables which may be shared — but not the order of stack accesses or other data known to be private to the thread.
- Examples of optimizations which cannot (in general) be used:
 - Register allocation
 - Code motion out of loops
 - Loop reordering
 - Software pipelining
- It's easy to imagine that these restrictions will slow down SC.
- The solution has often been to compile code for uniprocessors and use the `volatile` qualifier for shared data.
- Recall that `volatile` in C is different from `volatile` in Java!

Parallelizing Compilers

- If the compiler is doing the parallelization, these restrictions don't apply since the compiler writer hopefully knows what he or she is doing!
- After this course, however, you will probably not have too high hopes for automatic parallelization, except for numerical codes.
- In my view, parallelization of "normal" programs needs so drastic changes to the source code that automatic tools hardly can do that very well.
- But I may be wrong of course!

Cache Coherence Protocol States

- The cache coherence protocol maintains a state for each cache and memory block.
- The cache state can for instance be:
 - SHARED
 - INVALID
 - EXCLUSIVE — memory and this cache has a copy but it's not yet modified.
 - MODIFIED — only this cache has a copy and it's modified
- There are similar states for a memory block and also the bit-vector with info about which cache has a copy.
- In addition, a memory block can be in a so called **transient state** when acknowledgements are being collected, for instance.

Memory Access Penalty

- The time the processor is stalled due to waiting for the cache is called the memory access penalty.
- Waiting for read cache misses is difficult to avoid in any computer with caches (reminder: the Cell's SPU's have no caches...)
- Waiting for obtaining exclusive ownership of data at a write access is one of the disadvantages for SC
- How significant it is depends on the application
- Of course, once the CPU owns some data, it can modify it's cached copy without any further write access penalty, until some other CPU also wants to access that data, in which case the state becomes SHARED again.

Optimizing SC in Hardware

- Data prefetch can either fetch data in shared or exclusive mode
- By prefetching data in exclusive mode, the long delays of waiting for writes to complete can possibly be reduced or eliminated.
- Since exclusive mode prefetching invalidates other copies it can also increase the cache miss rate.
- Somewhat more complex cache coherence protocols can monitor the sharing behaviour and determine that it probably is a good idea to grant exclusive ownership directly instead of only a shared copy which is then likely followed by an ownership request.
- In superscalar processors it can be beneficial to permit speculative execution of memory reads. If the value was invalidated, the speculatively executed instructions (the read and subsequent instructions) are killed and the read is re-executed.

Optimizing SC in the Compiler

- Another approach is to have a memory read instruction which requests ownership as well.
- This can be done easily in optimizing compilers but needs a new instruction.
- It's very useful for data which moves from different caches and is first read and then written:

```
p->a += 1;  
p->b += 1;  
p->c += 1;
```

- Here the compiler can very easily determine that it's smarter to request ownership while doing the read.

Summary of Sequential Consistency

- Recall the two requirements of SC:
 - ① Program order of memory accesses
 - ② Write atomicity
- While SC is "nice" since it's easy to think about, it has some serious drawbacks:
 - The above two requirements... :-)
 - ...which limit compiler and hardware optimizations, and...
 - introduce a write access penalty
- The write access penalty is due to the processor cannot perform another memory access before the previous has completed.
- This is most notable for writes, since read misses are more difficult to optimize away by any method
- We will next look at relaxed memory consistency models

Relaxed Memory Models

- Relaxed memory models do not make programming any more complex (many would disagree however).
- However, you need to follow additional rules about how to synchronize threads.
- C11, Pthreads and Java are specified for relaxed memory models.
- In relaxed memory models, both the program order of memory references and the write atomicity requirements are removed and are replaced with different rules.
- For compiler writers and computer architects this means that more optimizations are permitted.
- For programmers it means two things:
 - 1 you must protect data with special system recognized synchronization primitives, e.g. locks, instead of normal variables used as flags.
 - 2 your code will almost certainly become faster, perhaps by between 10 and 20 percent, due to eliminating the write access penalty.

System Recognized Locks

- Under SC, you normally must protect your data using locks to avoid data races
- However, there are programs in which data races are acceptable
- Data races are forbidden in C11 and result in undefined behaviour.
- Under SC you can write your own locks by spinning on a variable `int flag` as you wish.
- Under relaxed memory models you should use whatever the system provides.

Pthreads in Linux

- On UNIX the system recognized locks normally mean Pthreads mutexes.
- On Linux Pthreads are implemented with a low-level lock called a Futex — fast mutex.
- These have the advantage of not involving the Linux kernel in the case when there is no contention for a lock
- If many threads want a lock and must wait, the futex mechanism will ask the kernel to make a context switch to a different thread.
- Short time busy waiting for a release of a lock is better than context switching but if the predicted waiting time is long enough to justify a context switch, that will be performed if you use a futex.

Write Your Own Locks!?!?

- You should not need to, but if you wish, you can write your own locks.
- That is not so difficult on SC.
- As we will see shortly, to do so on relaxed memory models requires the use of special **synchronization machine instructions** which must be executed as part of the lock and unlock.
- Use e.g. inlined assembler to specify them.

Relaxed Memory Models

- Recall that relaxed memory models relax the two constraints of memory accesses: program order and write atomicity.
- There are many different relaxed memory models and we will look only at a few.
- We have the following possibilities of relaxing SC.
- Relaxing A to B program order: we permit execution of B before A.
 - ① write to read program order to different addresses
 - ② write to write program order to different addresses
 - ③ read to write program order to different addresses
 - ④ read to read program order to different addresses
 - ⑤ read other CPU's write early
 - ⑥ read own write early
- Different relaxed memory models permit different subsets of these.

Assumptions

- All writes will eventually be visible to all CPUs.
- All writes are serialized (recall: this can be done at the memory by letting one write be handled at a time — other pending writes must be retried).
- Uniprocessor data and control dependences are enforced.

Relaxing the Write to Read Program Order Constraint

- Obviously different addresses are assumed!
- A read may be executed before a preceding write has completed.
- With it, Dekker's Algorithm fails, since both can enter the critical section. *How?*

```
bool    flag[2] = { false, false };
int     turn = 0;

void work(int i)
{
    for (;;) {
        flag[i] = true;
        while (flag[!i]) {
            if (turn != i) {
                flag[i] = false;
                while (turn != i)
                    ;
                flag[i] = true;
            }
        }

        /* critical section */

        turn = !i;
        flag[i] = false;

        /* not critical section */
    }
}
```

Some Models Which Permit Reordering a (Write, Read) Pair

- Processor consistency, James Goodman
- Weak Ordering, Michel Dubois
- Release Consistency, Kourosh Gharachorloo
- IBM 370, IBM Power
- Sun TSO, total store ordering
- Sun PSO, partial store ordering
- Sun RMO relaxed memory order
- Intel X86

Additionally Relaxing the Write to Write Program Order Constraint

- Recall the program below.

```
int a, f;
```

```
// called by T1
```

```
void v(void)
```

```
{
```

```
    a = u();
```

```
    f = 1;
```

```
}
```

```
// called by T2
```

```
void w(void)
```

```
{
```

```
    while (!f)
```

```
    ;
```

```
    printf("a = %d\n", a);
```

```
}
```

- By relaxing the write to write program order constraint, the write to `f` may be executed by T_1 even before the function call to `u`, resulting in somewhat unexpected output.

Some Models Which Permit Reordering a (Write, Write) Pair

- Weak Ordering
- Release Consistency
- IBM Power
- Sun PSO, partial store ordering
- Sun RMO relaxed memory order

Relaxing All Memory Ordering Constraints

- The only requirements left is the assumption of uniprocessor data and control dependences.
- Models which impose no reordering constraints for normal shared data include:
 - Weak Ordering
 - Release Consistency
 - IBM Power
 - Sun RMO relaxed memory order
- These consistency models permit very useful compiler and hardware optimizations and both Java and Pthreads (and other platforms) require from the programmer to understand and use them properly!
- In the preceding example, the two reads by T_2 are allowed to be reordered.
- The obvious question then becomes: how can you write a parallel program with these memory models???

What do you say?

Special Machine Instructions for Synchronization

- The short answer is that machines with relaxed memory models have special machine instructions for synchronization.
- Consider a machine with a **sync** instruction with the following semantics:
 - When executed, all memory access instructions issued **before** the sync must complete before the sync may complete.
 - All memory access instructions issued **after** the sync must wait (i.e. not execute) until the sync has completed.
 - Assume both T_1 and T_2 have cached a .

```
int a, f;
```

```
// called by T1
```

```
void v(void)
```

```
{
```

```
    a = u();
```

```
    asm("sync");
```

```
    f = 1;
```

```
}
```

```
// called by T2
```

```
void w(void)
```

```
{
```

```
    while (!f)
```

```
    ;
```

```
    asm("sync");
```

```
    printf("a = %d\n", a);
```

```
}
```

- With `asm` we can insert assembler code with gcc and most other compilers.
- The sync instructions are required, as explained next...

Sync Instruction Example 1(3)

```
int a, f;

// called by T1
void v(void)
{
    a = u();
    asm("sync");
    f = 1;
}

// called by T2
void w(void)
{
    while (!f)
        ;
    asm("sync");
    printf("a = %d\n", a);
}
```

- The write by T_1 to `a` results in an invalidation request being sent to T_2 .
- At the sync, T_1 must wait for an acknowledgement from T_2 .
- When T_2 receives the invalidation request, it acknowledges it directly and then puts it in a queue of incoming invalidations.
- When T_1 receives the acknowledgement, the write is complete and the sync can also complete, since there are no other pending memory accesses issued before the sync.

Sync Instruction Example 2(3)

```
int a, f;
```

```
// called by T1
```

```
void v(void)
```

```
{
```

```
    a = u();
```

```
    asm("sync");
```

```
    f = 1;
```

```
}
```

```
// called by T2
```

```
void w(void)
```

```
{
```

```
    while (!f)
```

```
        ;
```

```
    asm("sync");
```

```
    printf("a = %d\n", a);
```

```
}
```

- The write by T_1 to `f` also results in an invalidation request being sent to T_2 .
- When T_2 receives the invalidation request, it acknowledges it directly and then puts it in the queue of incoming invalidations.
- T_2 is spinning on `f` and therefore requests a copy of `f`.
- When that copy arrives, with value one, it must wait at the sync until the invalidation to `a` has been applied — it might already have been.
- Without the sync by T_2 its two reads could be reordered:
 - The compiler could have put `a` in a register before the while-loop.
 - The CPU could speculatively have read `a` from memory.
 - The incoming transactions may be reordered in a node.

Sync Instruction Example 3(3)

- Instead of this sync instruction used in order to be concrete, we can use the Linux kernel memory barrier:

```
int a, f;

// called by T1
void v(void)
{
    a = u();
    smp_mb();
    f = 1;
}

// called by T2
void w(void)
{
    while (!f)
        ;
    smp_mb();
    printf("a = %d\n", a);
}
```

- The memory barrier is a macro which will expand to a suitable machine instruction.

Weak Ordering

- The memory consistency model introduced with the sync instruction is called Weak Ordering, WO, and was invented by Michel Dubois.
- The key ideas for why it makes sense are the following:
 - Shared data structures are modified in critical sections.
 - Assume N writes to shared memory are needed in the critical section.
 - In SC the processor must wait for **each** of the N writes to complete in sequence.
 - In WO, the processor can pipeline the writes and only wait at the end of the critical section.
 - Sync instructions are then executed as part of both the lock and unlock calls.
- Of course, some or all of the N writes may be to already owned data in which case there is no write penalty.
- Measurements on different machines and applications show different values but 10-20 % percent of the execution time can be due to writes in SC.

Release Consistency

- Release Consistency, RC, is an extension to WO, and was invented for the Stanford DASH research project.
- Two different synchronization operations are identified.
- An **acquire** at a lock.
- A **release** at an unlock.
- An acquire orders all subsequent memory accesses, i.e. no read or write is allowed to execute before the acquire. Neither the compiler nor the hardware may move the access to before the acquire, and all acknowledged invalidations that have arrived before the acquire must be applied to the cache before the acquire can complete (i.e. leave the pipeline).
- A release orders all previous memory accesses. The processor must wait for all reads to have been performed (i.e. the write which produced the value read must have been acknowledged by all CPUs) and all writes made by itself must be acknowledged before the release can complete.

- Recall: by $A \rightarrow B$ we mean B may execute before A
- WO relaxation: data \rightarrow data
- RC relaxation:
 - data \rightarrow data
 - data \rightarrow acquire
 - release \rightarrow data
 - release \rightarrow acquire
- acquire \rightarrow data: forbidden
- data \rightarrow release: forbidden
- In practice not very significant difference on performance.

Succeeding in Programming with Relaxed Memory Models

- The bad news are that an execution might produce correct output 99% of the time and suffer painful data races only rarely.
- Recall Helgrind can only detect when things are wrong but not prove the code cannot have any data race!
- If you use system recognized locks such as Pthreads mutexes or Java synchronization you are likely to have more fun.