Lecture 08: Introduction to Memory Consistency

Vivek Kumar Computer Science and Engineering IIIT Delhi vivekk@iiitd.ac.in

CSE513: Parallel Runtimes for Modern Processors



- Reducing concurrent accesses to private deque
- Automatically controlling task granularity

Today's Class

- ► Memory consistency problem
 - Difference between coherence and consistency
 - Sequential consistency



Multicore System Overview



- Cache controller copes code/data to/from main memory
 - It runs the actual cache coherence algorithm
- LLC is shared with all the cores
- LLC is logically a memory side cache as it primarily serves to reduce the latency and increase bandwidth of memory accesses

Physical VS. Logical Cores



Lecture 08: Introduction to Memory Consistency

The Cache Coherence Problem



CSE513: Parallel Runtimes for Modern Processors

Cache Coherence (CC)

- CC protocols gives the illusion of an atomic memory system without any caches
- CC definition / requirements
 - Single-Writer-Multiple-Reader (write serialization)
 - At any given instant only one core can WRITE or WRITE+READ a given memory location, or multiple cores can only READ that location
 - Value of a memory location is *eventually* propagated to all readers (value propagation)

Epoch=1	Epoch=2	Epoch=3	Epoch=4	
Read (X) Core-1 & 2	Write (X) Core-1	Read (X) Core-3 & 4	Read-Write (X) Core-2	

Dividing a memory location's lifetime into epochs



CSE513: Parallel Runtimes for Modern Processors

Question

- Assume that for each memory locations in a CC system, there are N tokens, where N is the total number of cores
 - How to define memory access rule for a memory location X using these tokens?
 - If a core has all of the tokens, it may write the memory location. IF a core has one or more tokens, it may read the memory location



Multithreaded Programs

Initially A = B = 0

Thread 1Thread 2A = 4D = 4

$$\begin{array}{ll} A = 1 & B = 1 \\ \text{if } (B == 0) & \text{if } (A == 0) \\ \text{print "Hello";} & \text{print "World";} \end{array}$$

Question: What can be printed: "Hello", "World", or "Hello World"?



Why Memory Consistency Model?

- What does programmer think?
 - Cores atomically executes one instruction at a time in program order
 - Valid over a single core (sequential) processor



Why Memory Consistency Model?

- What does programmer think?
 - Cores atomically executes one instruction at a time in **program order**
 - Valid over a single core (sequential) processor

• What cores can do?

- Reorder memory operations to reduce memory access latency
 - Read-Read
 - Write-Write
 - Read-Write
 - Write-Read
- Memory consistency model defines what it means to "read a memory", "write at a memory", and the "respective ordering"



Today's Class

- Memory consistency problem
- Difference between coherence and consistency
 - Sequential consistency



Typical Memory Hierarchy



- Hiding Read (Load) latency
 - Prefetching loads that are delayed due to consistency model
 - Out-of-order execution of loads
 - Speculative execution to allow the processor to proceed even though consistency model requires memory accesses be delayed (e.g., Intel, AMD)
 - Further relaxing the memory consistency model (e.g., ARM, IBM)
- Hiding Write (Store) latency
 - Store Buffers available on almost all modern processors (more on this later)

CSE513: Parallel Runtimes for Modern Processors



- Memory latency continues to limit the performance of modern out-of-order cores
 - Loads and Stores are very expensive
 - Instruction pipelining is one simple technique for improving memory latency

Consistency VS. Coherence (Example)



Coherence assures that values written by one processor are made visible to other processors

Each student has a stale copy of shared variable "syllabus". They will **eventually** get the updated copy of "syllabus" due to coherence's value propagation property



CSE513: Parallel Runtimes for Modern Processors

Consistency VS. Coherence (Example)



Consistency VS. Coherence





Language v/s Hardware Memory Model



CSE513: Parallel Runtimes for Modern Processors

Memory Ordering at Different Stages



Picture (1) source: https://preshing.com/images/hardware-matters.png



CSE513: Parallel Runtimes for Modern Processors

Today's Class

- Memory consistency problem
- Difference between coherence and consistency
- Sequential consistency



- Programmer's productivity is high if the memory operations in his parallel program executes in the program order
 - Sequential consistency
 - Lamport defined SC (1979)
- Let's assume there is a processor that supports sequential consistency to provide high productivity for the programmer
 - Although it suffers on performance
- What that processor has to do for supporting sequential consistency?





















Summary: Sequential Consistency

- Each core execute its read/write in the program order
 Irrespective of whether they are at same or different addresses
- Every read from address A gets its value from the last write before it to the same address A
 - Operations must appear atomic—one operation must complete before the next one is issued



Don't start the next operation until the previous completes



The Problem with SC



- With a single view of memory, we can't run (2) until (1) has become visible to every other thread
 - On a modern CPU, that's a very expensive operation due to the cache hierarchy
- The only shared memory between the two cores is all the way back at the L3 cache, which often takes upwards of 90 cycles to access

CSE513: Parallel Runtimes for Modern Processors Credits:

Credits: James Bornholt, UW, CSE451

Sequential Consistency Supported!

- Is there support for sequential consistency in modern programming language?
 - Yes, **but** only for **one** particular case in almost every language!
 - Only for code block that is Data Race Free (DRF)
 - Even without using any special construct, you get it for free inside critical sections that are executed within mutex lock/unlock operations
 - We will see in lecture 10 about other features in C++11 to support DRF
 - No sequential consistency for rest of the program (racy code!)



Cores executing synchronized code blocks

o Can be seen as a "switch" running one instruction at a time



Cores executing synchronized code blocks



Cores executing synchronized code blocks



Cores executing synchronized code blocks



Cores executing synchronized code blocks

o Can be seen as a "switch" running one instruction at a time



Cores executing synchronized code blocks

o Can be seen as a "switch" running one instruction at a time



Cores executing synchronized code blocks



Reference Materials

https://www.youtube.com/watch?v=mDYLRX2pbFw

• Preshing

<u>https://preshing.com/20120515/memory-reordering-caught-in-the-act/</u>



Next Lecture (L #09)

- Hardware memory consistency model
- Extra lecture on Saturday (24/09) at 2.30pm to compensate for missing regular lecture of 30/09

