

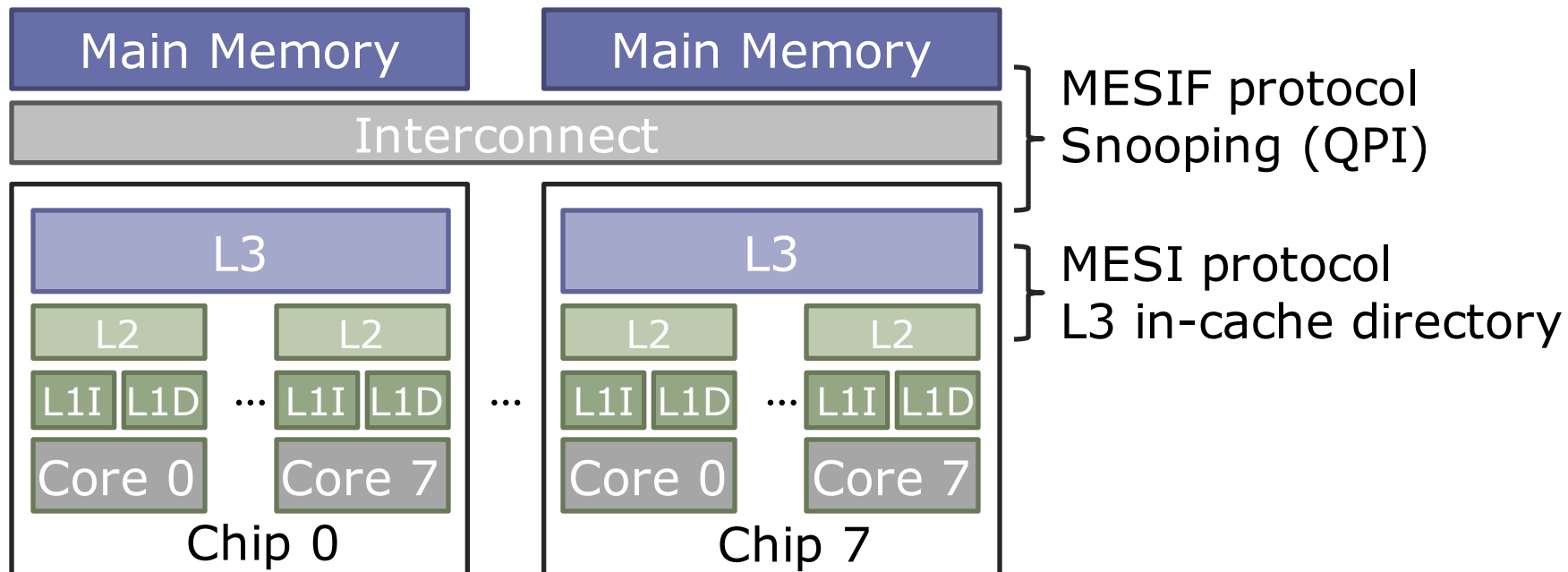
Consistency and Relaxed Memory Models

Joel Emer

Computer Science and Artificial Intelligence Lab
M.I.T.

Coherence in Multi-Level Hierarchies

- Can use the same or different protocols to keep coherence across multiple levels
- Key invariant: Ensure sufficient permissions in all intermediate levels
- Example: 8-socket Xeon E7 (8 cores/socket)



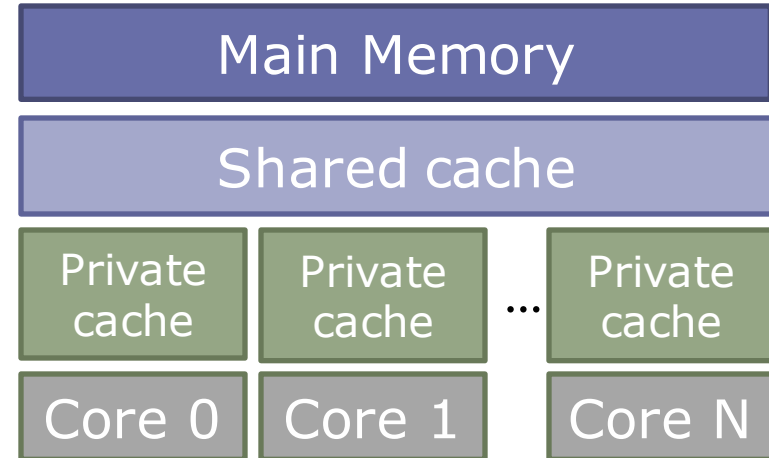
In-Cache Directories

- Common multicore memory hierarchy:

- 1+ levels of private caches
- A shared last-level cache
- Need to enforce coherence among private caches

- Idea: Embed the directory information in shared cache tags

- Shared cache must be inclusive



✓ Avoids tag overheads & separate lookups

✗ Can be inefficient if shared cache size \gg sum(private cache sizes)

Coherence vs Consistency

- Coherence: What values can a read return?
 - Concerns reads/writes to a single memory location
- Consistency: When do writes become visible to reads?
 - Concerns reads/writes to multiple memory locations

Why Consistency Matters

Initial memory contents

a: 0

flag: 0

Processor 1

Store (a), 10;

Store (flag), 1;

Processor 2

L: Load r1, (flag);

if r₁ == 0 goto L;

Load r2, (a);

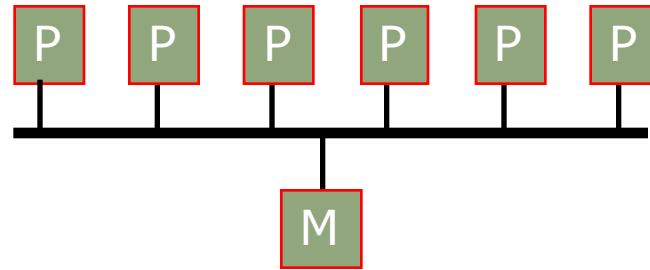
- What value does processor 1's r2 hold after both processors finish running this code?

It depends on the order in which processor 2 observes processor 1's stores!

10 if Store (flag) > Store (a); 0 or 10 otherwise

Sequential Consistency

A Straightforward Memory Model



“ A 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 the order specified by the program”

Leslie Lamport

Sequential Consistency =
arbitrary *order-preserving interleaving*
of memory references of sequential programs

Sequential Consistency

<p><i>Processor 1</i></p> <p>Store (a), 10;</p> <p>Store (flag), 1;</p>	<p><i>Processor 2</i></p> <p>L: Load r1, (flag);</p> <p>if r₁ == 0 goto L;</p> <p>Load r2, (a);</p>
---	--

initially flag = 0

- In-order instruction execution
- Atomic loads and stores

SC is easy to understand but architects and compiler writers want to violate it for performance

Memory Model Issues

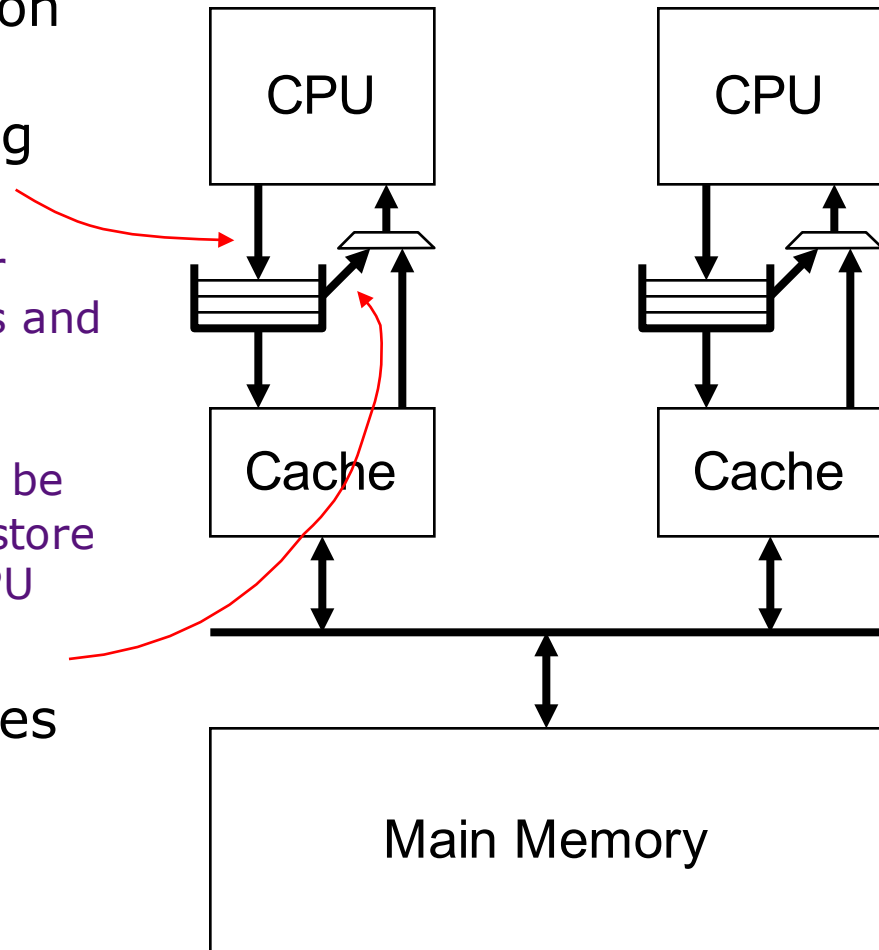
Architectural optimizations that are correct for uniprocessors, often violate sequential consistency and result in a new memory model for multiprocessors

Consistency Models

- Sequential Consistency
 - All reads and write in order
- Relaxed Consistency (one or more of the following)
 - Loads may be reordered after loads
 - e.g., PA-RISC, Power, Alpha
 - Loads may be reordered after stores
 - e.g., PA-RISC, Power, Alpha
 - Stores may be reordered after stores
 - e.g., PA-RISC, Power, Alpha, PSO
 - Stores may be reordered after loads
 - e.g., PA-RISC, Power, Alpha, PSO, TSO
 - Other more esoteric characteristics
 - e.g., Alpha

Committed Store Buffers

- CPU can continue execution while earlier committed stores are still propagating through memory system
 - Processor can commit other instructions (including loads and stores) while first store is committing to memory
 - Committed store buffer can be combined with speculative store buffer in an out-of-order CPU
- Local loads can bypass values from buffered stores to same address



Example 1: Store Buffers

Process 1

Store (flag₁), 1;
Load r₁, (flag₂);

Process 2

Store (flag₂), 1;
Load r₂, (flag₁);

Question: Is it possible that $r_1=0$ and $r_2=0$?

- *Sequential consistency: **No***
- *Suppose Loads can go ahead of Stores waiting in the store buffer: **Yes !***

Total Store Order (TSO):

Sun SPARC, IBM 370

Initially, all memory locations contain zeros

Example 2: Store-Load Bypassing

Process 1

Store (flag₁), 1;
 Load r₃, (flag₁);
 Load r₁, (flag₂);

Process 2

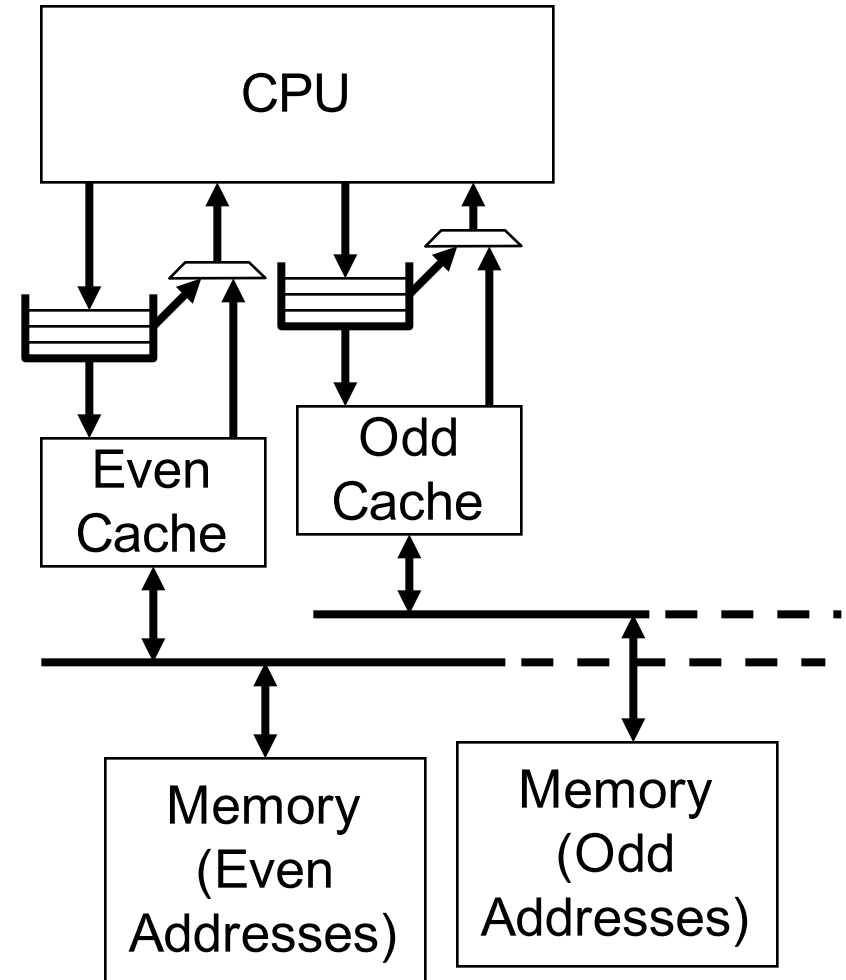
Store (flag₂), 1;
 Load r₄, (flag₂);
 Load r₂, (flag₁);

Question: Do extra Loads have any effect?

- *Sequential consistency: No*
- *Suppose Store-Load bypassing is permitted in the store buffer*
 - No effect in Sparc's TSO model, still not SC
 - In IBM 370, a load cannot return a written value until it is visible to other processors => implicity adds a memory fence, looks like SC

Interleaved Memory System

- Achieve greater throughput by spreading memory addresses across two or more parallel memory subsystems
 - In snooping system, can have two or more snoops in progress at same time (e.g., Sun UE10K system has four interleaved snooping busses)
 - Greater bandwidth from main memory system as two memory modules can be accessed in parallel



Example 3: Non-FIFO Store buffers

Process 1

Store (a), 1;

Store (flag), 1;

Process 2

Load r_1 , (flag);

Load r_2 , (a);

Question: Is it possible that $r_1=1$ but $r_2=0$?

- *Sequential consistency: No*
- *With non-FIFO store buffers: Yes*

Sparc's PSO memory model

Example 4: Non-Blocking Caches

Process 1

Store (a), 1;

Store (flag), 1;

Process 2

Load r_1 , (flag);

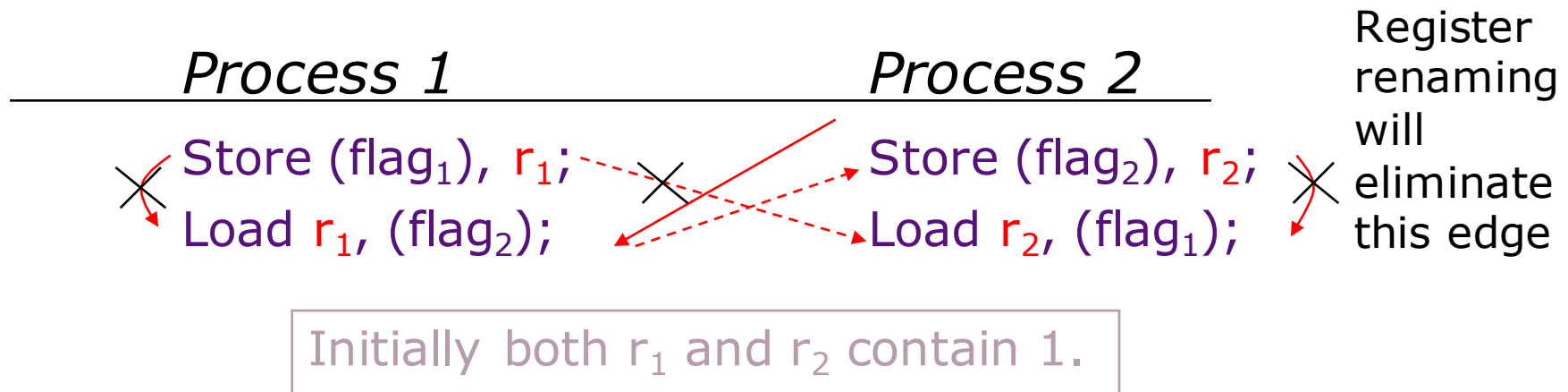
Load r_2 , (a);

Question: Is it possible that $r_1=1$ but $r_2=0$?

- *Sequential consistency: No*
- *Assuming stores are ordered: Yes because Loads can be reordered*

Alpha, Sparc's RMO, PowerPC's WO

Example 5: Register Renaming



Question: Is it possible that $r_1=0$ but $r_2=0$?

- *Sequential consistency: No*
- *Register renaming: Yes because it removes anti-dependencies*

Example 6: Speculative Execution

Process 1

Store (a), 1;
Store (flag), 1;

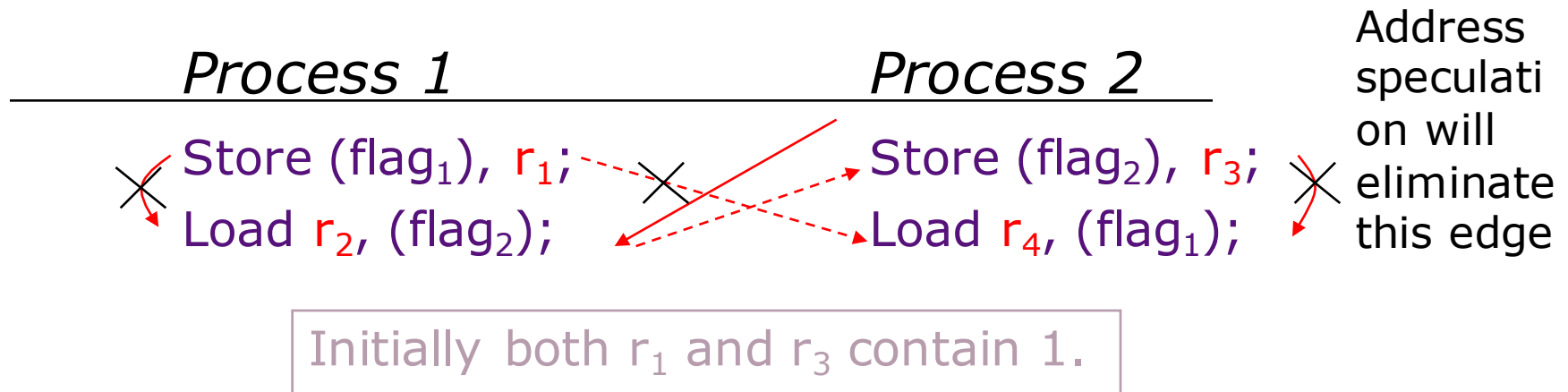
Process 2

L: Load r_1 , (flag);
if $r_1 == 0$ goto L;
Load r_2 , (a);

Question: Is it possible that $r_1=1$ but $r_2=0$?

- *Sequential consistency: No*
- *With speculative loads: Yes even if the stores are ordered*

Example 7: Address Speculation



Question: Is it possible that $r_2=0$ but $r_4=0$?

- *Sequential consistency: No*
- *Address speculation: Yes because it removes the dependencies between the stores and loads*

Flag₁ and flag₂ are registers pointing at memory locations

Example 8: Store Atomicity

<u>Process 1</u>	<u>Process 2</u>	<u>Process 3</u>	<u>Process 4</u>
Store (a),1;	Store (a),2;	Load $r_1, (a)$;	Load $r_3, (a)$;
		Load $r_2, (a)$;	Load $r_4, (a)$;

Question: Is it possible that $r_1=1$ and $r_2=2$ but $r_3=2$ and $r_4=1$?

- *Sequential consistency: No*
- *Even if Loads on a processor are ordered, the different ordering of stores can be observed if the Store operation is not atomic.*

Example 9: Causality

<u>Process 1</u>	<u>Process 2</u>	<u>Process 3</u>
Store (flag ₁), 1;	Load r ₁ , (flag ₁);	Load r ₂ , (flag ₂);
	Store (flag ₂), 1;	Load r ₃ , (flag ₁);

Question: Is it possible that $r_1=1$ and $r_2=1$ but $r_3=0$?

- *Sequential consistency: No*
- *With load/load reordering: Yes*

Alpha

Weaker Memory Models & Memory Fence Instructions

- Architectures with weaker memory models provide memory fence instructions to prevent otherwise permitted reorderings of loads and stores

Store (a_1), r2;

Fence_{wr}

Load r1, (a_2);

The Load and Store can be reordered if $a_1 \neq a_2$.
Insertion of **Fence_{wr}** will disallow this reordering

Similarly: Fence_{rr}; Fence_{rw}; Fence_{ww};

SUN's Sparc: MEMBAR;

MEMBARRR; MEMBARRW; MEMBARWR; MEMBARWW

PowerPC: Sync; EIEIO

Enforcing Ordering using Fences

Processor 1

Store (a),10;
Store (flag),1;

Processor 2

L: Load r₁, (flag);
if r₁ == 0 goto L;
Load r₂, (a);

Processor 1

Store (a),10;
Fence_{ww};
Store (flag),1;

Processor 2

L: Load r₁, (flag);
if r₁ == 0 goto L;
Fence_{rr};
Load r₂, (a);

Weak ordering

Weaker (Relaxed) Memory Models




- Hard to understand and remember
- Unstable - *Modèle de l'année*
- Abandon weaker memory models in favor of implementing SC.

Implementing SC

1. The memory operations of each individual processor appear to all processors in the order the requests are made to the memory.
 - *Provided by cache coherence, which ensures that all processors observe the same order of loads and stores to an address*
2. Any execution is the same as if the operations of all the processors were executed in some sequential order
 - *Provided by enforcing a dependence between each memory operation and the following one.*

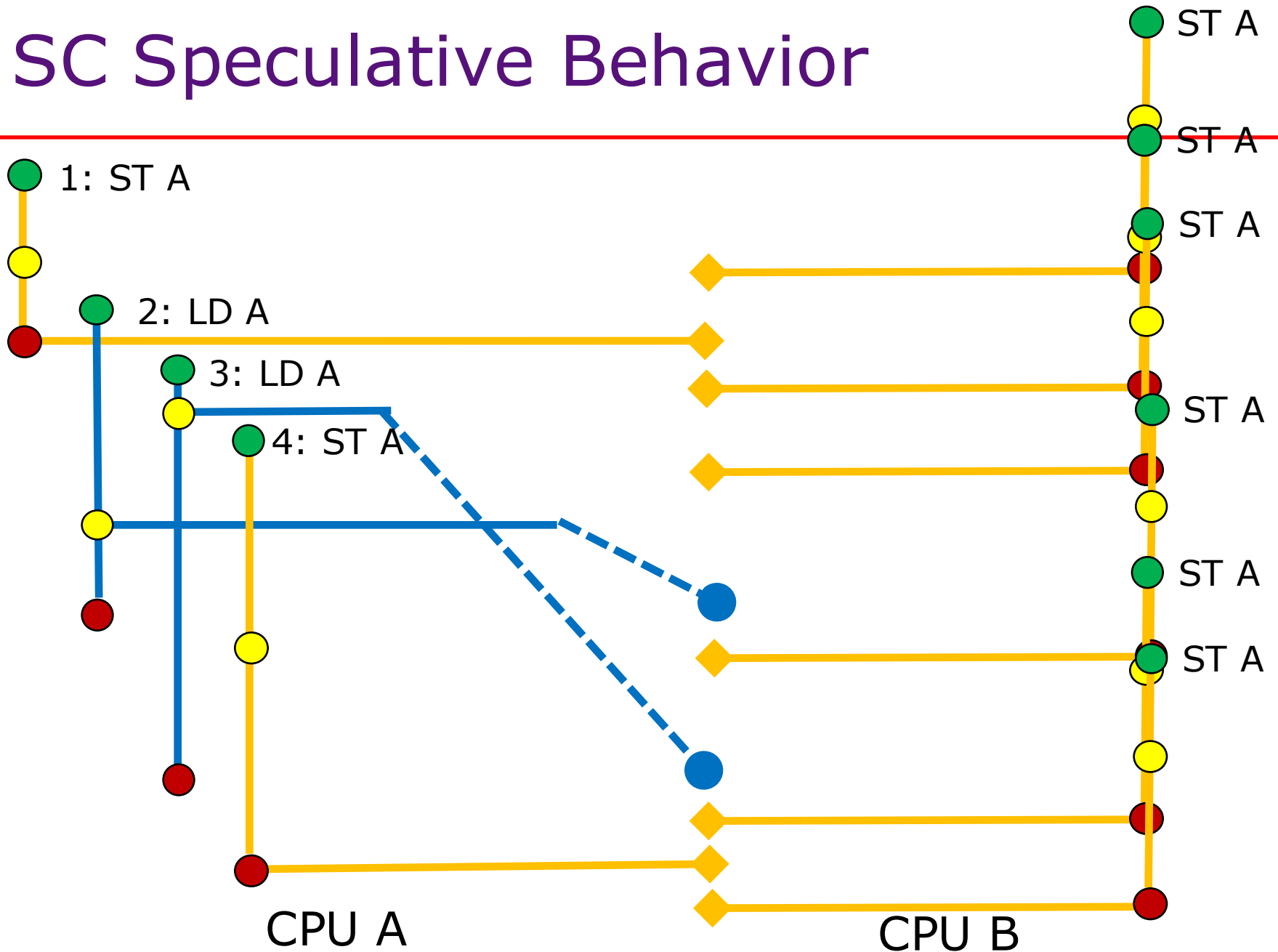
SC Data Dependence

- *Stall*
 - *Use in-order execution with blocking cache*
 - *Cache coherence plus allowing a processor to have only one request in flight at a time will provide SC*
- *Change architecture*  *Relaxed memory models*
 - *Use OOO and non-blocking caches*
 - *Cache coherence and allowing multiple requests (different addresses) concurrently gives high performance, then add fence operations to force ordering when needed*
- *Speculate...*

Sequential Consistency Speculation

- Local load-store ordering uses standard OOO mechanism
- Globally non-speculative stores
 - Stores execute at commit -> stores are in-order!
- Globally speculative loads
 - **Guess** at issue that the memory location used by a load will not change between issue and commit of the instruction
 - this is equivalent to loads happening in-order at commit
 - **Check** at commit by remembering all loads addresses starting at issue and watching for writes to that location.
 - **Data Management** for rollback relies on the basic out-of-order speculative data management used for uni-processor rollback and instruction re-execution.

SC Speculative Behavior



Properly Synchronized Programs

- Very few programmers do programming that relies on SC; instead higher-level synchronization primitives are used
 - locks, semaphores, monitors, atomic transactions
- A “properly synchronized program” is one where each shared writable variable is protected (say, by a lock) so that there is no race in updating the variable.
 - There is still race to get the lock
 - There is no way to check if a program is properly synchronized
- For properly synchronized programs, instruction reordering does not matter as long as updated values are committed before leaving a locked region.

Release Consistency

- Only care about inter-processor memory ordering at thread synchronization points, not in between
- Can treat all synchronization instructions as the only ordering points

...

Acquire(lock) // All following loads get most recent written values

... Read and write shared data ..

Release(lock) // All preceding writes are globally visible before
// lock is freed.

...

Takeaway

- SC is too low level a programming model. High-level programming should be based on critical sections & locks, atomic transactions, monitors, ...
- High-level parallel programming should be oblivious of memory model issues.
 - Programmer should not be affected by changes in the memory model
- ISA definition for Load, Store, Memory Fence, synchronization instructions should
 - Be precise
 - Permit maximum flexibility in hardware implementation
 - Permit efficient implementation of high-level parallel constructs.