



TECHNISCHE  
UNIVERSITÄT  
DRESDEN

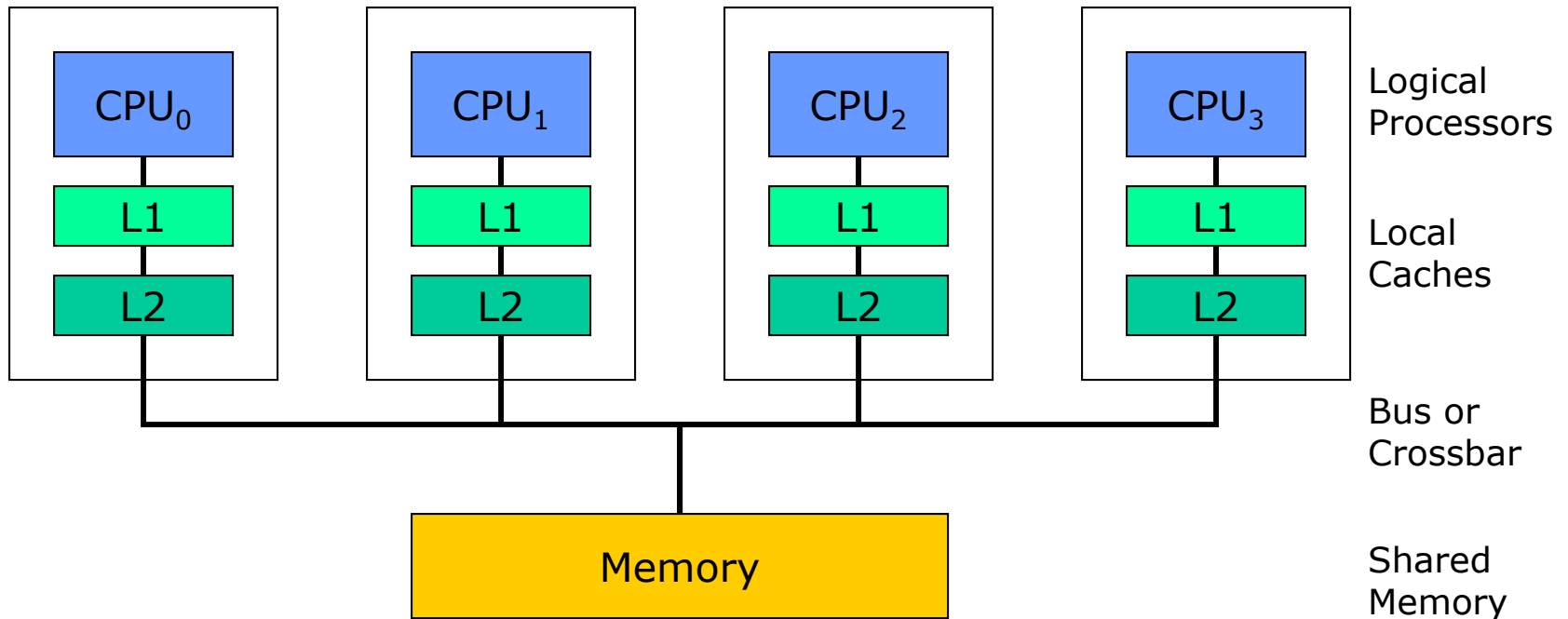
Department of Computer Science - Institute of System Architecture, Operating Systems Group

# Parallel Architectures

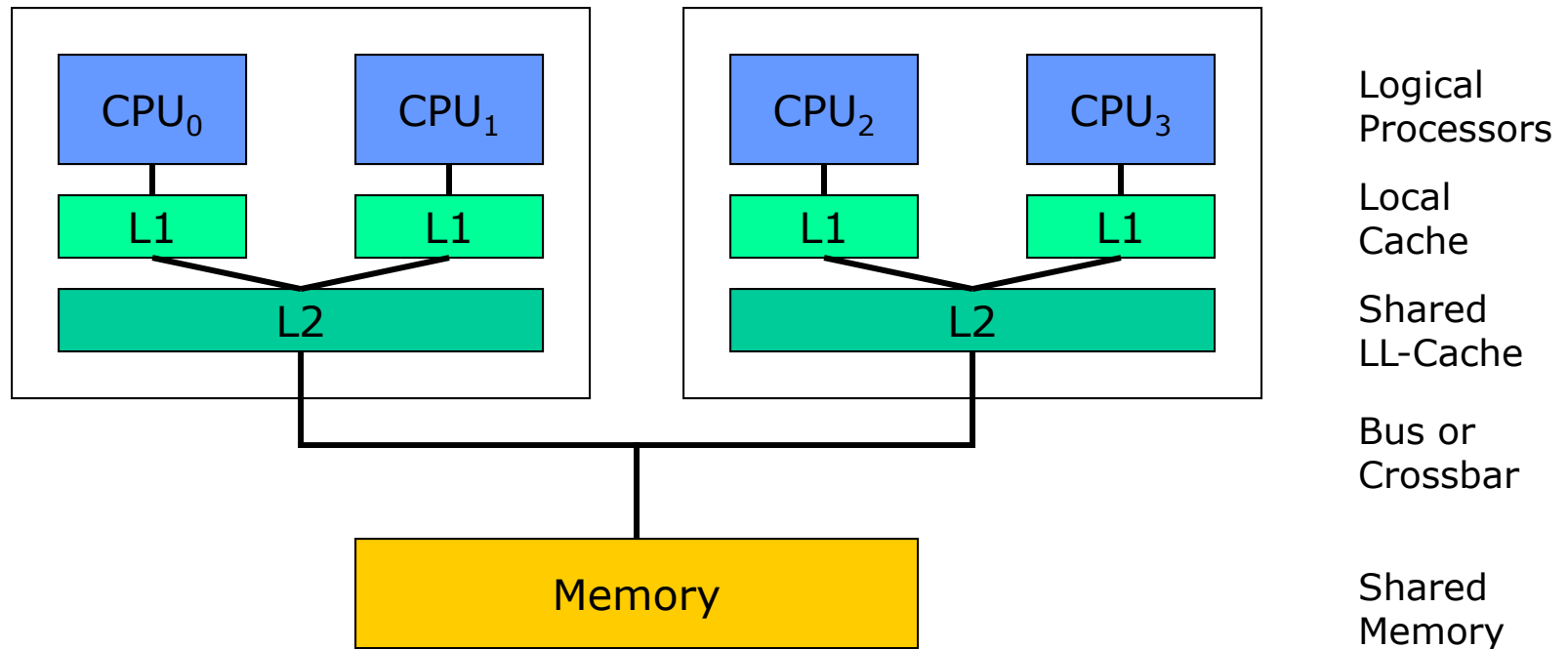
## Memory Consistency & Cache Coherency

Udo Steinberg

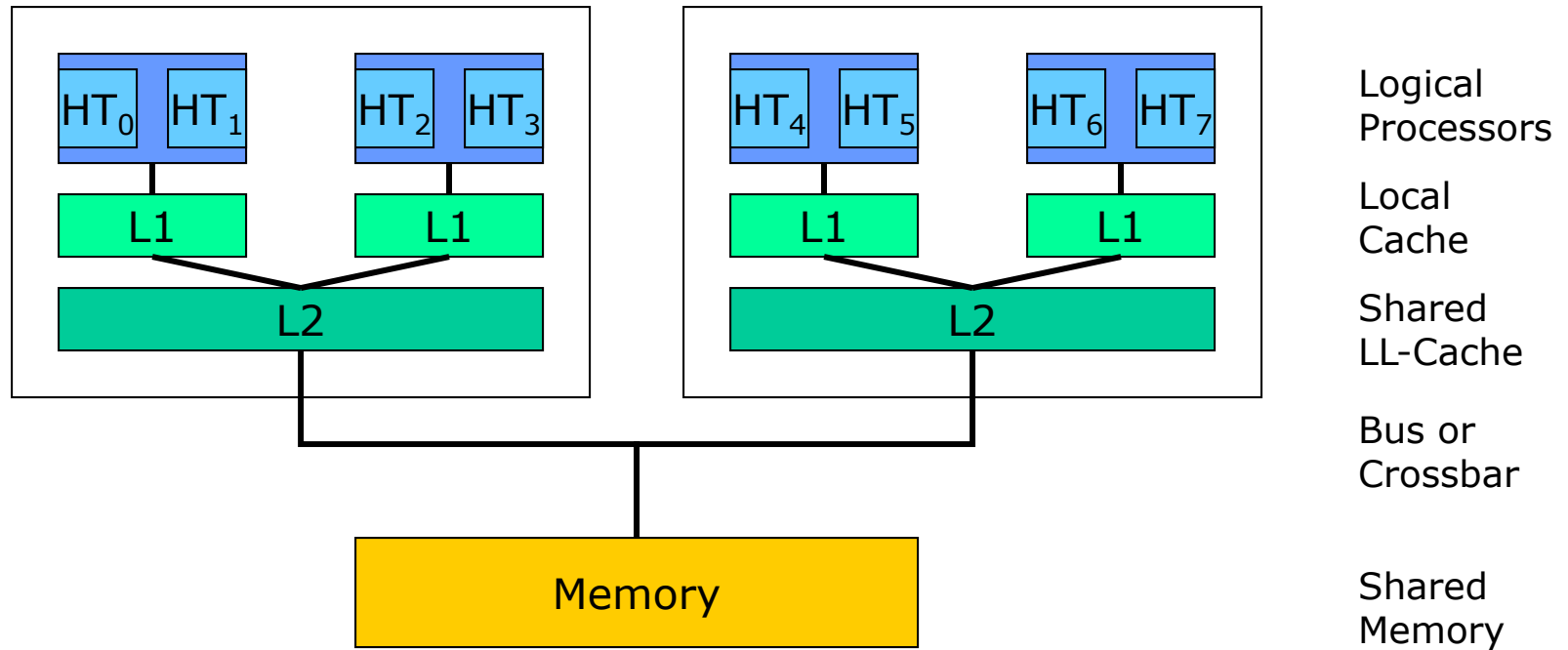
## Symmetric Multi-Processor (SMP)



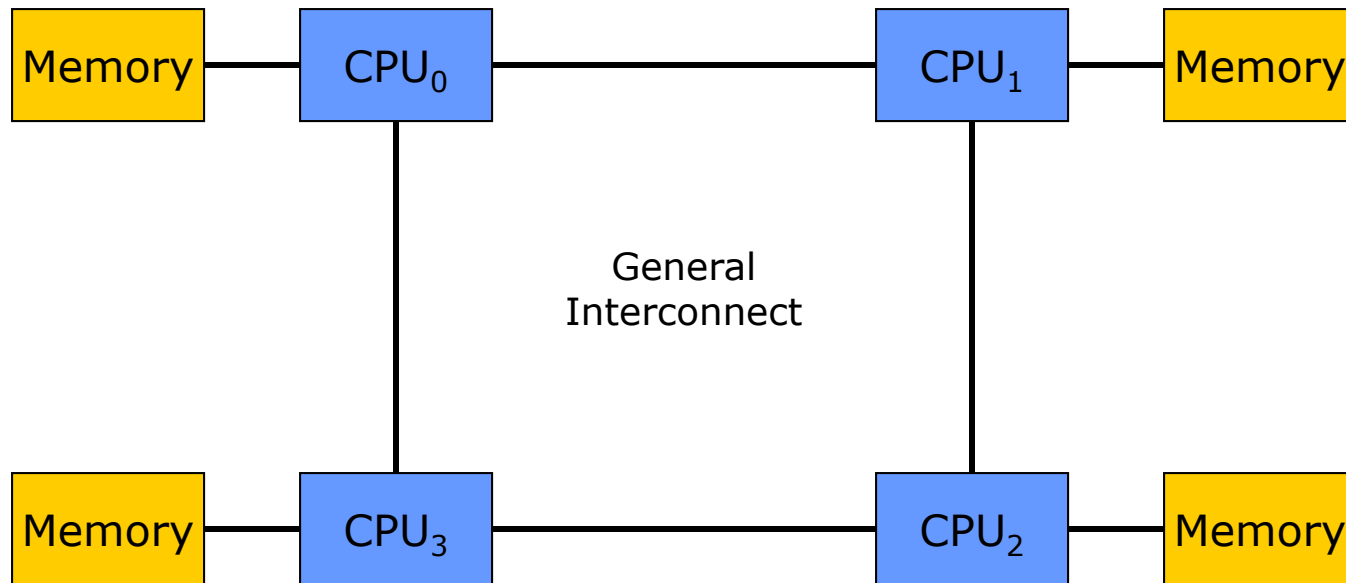
## Chip Multi-Processor (CMP), Multicore



## Symmetric Multi-Threading (SMT), Hyperthreading



## Non-Uniform Memory Access (NUMA)



## Multi-Processor Systems and Shared Memory

- Multiple processors share memory
- Memory managed by one or more memory controllers
  - UMA (Uniform Memory Access)
  - NUMA (Non-Uniform Memory Access)
- What is memory behavior under concurrent data access?
  - Reading a memory location should return last value written
  - „Last value written“ not clearly defined under concurrent access
- Defined by system`s memory consistency model
  - Defines in which order processors perceive concurrent accesses
  - Based on ordering, not timing of accesses

## Memory Consistency Models

- Different memory consistency models exist
  - Some platforms (e.g., SPARC) support multiple models
- More complex models attempt to expose more performance
- Terminology:
  - Program Order (of a processor's operations)
    - per-processor order of memory accesses determined by program (software)
  - Visibility Order (of all operations)
    - order of memory accesses observed by one or more processors
    - every read from a location returns value of most recent write

## Most Intuitive Model: Sequential Consistency

- 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. (Lamport 1979)
- Program Order Requirement
  - each CPU issues memory operations in program order
- Atomicity Requirement
  - Memory services operations one-at-a-time
  - all memory operations appear to execute atomically with respect to other memory operations



## Examples for Sequential Consistency

### **CPU<sub>1</sub>**

A = 1; (a<sub>1</sub>)

B = 1; (b<sub>1</sub>)

### **CPU<sub>2</sub>**

u = B; (a<sub>2</sub>)

v = A; (b<sub>2</sub>)

A,B ... Memory

u,v ... Registers

(u,v) = (1,1)

sequentially consistent

» example visibility order: a<sub>1</sub>,b<sub>1</sub>,a<sub>2</sub>,b<sub>2</sub>

(u,v) = (1,0)

sequentially inconsistent

» example visibility order: b<sub>1</sub>,a<sub>2</sub>,b<sub>2</sub>,a<sub>1</sub>

This visibility order violates program order on CPU<sub>1</sub>

No visibility order exists that satisfies program order on all CPUs and produces (u,v) = (1,0) result

## Examples for Sequential Consistency

### **CPU<sub>1</sub>**

A = 1; (a<sub>1</sub>)

u = B; (b<sub>1</sub>)

### **CPU<sub>2</sub>**

B = 1; (a<sub>2</sub>)

v = A; (b<sub>2</sub>)

(u,v) = (1,1)

sequentially consistent

» example visibility order: a<sub>1</sub>,a<sub>2</sub>,b<sub>1</sub>,b<sub>2</sub>

(u,v) = (0,0)

sequentially inconsistent

» example visibility order: b<sub>1</sub>,b<sub>2</sub>,a<sub>1</sub>,a<sub>2</sub>

This visibility order violates program order on CPU<sub>1/2</sub>

No visibility order exists that satisfies program order on all CPUs and produces (u,v) = (0,0) result

## Sequential Consistency vs. Architecture Optimizations

### **CPU<sub>1</sub>**

A = 1; (a<sub>1</sub>)

B = 1; (b<sub>1</sub>)

### **CPU<sub>2</sub>**

u = B; (a<sub>2</sub>)

v = A; (b<sub>2</sub>)

- Relaxing the Program Order:
  - Out-of-order execution may reorder operations (b<sub>2</sub>,a<sub>2</sub>)
  - Write Buffer may reorder writes (b<sub>1</sub>,a<sub>1</sub>)
  - produces sequentially inconsistent result (u,v) = (1,0)
- Maintaining Program Order:
  - May still produce sequentially inconsistent results
    - CPU<sub>1</sub> issues a<sub>1</sub>,b<sub>1</sub> in program order
    - but a<sub>1</sub> misses and b<sub>1</sub> hits in the cache (non-blocking cache)

## Causality

**CPU<sub>1</sub>**  
A = 1;

**CPU<sub>2</sub>**  
while (A == 0);  
B = 1;

**CPU<sub>3</sub>**  
while (B == 0);  
print A

- Relaxing the Atomicity of Writes:
  1. CPU<sub>1</sub> writes A = 1, generates update message to CPU<sub>2</sub> and CPU<sub>3</sub>
  2. CPU<sub>2</sub> receives update message for A from CPU<sub>1</sub>
  3. CPU<sub>2</sub> writes B = 1, generates update message to CPU<sub>3</sub>
  4. CPU<sub>3</sub> receives update message for B from CPU<sub>2</sub>
  5. CPU<sub>3</sub> prints A = 0
  6. CPU<sub>3</sub> receives update message for A from CPU<sub>1</sub>
    - Sequentially inconsistent result, because write to A not atomic wrt. other memory operations (e.g., write to B)

## Compiler Optimizations

**CPU<sub>1</sub>**

A = 1;

Flag = 1;

**CPU<sub>2</sub>**

while (Flag == 0);

u = A;

**CPU<sub>1</sub>**

A = 1;

Flag = 1;

**CPU<sub>2</sub>**

reg = Flag;

while (reg == 0);

u = A;

Programmer's Code

Compiler-Generated Code

- Compiler optimizations such as
  - register allocation and value caching
  - code motion
  - common sub-expression elimination
  - loop interchange

can reorder memory operations similar to architecture optimizations or even eliminate memory operations completely

## Relaxing Write-to-Read or Write-to-Write Order

- Write-to-Read (later reads can bypass earlier writes):
  - Write followed by a read can execute out-of-order
  - Typical hardware usage: Write Buffer
    - Writes must wait for ownership of cache line
    - Reads can bypass writes in the write buffer
    - Hides write latency
- Write-to-Write (later writes can bypass earlier writes) :
  - Write followed by another write can execute out-of-order
  - Typical hardware usage: Non-Blocking Cache, Write Coalescing
    - Writes must wait for ownership of cache line
    - Latency for obtaining ownership depends on hop count to cache line owner

## IBM-370 (zSeries)

- In-order memory operations:
  - Read-to-Read
  - Read-to-Write
  - Write-to-Write
- Out-of-order memory operations:
  - Write-to-Read (later reads can bypass earlier writes)
    - unless both are to the same memory location
    - breaks Dekker's algorithm for mutual exclusion
  - Write-to-Read to same location must execute in-order
    - no forwarding of pending writes from the write buffer

## Dekker's Algorithm on IBM-370 (zSeries)

```
bool flag0 = false, flag1 = false; // intention to enter crit. section
int turn = 0;                       // whose turn is it?
```

### CPU #0

```
P: flag0 = true; // Buffered
  while (flag1) {
    if (turn == 1) {
      flag0 = false;
      goto P;
    }
  }
  // critical section
  flag0 = false;
  turn = 1;
```

### CPU #1

```
P: flag1 = true; // Buffered
  while (flag0) {
    if (turn == 0) {
      flag1 = false;
      goto P;
    }
  }
  // critical section
  flag1 = false;
  turn = 0;
```



## SPARC V8 Total Store Order (TSO)

- In-order memory operations:
  - Read-to-Read
  - Read-to-Write
  - Write-to-Write
- Out-of-order memory operations:
  - Write-to-Read (later reads can bypass earlier writes)
    - Forwarding of pending writes in the write buffer to successive read operations of the same location
      - Writes become visible to writing processor first
    - Breaks Peterson`s algorithm for mutual exclusion

## Peterson's Algorithm on SPARC V8 TSO

```
bool flag0 = false, flag1 = false; // intention to enter crit. section
int turn = 0;                       // whose turn is it?
```

### CPU #0

```
flag0 = true; // Buffered
turn = 1;     // Buffered
while (turn == 1 && flag1) {};
// critical section
flag0 = false;
```

### CPU #1

```
flag1 = true; // Buffered
turn = 0;     // Buffered
while (turn == 0 && flag0) {};
// critical section
flag1 = false;
```

## Total Store Order (TSO) vs. SC and IBM-370

### **CPU<sub>1</sub>**

A = 1; (a<sub>1</sub>)

u = A; (b<sub>1</sub>)

w = B; (c<sub>1</sub>)

### **CPU<sub>2</sub>**

B = 1; (a<sub>2</sub>)

v = B; (b<sub>2</sub>)

x = A; (c<sub>2</sub>)

- (u,v,w,x) = (1,1,0,0)
  - not possible with Sequential Consistency (SC) and IBM-370
  - but possible with Total Store Order (TSO)
    - Example total order: b<sub>1</sub>, b<sub>2</sub>, c<sub>1</sub>, c<sub>2</sub>, a<sub>1</sub>, a<sub>2</sub>
    - b<sub>1</sub> reads A=1 from write buffer
    - b<sub>2</sub> reads B=1 from write buffer

## Processor Consistency (PC)

- Similar to Total Store Order (TSO)
- But model additionally supports multiple cached memory copies
  - Relaxed atomicity for write operations
    - Each write operation broken into sub-operations to update cached copies of other CPUs
  - Non-unique write order, requires per-CPU visibility order
  - Additional Coherency Requirement:
    - All writes sub-operations to the same memory location complete in the same order across all memory copies (or in other words: every processor should see writes to the same location in the same order)
    - If one CPU observes writes to X in the order  $W_1(X)$  before  $W_2(X)$ , another CPU must not see  $W_2(X)$  before  $W_1(X)$

## Processor Consistency (PC) vs. SC, IBM-370 and TSO

**CPU<sub>1</sub>**

A = 1; (a<sub>1</sub>)

**CPU<sub>2</sub>**

u = A; (a<sub>2</sub>)

B = 1; (b<sub>2</sub>)

**CPU<sub>3</sub>**

v = B; (a<sub>3</sub>)

w = A; (b<sub>3</sub>)

- (u,v,w) = (1,1,0)
  - not possible with SC, IBM-370 and TSO
  - but possible with Processor Consistency (PC)
    - CPU<sub>1</sub> sets A = 1, sends W<sub>1</sub>(A) to other CPUs
    - CPU<sub>2</sub> observes W<sub>1</sub>(A), sets B = 1, sends W<sub>2</sub>(B) to other CPUs
    - CPU<sub>3</sub> observes W<sub>2</sub>(B) ... but has not yet received W<sub>1</sub>(A)
  - Single memory bus enforces single visibility order
  - Multiple visibility orders possible with other topologies

## SPARC V8 Partial Store Order (PSO)

- In-order memory operations:
  - Read-to-Read
  - Read-to-Write
- Out-of-order memory operations:
  - Write-to-Read (later reads can bypass earlier writes)
    - Forwarding of pending writes in the write buffer to successive read operations of the same location
  - Write-to-Write (later writes can bypass earlier writes)
    - unless both are to the same memory location
    - breaks Producer-Consumer Code
- Write Atomicity is maintained -> single visibility order

## Partial Store Order (PSO) vs. SC, IBM-370, TSO and PC

### **CPU<sub>1</sub>**

A = 1;                   (a<sub>1</sub>)

B = 1;                   (b<sub>1</sub>)

Flag = 1;               (c<sub>1</sub>)

### **CPU<sub>2</sub>**

while (Flag == 0);       (a<sub>2</sub>)

u = A;                   (b<sub>2</sub>)

v = B;                   (c<sub>2</sub>)

- (u,v) = (0,0) or (0,1) or (1,0)
  - not possible with SC, IBM-370, TSO and PC
  - but possible with Partial Store Order (PSO)
    - Example total order: c<sub>1</sub>,a<sub>2</sub>,b<sub>2</sub>,c<sub>2</sub>,b<sub>1</sub>,a<sub>1</sub>
    - Store barrier (STBAR) before c<sub>1</sub> ensures sequentially consistent result (u,v) = (1,1)

## Relaxing all Program Orders

- In addition to previous relaxations:
  - Read-to-Read (later reads can bypass earlier reads) :
    - Read followed by a read can execute out-of-order
  - Read-to-Write (later writes can bypass earlier reads):
    - Read followed by a write can execute out-of-order
- Examples:
  - Weak Ordering (WO)
  - Release Consistency (RC)
  - DEC Alpha
  - SPARC V9 Relaxed Memory Order (RMO)
  - PowerPC
  - Itanium (IA64)



## Weak Ordering (WO)

- Conceptually similar to Processor Consistency (PC)
  - including coherency requirement
- Classifies memory operations into two categories:
  - data operations
  - synchronization operations
- Reordering of memory accesses between synchronization operations typically does not affect correctness of a program
- Program order only maintained at synchronization points
  - between synchronization operations

## Release Consistency (RC)

- Distinguishes memory operations as
  - ordinary (data)
  - special
    - sync (synchronization)
    - nsync (asynchronous data)
- Sync operations classified as
  - acquire
    - read operation for gaining access to a shared resource
    - e.g., spinning on a flag to be set
  - release
    - write operation for granting permission to a shared resource
    - e.g., setting a synchronization flag

## Flavors of Release Consistency (RC)

- $RC_{SC}$ 
  - Sequential consistency between special operations
  - Program order enforced between:
    - acquire -> all
    - all -> release
    - special -> special
- $RC_{PC}$ 
  - Processor consistency between special operations
  - Program order enforced between:
    - acquire -> all
    - all -> release
    - special -> special
      - except special write followed by special read
      - can use read-modify-write instruction to achieve effect

## Memory Consistency in Modern Architectures

Reordered Memory Accesses	Read-to-Read	Read-to-Write	Write-to-Write	Write-to-Read	Atomic OPs and Reads	Atomic OPs and Writes
Alpha	Y	Y	Y	Y	Y	Y
AMD64	Y			Y		
IA64	Y	Y	Y	Y	Y	Y
PA-RISC	(Y)	(Y)	(Y)	(Y)		
POWER*	Y	Y	Y	Y	Y	Y
SPARC RMO	Y	Y	Y	Y	Y	Y
SPARC PSO			Y	Y		Y
SPARC TSO				Y		
IA32*	Y	Y	(Y)	Y		

\*write atomicity relaxed

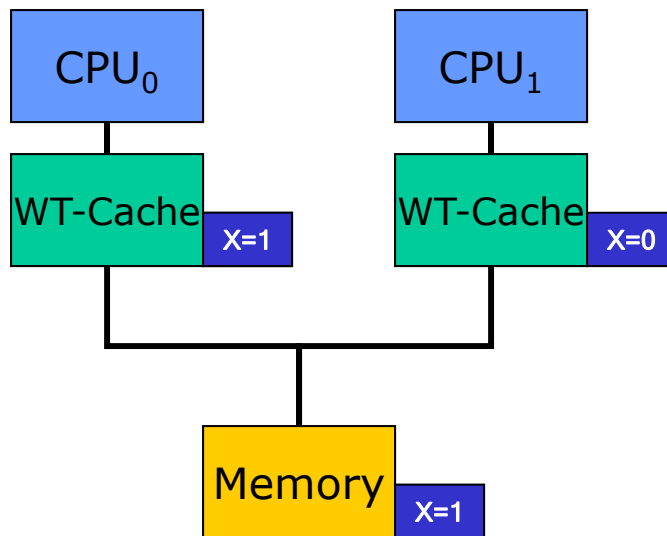
## Enforcing Ordering: Synchronization Instructions

- IA32/AMD64:
  - lfence (load fence), sfence (store fence), mfence (memory fence)
- Alpha:
  - mb (memory barrier), wmb (write memory barrier)
- SPARC (PSO)
  - stbar (store barrier)
- SPARC (RMO)
  - membar (4-bit encoding for r-r, r-w, w-r, w-w ordering)
- PowerPC
  - sync (similar to Alpha mb, except for r-r), lwsync

## Cache Coherency

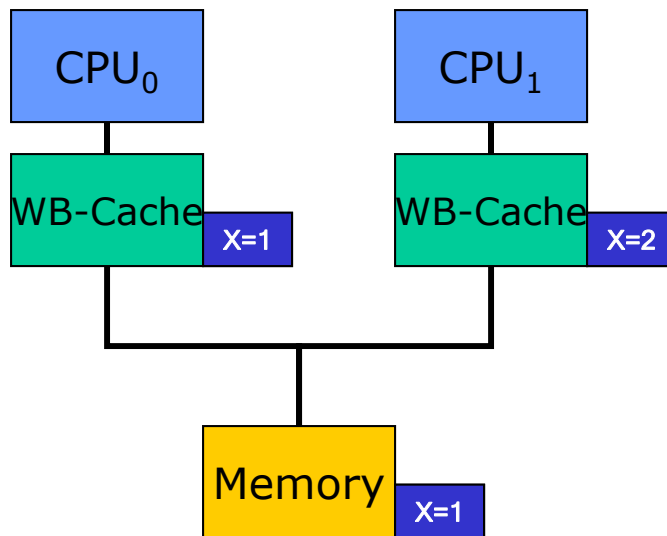
- Caching leads to presence of multiple copies for a memory location
- Cache coherency is a mechanism for keeping copies up-to-date
  - locate all cached copies of a memory location
  - eliminate stale copies (invalidate/update)
- Requirements:
  - Write Propagation: Writes must eventually become visible to all processors
  - Write Serialization: Every processor should see the writes to the **same** location in the same order

## Incoherency Example (1)



1. CPU<sub>0</sub> reads X from memory
    - stores X=0 into its cache
  2. CPU<sub>1</sub> reads X from memory
    - stores X=0 into its cache
  3. CPU<sub>0</sub> writes X=1
    - stores X=1 in its cache
    - stores X=1 in memory
  4. CPU<sub>1</sub> reads X from its cache
    - loads X=0 from its cache
- Incoherent value for X on CPU<sub>1</sub>

## Incoherency Example (2)



1. CPU<sub>0</sub> reads X from memory
    - loads X=0 into its cache
  2. CPU<sub>1</sub> reads X from memory
    - loads X=0 into its cache
  3. CPU<sub>0</sub> writes X=1
    - stores X=1 in its cache
  4. CPU<sub>1</sub> writes X=2
    - stores X=2 in its cache
  5. CPU<sub>1</sub> writes back cache line
    - stores X=2 in memory
  6. CPU<sub>0</sub> writes back cache line
    - stores X=1 in memory
- Later store X=2 from CPU<sub>1</sub> lost



## Cache Coherency: Problems and Solutions

- Problem 1:  
CPU<sub>1</sub> used stale value that had already been modified by CPU<sub>0</sub>
  - Solution:  
Invalidate all copies before allowing a write to proceed
- Problem 2:  
Incorrect writeback order of modified cache lines
  - Solution:  
Disallow more than one modified copy

## Coherency Protocol Approaches

- Invalidation-based
  - all coherency-related traffic broadcast to all CPUs
  - each processor snoops traffic and reacts accordingly
    - invalidate lines written to by another CPU
    - signal sharing for cache lines currently in cache
  - straightforward solution for bus-based systems
  - suited for small-scale systems
- Update-based
  - Uses central directory for cache line ownership
  - Write operation updates copies in other caches
    - can update all other CPUs at once (less bus traffic)
    - but: multiple writes cause multiple updates (more bus traffic)
  - suited for large-scale systems

## Invalidation vs. Update Protocols

- Invalidation-based
  - only write misses hit the bus (suited for write-back caches)
  - subsequent writes to same cache-line are write-hits
  - Good for multiple writes to the same cache line by the same CPU
- Update-based
  - all sharers of the cache line continue to hit in the cache after a write by one cache
  - Good for large-scale producer-consumer code
  - Otherwise lots of useless updates (wastes bandwidth)
- Hybrid forms are possible

## MESI Cache Coherency Protocol

- Modified (M)
  - No copies exist in other caches; local copy is modified
  - Memory is stale
- Exclusive (E)
  - No copies exist in other caches
  - Memory is up-to-date
- Shared (S)
  - Unmodified copies may exist in other caches
  - Memory is up-to-date
- Invalid (I)
  - Not in Cache

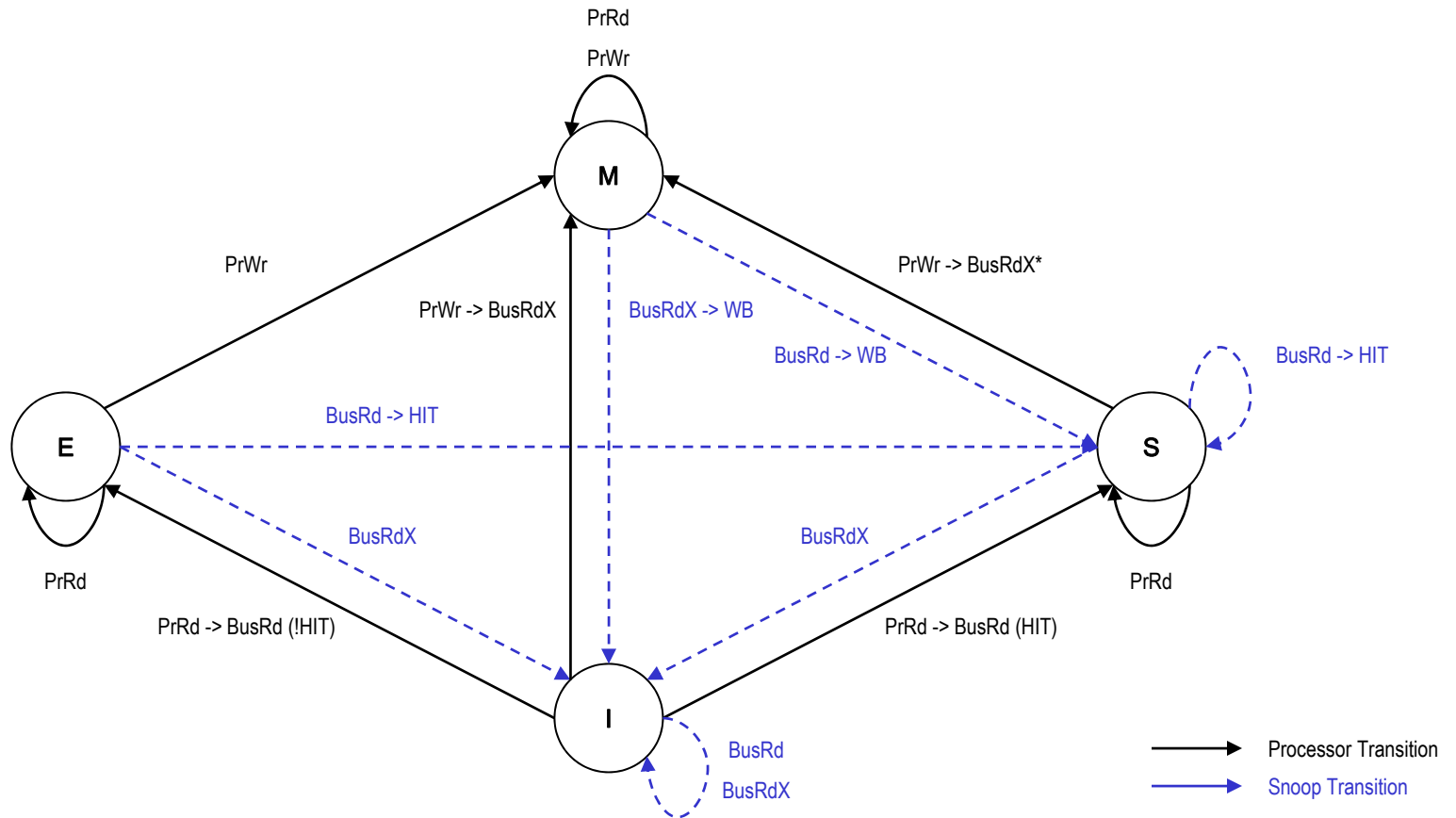
## MESI Cache Coherency Protocol (Processor Transitions)

- State is I, CPU reads (PrRd)
  - Generate bus read request (BusRd), other caches signal sharing
  - If cache line in another cache go to S, otherwise transition to E
- State is S, E or M, CPU reads (PrRd)
  - No bus transaction, cache line already cached
- State is I, CPU writes (PrWr)
  - Generate bus read request for exclusive ownership (BusRdX)
  - transition to M
- State is S, CPU writes (PrWr)
  - Cache line already cached, but need to upgrade it for exclusive ownership (BusRdX\*), transition to M
- State is E or M, CPU writes (PrWr)
  - No bus transaction, cache line already exclusively cached
  - transition to M

## MESI Cache Coherency Protocol (Snoop Transitions)

- Receiving a read snoop (BusRd) for a cache line
  - If cache line is in cache (E or S), tell the requesting cache that the line is going to be shared (HIT signal) and transition to S
  - If cache line is modified in cache (M), write the cache line back to memory (WB) and transition to S
- Receiving a read for exclusive ownership snoop (BusRdX) for a cache line
  - If cache line is modified in cache (M), write the cache line back to memory (WB), discard it and transition to I
  - If cache line is unmodified (E or S), discard it and transition to I

# MESI Cache Coherency Protocol



## MOESI Cache Coherency Protocol

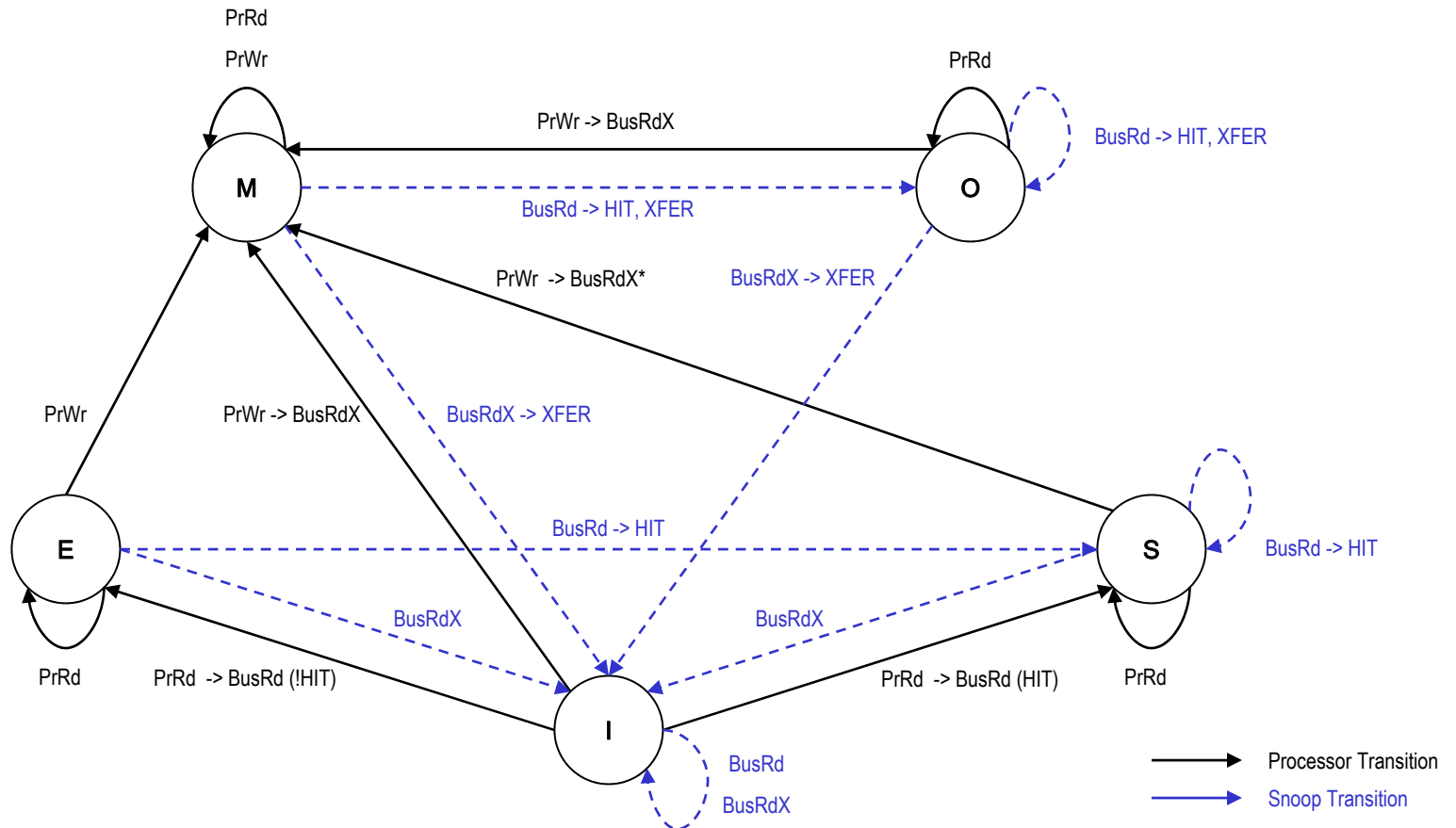
- Modified (M) Modified-Exclusive
  - No copies exist in other caches; local copy is modified
  - Memory is stale; Cache supplies copy instead of memory
- Owner (O) Modified-Shared
  - Unmodified copies may exist in other caches; local copy is modified
  - Memory is stale; Cache supplies copy instead of memory
- Exclusive (E)
  - No copies exist in other caches
  - Memory is up-to-date
- Shared (S)
  - Unmodified copies may exist in other caches
  - Memory is up-to-date unless a processor holds copy in O state
- Invalid (I)
  - Not in cache



## MOESI Cache Coherency Protocol (Transitions)

- Similar to MESI, with some extensions
- Cache-to-Cache transfers of modified cache lines
  - Cache in M or O state always transfers (XFER) cache line to requesting cache instead of memory supplying the cache line
- Avoids write-back to memory when another processor accesses the cache line
  - Beneficial when cache-to-cache latency/bandwidth is better than cache-to-memory latency/bandwidth
    - E.g., multi-core CPU with shared last-level cache

# MOESI Cache Coherency Protocol



## Coherency in Multi-Level Caches

- Bus only connected to last-level cache (e.g., L2)
- Problem:
  - Snoop requests are relevant to inner-level caches (e.g, L1)
  - Modifications made in L1 may not be visible to L2 (and the bus)
- L1 intervention:
  - on BusRd check if cache line is M in L1 (may be E or S in L2)
  - on BusRdX send invalidation to L1
- Some interventions not needed when L1 is write-through
  - but causes more write traffic to L2