

ECE/CS 757: Advanced Computer Architecture II

Instructor: Mikko H Lipasti

Spring 2009
University of Wisconsin-Madison

Lecture notes based on slides created by John Shen, Mark Hill, David Wood, Guri Sohi, and Jim Smith, Natalie Enright Jerger, and probably others

Cache Coherence

- Coherence States
- Snoopy bus-based Invalidate Protocols
- Invalidate protocol optimizations
- Update Protocols (Dragon/Firefly)
- Directory protocols
- Implementation issues

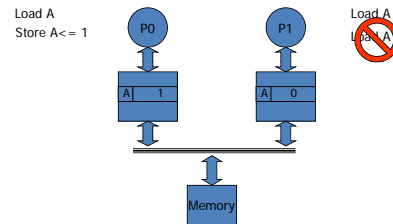
2

Readings

- Readings:
 - Firefly
 - Archibald
 - Sweazey/Smith
 - Laudon/Lenoski: Origin 2000
 - Opteron
 - Gigaplane
 - Power5
 - Intel 870

3

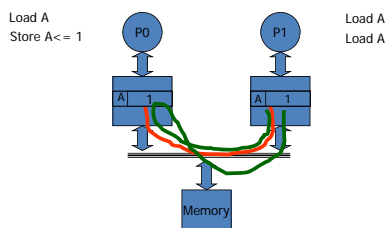
Cache Coherence Problem



© 2005 Mikko Lipasti

4

Cache Coherence Problem



© 2005 Mikko Lipasti

5

Possible Causes of Incoherence

- Sharing of writeable data
 - Cause most commonly considered
- Process migration
 - Can occur even if independent jobs are executing
- I/O
 - Often fixed via O/S cache flushes

02/07

ECE/CS 757; copyright J. E. Smith, 2007

6

Cache Coherence

- Informally, with coherent caches: accesses to a memory location *appear* to occur simultaneously in all copies of the memory location
 - "copies" \Rightarrow caches
- Cache coherence suggests an absolute time scale -- this is not necessary
 - What is required is the "appearance" of coherence... not absolute coherence
 - E.g. temporary incoherence between memory and a write-back cache may be OK.

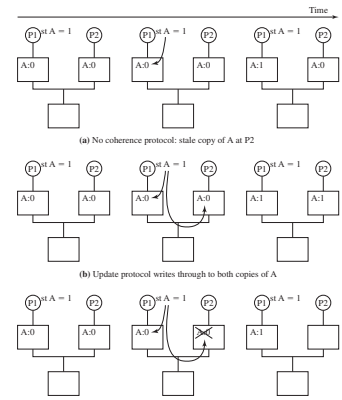
02/07

ECE/CS 757, copyright J. E. Smith, 2007

7

Update vs. Invalidation Protocols

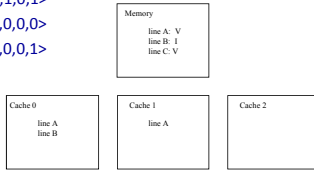
- Coherent Shared Memory
 - All processors see the effects of others' writes
- How/when writes are propagated
 - Determine by coherence protocol



© 2005 Mikko Lipasti. Invalidate protocol eliminates stale cache copy

Global Coherence States

- A memory line can be present (valid) in any of the caches and/or memory
- Represent global state with an N+1 element vector
 - First N components \Rightarrow cache states (valid/invalid)
 - N+1st component \Rightarrow memory state (valid/invalid)
- Example:
 - Line A: $\langle 1, 1, 0, 1 \rangle$
 - Line B: $\langle 1, 0, 0, 0 \rangle$
 - Line C: $\langle 0, 0, 0, 1 \rangle$



02/07

ECE/CS 757, copyright J. E. Smith, 2007

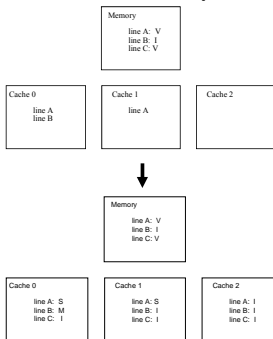
Local Coherence States

- Individual caches can maintain a summary of the state of memory lines, from a "local" perspective
 - Reduces storage for maintaining state
 - May have only partial information
- Invalid (I): $\langle 0, X, X, \dots, X \rangle$ -- local cache does not have a valid copy; (cache miss)
 - Don't confuse invalid state with empty frame
- Shared (S): $\langle 1, X, X, \dots, 1 \rangle$ -- local cache has a valid copy, main memory has a valid copy, other caches ??
- Modified (M): $\langle 1, 0, 0, \dots, 0 \rangle$ -- local cache has only valid copy.
- Exclusive (E): $\langle 1, 0, 0, \dots, 1 \rangle$ -- local cache has a valid copy, no other caches do, main memory has a valid copy.
- Owned (O): $\langle 1, X, X, \dots, X \rangle$ -- local cache has a valid copy, all other caches and memory may have a valid copy.
 - Only one cache can be in O state
 - $\langle 1, X, 1, X, \dots, 0 \rangle$ is included in O, but not included in any of the others.

02/07

ECE/CS 757, copyright J. E. Smith, 2007

Example



02/07

ECE/CS 757, copyright J. E. Smith, 2007

11

Snoopy Cache Coherence

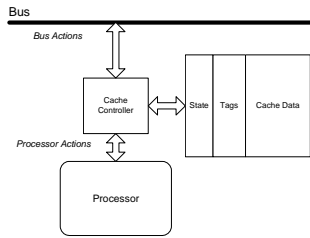
- All requests broadcast on bus
- All processors and memory snoop and respond
- Cache blocks writeable at one processor or read-only at several
 - Single-writer protocol
- Snoops that hit dirty lines?
 - Flush modified data out of cache
 - Either write back to memory, then satisfy remote miss from memory, or
 - Provide dirty data directly to requestor
 - Big problem in MP systems
 - Dirty/coherence/sharing misses

© 2005 Mikko Lipasti

12

Bus-Based Protocols

- Protocol consists of states and actions (state transitions)
- Actions can be invoked from processor or bus

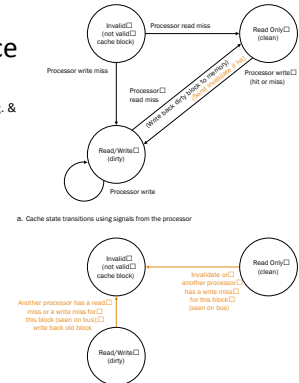


02/07

ECE/CS 757; copyright J. E. Smith, 2007

Minimal Coherence Protocol FSM

[Source: Patterson/Hennessy, Comp. Org. & Design]



© 2005 Mikko Lipasti

14

MSI Protocol

Current State	Action and Next State					
	Processor Read	Processor Write	Eviction	Cache Read	Cache Read&M	Cache Upgrade
I	Cache Read Acquire Copy → S	Cache Read&M Acquire Copy → M		No Action → I	No Action → I	No Action → I
S	No Action → S	Cache Upgrade → M	No Action → I	No Action → S	Invalidate Frame → I	Invalidate Frame → I
M	No Action → M	No Action → M	Cache Write back → I	Memory inhibit; Supply data; → S	Invalidate Frame; Memory inhibit; Supply data; → I	

02/07

ECE/CS 757; copyright J. E. Smith, 2007

MSI Example

Thread Event	Bus Action	Data From	Global State	Local States:		
				C0	C1	C2
0. Initially:			<0,0,1>	I	I	I
1. T0 read→	CR	Memory	<1,0,1>	S	I	I
2. T0 write→	CU		<1,0,0>	M	I	I
3. T2 read→	CR	C0	<1,0,1>	S	I	S
4. T1 write→	CRM	Memory	<0,1,0,0>	I	M	I

- If line is in no cache
 - Read, modify, Write requires 2 bus transactions
 - Optimization: add Exclusive state

02/07

ECE/CS 757; copyright J. E. Smith, 2007

Invalidate Protocol Optimizations

- Observation: data can be read shared
 - Add S (shared) state to protocol: MSI
- State transitions:
 - Local read: I->S, fetch shared
 - Local write: I->M, fetch modified; S->M, invalidate other copies
 - Remote read: M->S, supply data
 - Remote write: M->I, supply data; S->I, invalidate local copy
- Observation: data can be write-private (e.g. stack frame)
 - Avoid invalidate messages in that case
 - Add E (exclusive) state to protocol: MESI
- State transitions:
 - Local read: I->E if only copy, I->S if other copies exist
 - Local write: E->M silently, S->M, invalidate other copies

© 2005 Mikko Lipasti

17

MESI Protocol

- Variation used in Pentium Pro/2/3
 - (cache to cache transfer)
- 4-State Protocol
 - Modified: <1,0,0,...,0>
 - Exclusive: <1,0,0,...,1>
 - Shared: <1,X,X,...,1>
 - Invalid: <0,X,X,...,X>
- Bus/Processor Actions
 - Same as MSI
- Adds *shared* signal to indicate if other caches have a copy

02/07

ECE/CS 757; copyright J. E. Smith, 2007

18

MESI Protocol

Action and Next State						
Current State	Processor Read	Processor Write	Eviction	Cache Read	Cache Read&M	Cache Upgrade
I	Cache Read If no sharers: → E If sharers: → S	Cache Read&M → M		No Action → I	No Action → I	No Action → I
S	No Action → S	Cache Upgrade → M	No Action → I	Respond Shared: → S	No Action → I	No Action → I
E	No Action → E	No Action → M	No Action → I	Respond Shared: → S	No Action → I	
M	No Action → M	No Action → M	Cache Write-back → I	Respond dirty; Write back data; → S	Respond dirty; Write back data; → I	

02/07

ECE/CS 757, copyright J. E. Smith, 2007

MESI Example

- Optimization – cache-to-cache transfer for shared data (but only one sharer can respond)
- If modified in some cache and another reads then writeback to memory (w/ snarf)
 - Optimization: add **Owned** state (and perform cache-to-cache transfer)

Thread Event	Bus Action	Data From	Global State	Local States:		
				C0	C1	C2
0. Initially:			<0,0,0,1>	I	I	I
1. T0 read→	CR	Memory	<1,0,0,1>	E	I	I
2. T0 write→	none		<1,0,0,0>	M	I	I

02/07

ECE/CS 757, copyright J. E. Smith, 2007

MOESI Optimization

- Observation: shared ownership complicates/delays sourcing data
 - Owner is responsible for sourcing data to any requestor
 - Add O (owner) state to protocol: MOSI/MOESI
 - Last requestor becomes the owner
 - Ownership can be on per-node basis in hierarchically structured system
 - Avoid writeback (to memory) of dirty data
 - Also called *shared-dirty* state, since memory is stale

© 2005 Mikko Lipasti

21

MOESI Protocol

- Used in AMD Opteron
- 5-State Protocol
 - Modified: <1,0,0...0>
 - Exclusive: <1,0,0,...,1>
 - Shared: <1,X,X,...,X>
 - Invalid: <0,X,X,...,X>
 - Owned: <1,X,X,X,X> ; *only one owner*
- Owner can supply data, so memory does not have to
 - Avoids lengthy memory access

02/07

ECE/CS 757, copyright J. E. Smith, 2007

22

MOESI Protocol

Action and Next State						
Current State	Processor Read	Processor Write	Eviction	Cache Read	Cache Read&M	Cache Upgrade
I	Cache Read If no sharers: → E If sharers: → S	Cache Read&M → M		No Action → I	No Action → I	No Action → I
S	No Action → S	Cache Upgrade → M	No Action → I	Respond shared: → S	No Action → I	No Action → I
E	No Action → E	No Action → M	No Action → I	Respond shared; Supply data: → S	Respond shared; Supply data: → I	
O	No Action → O	Cache Upgrade → M	Cache Write-back → I	Respond shared; Supply data: → O	Respond shared; Supply data: → I	
M	No Action → M	No Action → M	Cache Write-back → I	Respond shared; Supply data: → O	Respond shared; Supply data: → I	

02/07

ECE/CS 757, copyright J. E. Smith, 2007

MOESI Example

Thread Event	Bus Action	Data From	Global State	local states		
				C0	C1	C2
0. Initially:			<0,0,0,1>	I	I	I
1. T0 read→	CR	Memory	<1,0,0,1>	E	I	I
2. T0 write→	none		<1,0,0,0>	M	I	I
3. T2 read→	CR	C0	<1,0,1,0>	O	I	S
4. T1 write→	CRM	C0	<0,1,0,0>	I	M	I

02/07

ECE/CS 757, copyright J. E. Smith, 2007

Update Protocols

- Basic idea:
 - All writes (updates) are made visible to all caches:
 - (address,value) tuples sent "everywhere"
 - Similar to write-through protocol for uniprocessor caches
 - Obviously not scalable beyond a few processors
 - No one actually builds machines this way
- Simple optimization
 - Send updates to memory/directory
 - Directory propagates updates to all known copies: less bandwidth
- Further optimizations: combine & delay
 - Write-combining of adjacent updates (if consistency model allows)
 - Send write-combined data
 - Delay sending write-combined data until requested
- Logical end result
 - Writes are combined into larger units, updates are delayed until needed
 - Effectively the same as invalidate protocol

© 2005 Mikko Lipasti

25

Update Protocol: Dragon

- Dragon (developed at Xerox PARC)
- 5-State Protocol
 - Invalid: $\langle 0, X, X, \dots, X \rangle$
 - Some say no invalid state – due to confusion regarding empty frame versus invalid line state
 - Exclusive: $\langle 1, 0, 0, \dots, 1 \rangle$
 - Shared-Clean (Sc): $\langle 1, X, X, \dots, X \rangle$ memory may not be up-to-date
 - Shared-Modified (Sm): $\langle 1, X, X, X, \dots, 0 \rangle$ memory not up-to-date; only one copy in Sm
 - Modified: $\langle 1, 0, 0, \dots, 0 \rangle$
- Includes Cache Update action
- Includes Cache Writeback action
- Bus includes Shared flag
 - Appears to also require memory inhibit signal
 - Distinguish shared case where cache (not memory) supplies data

02/07

ECE/CS 757, copyright J. E. Smith, 2007

26

Dragon State Diagram

Action and Next State					
Current State	Processor Read	Processor Write	Eviction	Cache Read	Cache Update
I	Cache Read If no sharers: → E If sharers: → Sc	Cache Read If no sharers: → M If sharers: → Sm Cache Update → Sm		→ I	→ I
Sc	No Action → Sc	Cache Update If no sharers: → M If sharers: → Sm	No Action → I	Respond Shared; → Sc	Respond shared; Update copy; → Sc
E	No Action → E	No Action → M	No Action → I	Respond shared; Supply data → Sc	
Sm	No Action → Sm	Cache Update If no sharers: → M If sharers: → Sm	Cache Write-back → I	Respond shared; Supply data; → Sm	Respond shared; Update copy; → Sc
M	No Action → M	No Action → M	Cache Write-back → I	Respond shared; Supply data; → Sm	

02/07

ECE/CS 757, copyright J. E. Smith, 2007

Example

Thread Event	Bus Action	Data From	Global State	local states C0 C1 C2
0. Initially:			$\langle 0, 0, 0, 1 \rangle$	I I I
1. T0 read→	CR	Memory	$\langle 1, 0, 0, 1 \rangle$	E I I
2. T0 write→	none		$\langle 1, 0, 0, 0 \rangle$	M I I
3. T2 read→	CR	C0	$\langle 1, 0, 1, 0 \rangle$	Sm I Sc
4. T1 write→	CR, CU	C0	$\langle 1, 1, 1, 0 \rangle$	Sc Sm Sc
5. T0 read→	none (hit)	C0	$\langle 1, 1, 1, 0 \rangle$	Sc Sm Sc

- Appears to require atomic bus cycles CR, CU on write to invalid line

02/07

ECE/CS 757, copyright J. E. Smith, 2007

Update Protocol: Firefly

- Developed at DEC by ex-Xerox people
- 5-State Protocol
 - Similar to Dragon – different state naming based on shared/exclusive and clean/dirty
 - Invalid: $\langle 0, X, X, \dots, X \rangle$
 - EC: $\langle 1, 0, 0, \dots, 1 \rangle$
 - SC: $\langle 1, X, X, \dots, X \rangle$ memory may not be up-to-date
 - EM: $\langle 1, 0, 0, \dots, 0 \rangle$
 - SM: $\langle 1, X, X, X, \dots, 0 \rangle$ memory not up-to-date; only one copy in Sm
- Performs write-through updates (different from Dragon)

02/07

ECE/CS 757, copyright J. E. Smith, 2007

29

Firefly State Diagram

Action and Next State					
Current State	Processor Read	Processor Write	Eviction	Cache Read	Cache Read&M
I	Cache Read If no sharers: → E If sharers: → Sc	Cache Read If no sharers: → Em If sharers: → Sm Cache Update → Sm		No Action → I	No Action → I
Sc	No Action → Sc	Cache Read&M If no sharers: → E If sharers: → Sc	No Action → I	Respond Shared; → Sc	Respond Shared → Sc
Ec	No Action → E	No Action → Em	No Action → I	Respond shared; → Sc	Respond Shared → Sc
Sm	No Action → Sm	Cache Read&M If no sharers: → E If sharers: → Sc	Cache Write-back → I	Respond shared; Supply data; → Sm	Respond shared; Supply data; → Sc
Em	No Action → Em	No Action → Em	Cache Write-back → I	Respond shared; Supply data; → Sm	Respond shared; Supply data; → Sc

02/07

ECE/CS 757, copyright J. E. Smith, 2007

Update vs Invalidate

- [Weber & Gupta, ASPLOS3]
 - Consider sharing patterns
- No Sharing
 - Independent threads
 - Coherence due to thread migration
 - Update protocol performs many wasteful updates
- Read-Only
 - No significant coherence issues; most protocols work well
- Migratory Objects
 - Manipulated by one processor at a time
 - Often protected by a lock
 - Usually a write causes only a single invalidation
 - E state useful for Read-modify-Write patterns
 - Update protocol could proliferate copies

02/07

ECE/CS 757, copyright J. E. Smith, 2007

31

Update vs Invalidate, contd.

- Synchronization Objects
 - Locks
 - Update could reduce spin traffic invalidations
 - Test&Test&Set w/ invalidate protocol would work well
- Many Readers, One Writer
 - Update protocol may work well, but writes are relatively rare
- Many Writers/Readers
 - Invalidate probably works better
 - Update will proliferate copies
- What is used today?
 - Invalidate is dominant
 - CMP *may* change this assessment
 - more on-chip bandwidth

02/07

ECE/CS 757, copyright J. E. Smith, 2007

32

Nasty Realities

- State diagram is for (ideal) protocol assuming instantaneous and actions
- In reality controller implements more complex diagrams
 - A protocol state transition may be started by controller when bus activity changes local state
 - Example: an upgrade pending (for bus) when an invalidate for same line arrives

02/07

ECE/CS 757, copyright J. E. Smith, 2007

35

Partial Implementation State Table

Current State	Action and Next State						
	Processor Read	Processor Write	Bus Grant	Bus Response	Cache Read	Cache Read&M	Cache Upgrade
I	Request Bus → IR	Request Bus → IW			No Action → I	No Action → I	No Action → I
S	No Action → S	Request Bus → SW			Respond Shared: → S	No Action → I	No Action → I
E	No Action → E	No Action → M			Respond Shared: → S	No Action → I	
M	No Action → M	No Action → M			Respond dirty: Write back data; → S	Respond dirty: Write back data; → I	
IR			Cache Read → IRK				
IW			Cache Read&M → IWR				
IRK				If no shatters: → E If shatters: → S Load line			
IWR				→ M Load line			
SW			Cache Upgrade → M		Respond Shared: → SW	No Action → IW	No Action → IW

Further Optimizations

- Observation: Shared blocks should only be fetched from memory once
 - If I find a shared block on chip, forward the block
 - Problem: multiple shared blocks possible, who forwards?
 - Everyone? Power/bandwidth wasted
 - Single forwarder, but who?
 - Last one to receive block: F state
 - I->F for requestor, F->S for forwarder
 - What if F block is evicted?
 - Favor F blocks in replacement?
 - Don't allow silent eviction (force some other node to be F)
 - Fall back on memory copy if can't find F copy
- Very old idea (IBM machines have done this for a long time), but recent Intel patent issued anyway [Hum/Goodman]

© 2005 Mikko Lipasti

35

Further Optimizations

- Observation: migratory data often "flies by"
 - Add T (transition) state to protocol
 - Tag is still valid, data isn't
 - Data can be snarfed as it flies by
 - Only works with certain kinds of interconnect networks
 - Replacement policy issues
- Many other optimizations are possible
 - Literature extends 25 years back
 - Many unpublished (but implemented) techniques as well

© 2005 Mikko Lipasti

36

Implementing Cache Coherence

- Snooping implementation
 - Origins in shared-memory-bus systems
 - All CPUs could observe all other CPUs requests on the bus; hence “snooping”
 - Bus Read, Bus Write, Bus Upgrade
 - React appropriately to snooped commands
 - Invalidate shared copies
 - Provide up-to-date copies of dirty lines
 - Flush (writeback) to memory, or
 - Direct intervention (*modified intervention* or *dirty miss*)
- Snooping suffers from:
 - Scalability: shared busses not practical
 - Ordering of requests without a shared bus
 - Lots of recent and on-going work on scaling snoop-based systems

© 2005 Mikko Lipasti 37

Snooping Cache Coherence

- Basic idea: broadcast snoop to all caches to find owner
- Not scalable?
 - Address traffic roughly proportional to square of number of processors
 - Current implementations scale to 64/128-way (Sun/IBM) with multiple address-interleaved broadcast networks
- Inbound snoop bandwidth: big problem

$$OutboundSnoopRate = s_o = \langle CacheMissRate \rangle + \langle BusUpgradeRate \rangle$$

$$InboundSnoopRate = s_i = n \times s_o$$

© 2005 Mikko Lipasti 38

Snoop Bandwidth

- Snoop filtering of various kinds is possible
- Filter snoops at sink: Jetty filter [Moshovos et al., HPCA 2001]
 - Check small “filter cache” that summarizes contents of local cache
 - Avoid power-hungry lookups in each tag array
- Filter snoops at source: Multicast snooping [Bilir et al., ISCA 1999]
 - Predict likely sharing set, snoop only predicted sharers
 - Double-check at directory to make sure
- Filter snoops at source: Region coherence
 - Concurrent work: [Cantin/Smith/Lipasti, ISCA 2005; Moshovos, ISCA 2005]
 - Check larger region of memory on every snoop; remember when no sharers
 - Snoop only on first reference to region, or when region is shared
 - Eliminate 60%+ of all snoops

© 2005 Mikko Lipasti 39

Snoop Latency

- Snoop latency:
 - Must reach all nodes, return and combine responses
 - Probably the bigger scalability issue in the future
 - Topology matters: ring, mesh, torus, hypercube
 - No obvious solutions
- Parallelism: fundamental advantage of snooping
 - Broadcast exposes parallelism, enables speculative latency reduction

LDir	XSnp	RDir	XRsp	CRsp	XRd	RDat	XDat	UDat
		RDat	XDat			UDat		
		RDat	XDat	UDat				
		RDat		XDat	UDat			

© 2005 Mikko Lipasti 40

Scaleable Cache Coherence

- Eschew physical bus but still snoop
 - Point-to-point tree structure (indirect) or ring
 - Root of tree or ring provide ordering point
 - Use some scalable network for data (ordering less important)
- Or, use level of indirection through directory
 - Directory at memory remembers:
 - Which processor is “single writer”
 - Which processors are “shared readers”
 - Level of indirection has a price
 - Dirty misses require 3 hops instead of two
 - Snoop: Requestor->Owner
 - Directory: Requestor->Directory->Owner

© 2005 Mikko Lipasti 41

Implementing Cache Coherence

- Directory implementation
 - Extra bits stored in memory (directory) record state of line
 - Memory controller maintains coherence based on the current state
 - Other CPUs’ commands are not snooped, instead:
 - Directory forwards relevant commands
 - Powerful filtering effect: only observe commands that you need to observe
 - Meanwhile, bandwidth at directory scales by adding memory controllers as you increase size of the system
 - Leads to very scalable designs (100s to 1000s of CPUs)
- Directory shortcomings
 - Indirection through directory has latency penalty
 - If shared line is dirty in other CPU’s cache, directory must forward request, adding latency
 - This can severely impact performance of applications with heavy sharing (e.g. relational databases)

© 2005 Mikko Lipasti 42

Directory Protocol Implementation

- Basic idea: Centralized directory keeps track of data location(s)
- Scalable
 - Address traffic roughly proportional to number of processors
 - Directory & traffic can be distributed with memory banks (interleaved)
 - Directory cost (SRAM) or latency (DRAM) can be prohibitive
- Presence bits track sharers
 - Full map (N processors, N bits): cost/scalability
 - Limited map (limits number of sharers)
 - Coarse map (identifies board/node/cluster; must use broadcast)
- Vectors track sharers
 - Point to shared copies
 - Fixed number, linked lists (SCI), caches chained together
 - Latency vs. cost vs. scalability

© 2005 Mikko Lipasti 43

Directory Protocol Latency

LDir	XSnp	RDir	XRd	RDat	XDat	UDat
-------------	-------------	-------------	------------	-------------	-------------	-------------

- Access to non-shared data
 - Overlap directory read with data read
 - Best possible latency given NUMA arrangement
- Access to shared data
 - Dirty miss, modified intervention
 - Shared intervention?
 - If DRAM directory, no gain
 - If directory cache, possible gain (use F state)
 - No inherent parallelism
 - Indirection adds latency
 - Minimum 3 hops, often 4 hops

© 2005 Mikko Lipasti 44

Directory-based Cache Coherence

- An alternative for large, scalable MPs
- Can be based on any of the protocols discussed thus far
 - We will use MSI
- Memory Controller becomes an active participant
- Sharing info held in memory directory
 - Directory may be distributed
- Use point-to-point messages
- Network is not totally ordered

02/07 ECE/CS 757; copyright J. E. Smith, 2007

Example: Simple Directory Protocol

- Local cache controller states
 - M, S, I as before
- Local directory states
 - Shared: <1,X,X,...,1>; one or more proc. has copy; memory is up-to-date
 - Modified: <0,1,0,...,0> one processor has copy; memory does not have a valid copy
 - Uncached: <0,0,...,0,1> none of the processors has a valid copy
- Directory also keeps track of sharers
 - Can keep global state vector in full
 - e.g. via a bit vector

02/07 ECE/CS 757; copyright J. E. Smith, 2007 46

Example

- Local cache suffers load miss
- Line in remote cache in M state
 - It is the owner
- Four messages send over network
 - Cache read from local controller to home memory controller
 - Memory read to remote cache controller
 - Owner data back to memory controller; change state to S
 - Memory data back to local cache; change state to S

02/07 ECE/CS 757; copyright J. E. Smith, 2007

Cache Controller State Diagram

- Intermediate States for clarity

Current State	Cache Controller Actions and Next States							
	from Processor Side			from Memory Side				
	Processor Read	Processor Write	Eviction	Memory Read	Memory Read&M	Memory Invalidate	Memory Upgrade	Memory Data
I	Cache Read → F	Cache Read&M → F*				No Action → I		
S	No Action → S	Cache Upgrade → S*	No Action → I			Invalidate Frame: Cache ACK; → I		
M	No Action → M	No Action → M	Cache Write-back → I	Owner Data; → S	Owner Data; → I	Invalidate Frame: Cache ACK; → I		
F*								Fill Cache → S
F*								Fill Cache → M
S*							No Action → M	

02/07 ECE/CS 757; copyright J. E. Smith, 2007 48

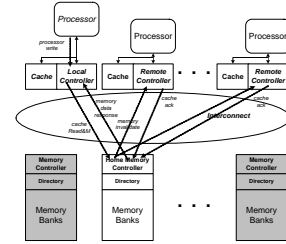
Memory Controller State Diagram

Current Directory State	Memory Controller Actions and Next States			response from Remote Cache Controller		
	Cache Read	Cache Read&M	Cache Upgrade	Data Write-back	Cache ACK	Owner Data
<i>I</i>	Memory Data; Add Requestor to Sharers; → S	Memory Data; Add Requestor to Sharers; → M				
<i>S</i>	Memory Data; Add Requestor to Sharers; → S	Memory Invalidate All Sharers; → M'	Memory Upgrade All Sharers; → M'	No Action → I		
<i>M</i>	Memory Read from Owner; → S'	Memory Read&M; to Owner; → M'		Make Sharers Empty; → I		
<i>S'</i>						Memory Data to Requestor; Write memory; Add Requestor to Sharers; → S
<i>M'</i>					When all ACKS Memory Data; → M	Memory Data to Requestor; → M
<i>M''</i>					When all ACKS then → M	

49

Another Example

- Local write (miss) to shared line
- Requires invalidations and acks



02/07

ECE/CS 757, copyright J. E. Smith, 2007

Example Sequence

- Similar to earlier sequences

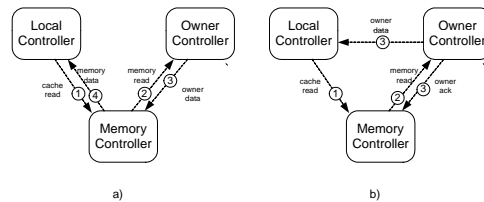
Thread Event	Controller Actions	Data From	global state	local states: C0 C1 C2		
0. Initially:			<0,0,0,1>	I	I	I
1. T0 read→	CR,MD	Memory	<1,0,0,1>	S	I	I
2. T0 write→	CU, MU*,MD		<1,0,0,0>	M	I	I
3. T2 read→	CR,MR,MD	C0	<1,0,1,1>	S	I	S
4. T1 write→	CRM,MI,CAM D	Memory	<0,1,0,0>	I	M	I

02/07

ECE/CS 757, copyright J. E. Smith, 2007

Variation: Three Hop Protocol

- Have owner send data directly to local controller
- Owner Acks to Memory Controller in parallel



02/07

ECE/CS 757, copyright J. E. Smith, 2007

Directory Protocol Optimizations

- Remove dead blocks from cache:
 - Eliminate 3- or 4-hop latency
 - Dynamic Self-Invalidation [Lebeck/Wood, ISCA 1995]
 - Last touch prediction [Lai/Falsafi, ISCA 2000]
 - Dead block prediction [Lai/Fide/Falsafi, ISCA 2001]
- Predict sharers
 - Prediction in coherence protocols [Mukherjee/Hill, ISCA 1998]
 - Instruction-based prediction [Kaxiras/Goodman, ISCA 1999]
 - Sharing prediction [Lai/Falsafi, ISCA 1999]
- Hybrid snooping/directory protocols
 - Improve latency by snooping, conserve bandwidth with directory
 - Multicast snooping [Billir et al., ISCA 1999; Martin et al., ISCA 2003]
 - Bandwidth-adaptive hybrid [Martin et al., HPCA 2002]
 - Token Coherence [Martin et al., ISCA 2003]
 - VCT Multicasting [Enright Jerger thesis 2008]

© 2005 Mikko Lipasti

53

Protocol Races

- Atomic bus
 - Only stable states in protocol (e.g. M, S, I)
 - All state transitions are atomic (I→M)
 - No conflicting requests can interfere since bus is held till transaction completes
 - Distinguish coherence transaction from data transfer
 - Data transfer can still occur much later; easier to handle this case
- Atomic buses don't scale
 - At minimum, separate bus request/response
- Large systems have broadly variable delays
 - Req/resp separated by dozens of cycles
 - Conflicting requests can and do get issued
 - Messages may get reordered in the interconnect
- How do we resolve them?

54

Resolving Protocol Races

- Req/resp decoupling introduces transient states
 - E.g. I->S is now I->ItoX->ItoS_nodata->S
- Conflicting requests to blocks in transient states
 - NAK – ugly; livelock, starvation potential
 - Keep adding more transient states ...
- Directory protocol makes this a bit easier
 - Can order at directory, which has full state info
 - Even so, messages may get reordered

55

Common Protocol Races

- Read strings: P0 read, P1 read, P2 read
 - Easy, since read is nondestructive
 - Can rely on F state to reduce DRAM accesses
 - Forward reads to previous requestor (F)
- Write strings: P0 write, P1 write, P2 write
 - Forward P1 write req to P0 (M)
 - P0 completes write then forwards M block to P1
 - Build string of writes (write string forwarding)
- Read after write (similar to prev. WAW)
- Writeback race: P0 evicts dirty block, P1 reads
 - Dirty block is in the network (no copy at P0 or at dir)
 - NAK P1, or force P0 to keep copy till dir ACKs WB
- Many others crop up, esp. with optimizations

56

Summary

- Coherence States
- Snoopy bus-based Invalidate Protocols
- Invalidate protocol optimizations
- Update Protocols (Dragon/Firefly)
- Directory protocols
- Implementation issues

57