

332 Advanced Computer Architecture Chapter 6

Parallel architectures, shared memory, and cache coherency

March 2006 Paul H J Kelly

These lecture notes are partly based on the course text, Hennessy and Patterson's *Computer Architecture*, a quantitative approach (3^{d} *ed*), and on the lecture slides of David Patterson, John Kubiatowicz and Yujia Jin at Berkeley

Why add another processor?

- How should they be connected I/O, memory bus, L2 cache, registers?
- le Cache coherency the problem
- M Coherency what are the rules?
- Coherency using broadcast update/invalidate
- Coherency using multicast SMP vs ccNUMA
- Distributed directories, home nodes, ownership; the ccNUMA design space
- Beyond ccNUMA; COMA and Simple COMA
- Mybrids and clusters

Advanced Computer Architecture Chapter 6.2

London e-Science Centre's "Mars" cluster



408 *C*PUs:

- 72 dual 1.8GHz Opteron processors, 2Gb memory, 80Gb S-ATA disk, Infiniband
- 40 dual 1.8GHz Opteron processors, 4Gb memory, 80Gb S-ATA disk, Gigabit Ethernet
- 88 dual 1.8GHz Opteron processors, 2Gb memory, 80Gb S-STA disk, Gigabit Ethernet
- 4 dual 2.2GHz Opteron processors, 4Gb memory, 36Gb SCSI disk, Gigabit Ethernet
- I Mellanox MTS9600 Infiniband switch configured with 72 ports
- 5 Extreme Networks Summit 400-48t 48 port Gigabit Switches

NPACI Blue Horizon, San Diego Supercomputer Center

- 144 eight-processor SMP High Nodes perform the primary compute functions.
- 12 two-processor SMP High Nodes perform service functions.
- 1,152 Power3 processors in the compute nodes run at 375 MHz. Each processor has a peak performance of 1.500 Mflops (millions of floating-point operations per second)
- 576 gigabytes of total system memory, at 4 GB per compute node
- 5.1 terabytes -- 5,100 gigabytes of associated disk



- 42 towers, or frames, house the compute and service nodes. Four racks house the disk.
- ▶ 1,500 square feet
- Programmed using MPI ("Message Passing Interface")

HPCx - Supercomputing service for UK science

- Located at Daresbury Labs, Cheshire, operated by University of Edinburgh Parallel Computing Centre
- **IBM eServer 575 nodes**, each 16-way SMP Power5
- ▶ 96 nodes for compute jobs for users, a total 1536 processors
- Four chips (8 processors) are integrated into a multi-chip module (MCM)
- Two MCMs (16 processors) comprise one frame
- Each MCM is configured with 128 MB of L3 cache and 16 GB of main memory
- Total main memory of 32 GB per frame shared between 16 processors of frame
- frames connected via IBM's High Performance Switch (HPS)



Advanced Computer Architecture Chapter 6.

SGI Origin 3800 at SARA, Netherlands



- 1024-CPU system consisting of two 512-CPU <u>SGI Origin 3800</u> systems. Peak performance of 1 TFlops (10¹² floating point operations) per second. 500MHz R14000 CPUs organized in 256 4-CPU nodes
- 🕨 1 TByte of RAM. 10 TByte of on-line storage & 100 TByte near-line storage
- 🕨 45 racks, 32 racks containing CPUs & routers, 8 I/O racks & 5 racks for disks
- ▶ Each 512-CPU machine offers application program a single, shared memory image

Advanced Computer Architecture Chapter 6.6





Livermore Labs "classified service in support of the National Nuclear Security Administration's stockpile science mission"

Bluegene/L at LLNL

- 1024 nodes/cabinet, 2 CPUs per node
- 🕨 65,536 nodes in total
- 32 x 32 x 64 3D torus interconnect
- 1,024 gigabit-per-second links to a global parallel file system to
- support fast input/output to disk 1.5MWatts





-							
	Manufacturer	Computer	[TF/s]	Installation Site	Country	Year	#Proc
1	IBM	BlueGene/L eServer Blue Gene	280.6	DOE/NNSA/LLNL	USA	2005	131072
2	IBM	BGW eServer Blue Gene	91.29	IBM Thomas Watson	USA	2005	40960
3	IBM	ASC Purple eServer pSeries p575	63.39	DOE/NNSA/LLNL	USA	2005	10240
4	SGI	Columbia Altix, Infiniband	51.87	NASA Ames	USA	2004	10160
5	Dell	Thunderbird	38.27	Sandia	USA	2005	8000
6	Cray	Red Storm Cray XT3	36.19	Sandia	USA	2005	10880
7	NEC	Earth-Simulator	35.86	Earth Simulator Center	Japan	2002	5120
8	IBM	MareNostrum BladeCenter JS20, Myrinet	27.91	Barcelona Supercomputer Center	Spain	2005	4800
9	IBM	eServer Blue Gene	27.45	ASTRON University Groningen	Netherlands	2005	12288
10	Cray	Jaguar Cray XT3	20.53	Oak Ridge National Lab	USA	2005	5200

What are parallel computers used for?

bExecuting loops in parallel

- Improve performance of single application
- Barrier synchronisation at end of loop
- Iteration i of loop 2 may read data produced by iteration i of loop 1 - but perhaps also from other iterations
- Example: NaSt3DGP

WHigh-throughput servers

- Eg. database, transaction processing, web server, ecommerce
- Improve performance of single application
- Consists of many mostly-independent transactions
- Sharing data
- Synchronising to ensure consistency
- Transaction roll-back

Mixed, multiprocessing workloads

Advanced Computer Architecture Chapter 6.10

Why add another processor? Further simultaneous instruction issue slots rarely usable in real code Smallest working CPU Number of transistors

- Increasing the complexity of a single CPU leads to diminishing returns
 - Due to lack of instruction-level parallelism
 - ➡ Too many simultaneous accesses to one register file
 - Forwarding wires between functional units too long inter-cluster communication takes >1 cycle

Architectural effectiveness of Intel processors



- Architectural effectiveness shows performance gained through architecture rather than clock rate
- Extra transistors largely devoted to cache, which of course is essential in allowing high clock rate Source: http://www.galaxycentre.com/intelcpu.htm and Intel

Page 3



42M transistors



Avances Computer Architecture Chapter (

How to add another processor?

- Idea: instead of trying to exploit more instructionlevel parallelism by building a bigger CPU, build two or more
- This only makes sense if the application parallelism exists...
- Why might it be better?
 - No need for multiported register file
 - ➡ No need for long-range forwarding
 - ➡ CPUs can take independent control paths
 - Still need to synchronise and communicate
 - Program has to be structured appropriately...

How to add another processor?

- In How should the CPUs be connected?
- Idea: systems linked by network connected via I/O bus
 Eg Fujitsu AP3000, Myrinet, Quadrics
- Idea: CPU/memory packages linked by network connecting main memory units
 - Eg SGI Origin
- Idea: CPUs share main memory
 - Eg Intel Xeon SMP
- Idea: CPUs share L2/L3 cache
 - Eg IBM Power4
- Idea: CPUs share L1 cache
- Model: CPUs share registers, functional units
 - Cray/Tera MTA (multithreaded architecture), Symmetric multithreading (SMT), as in Hyperthreaded Pentium 4, Alpha 21464, etc

Advanced Computer Architecture Chapter 6.15

How to program a parallel computer?





- First loop operates on rows in parallel
- Second loop operates on columns in parallel
- With distributed memory we would have to program message passing to transpose the array in between
- With shared memory... no problem!

}

ed Computer Architecture Chapter 6.17

Shared-memory parallel - OpenMP

- OpenMP is a standard design for language extensions for shared-memory parallel programming
- Language bindings exist for Fortran, C and to some extent (eg research prototypes) for C++, Java and C#
- Implementation requires compiler support

▶ Example:

```
for(i=0; I<N; ++i)
#pragma omp parallel for
for(j=0; j<M; ++j)
A[i,j] = (A[i-1,j] + A[i,j])*0.5;
#pragma omp parallel for
par_for(i=0; I<N; ++i)
for(j=0; j<M; ++j)
A[i,j] = (A[i,j-1] + A[i,j])*0.5;</pre>
```

```
Advanced Computer Architecture Chapter 6.18
```

Implementing shared-memory parallel loop



More OpenMP

#pragma omp parallel for \
 default(shared) private(i) \
 schedule(static,chunk) \
 reduction(+:result)
for (i=0; i < n; i++)
 result = result + (a[i] * b[i]);</pre>

Mefault(shared) private(i):

All variables except i are shared by all threads.

w schedule(static,chunk):

Iterations of the parallel loop will be distributed in equal sized blocks to each thread in the "team"

le reduction(+:result):

performs a reduction on the variables that appear in its argument list

A private copy for each variable is created for each thread. At the end of the reduction, the reduction operator is applied to all private copies of the shared variable, and the final result is written to the global shared variable.

http://www.llnl.gov/computing/tutorials/openMP/#REDUCTION

Distributed-memory parallel - MPI

- MPI ("Message-passing Interface) is a standard API for parallel programming using message passing
- **W** Usually implemented as library
- Six basic calls:
 - MPI_Init Initialize MPI
 - MPI_Comm_size Find out how many processes there are
 - MPI_Comm_rank Find out which process I am
 - MPI_Send Send a message
 - MPI_Recv Receive a message
 - MPI_Finalize Terminate MPI
- **W** Key idea: collective operations
 - MPI_Bcast broadcast data from the process with rank "root" to all other processes of the group.
 - MPI_Reduce combine values on all processes into a single value using the operation defined by the parameter op.

Advanced Computer Architecture Chapter 6.21

MPI Example: initialisation

SPMD

- "Single Program, Multiple Data"
- Each processor executes the program
- First has to work out what part it is to play
- "myrank" is index of this CPU
- "comm" is MPI "communicator" abstract index space of p processors
- In this example, array is partitioned into slices



! Compute number of processes and myrank CALL MPI COMM SIZE(comm, p. jerr) CALL MPI COMM RANK(comm, myrank, jerr) ! compute size of local block m = n/pIF (myrank.LT.(n-p*m)) THEN m = m+1END IF ! Compute neighbors IF (myrank, EQ.0) THEN left = MPI_PROC_NULL ELSE left = myrank - 1 END IF IF (myrank.EQ.p-1)THEN right = MPI PROC NULL ELSE right = myrank+1 END IF ! Allocate local arrays ALLOCATE (A(0:n+1,0:m+1), B(n,m))

ttp://www.netlib.org/utk/papers/mpi-book/node51.html

Advanced Computer Architecture Chapter 6.22

Main Loop



How to connect processors...

🕨 Tradeoffs:

- close coupling to minimise delays incurred when processors interact
- separation to avoid contention for shared resources
- 🗽 Result:
 - spectrum of alternative approaches based on application requirements, cost, and packaging/integration issues
- **b** Currently:
 - just possible to integrate 2 full-scale CPUs on one chip together with large shared L2 cache
 - common to link multiple CPUs on same motherboard with shared bus connecting to main memory
 - more aggressive designs use richer interconnection network, perhaps with cacheto-cache transfer capability



Multiple caches... and trouble

▶ Suppose processor 0 loads memory location ×

▶ × is fetched from main memory and allocated into processor O's cache(s)

Advanced Computer Architecture Chapter 6.25





In Suppose processor 1 loads memory location X

▶ × is fetched from main memory and allocated into processor 1's cache(s) as well



- Processor O's cached copy of X is updated
- Processor 1 continues to used the old value of X

Multiple caches... and trouble



▶ Suppose processor 2 loads memory location ×

How does it know whether to get x from main memory, processor 0 or processor 1?

Implementing distributed, shared memory

- 🕨 Two issues:
- 1. How do you know where to find the latest version of the cache line?
- 2. How do you know when you can use your cached copy - and when you have to look for a more up-to-date version?
- We will find answers to this after first thinking about what a distributed shared memory implementation is supposed to do...

Cache consistency (aka cache coherency)

- ▶ Goal (?):
 - "Processors should not continue to use out-of-date data indefinitely"
- **Goal (?)**:
 - "Every load instruction should yield the result of the most recent store to that address"

Goal (?): (definition: Sequential Consistency)

"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"

(Leslie Lamport, "How to make a multiprocessor computer that correctly executes multiprocess programs" (IEEE Trans Computers Vol.C-28(9) Sept 1979)

Implementing Strong Consistency: update

Idea #1: when a store to address x occurs, update all the remote cached copies

Interview Provide States and Sta

- To broadcast every store to every remote cache
- Or to keep a list of which remote caches hold the cache line
- Or at least keep a note of whether there are any remote cached copies of this line
- But first...how well does this update idea work?

Implementing Strong Consistency: update...

Problems with update

- 1. What about if the cache line is several words long?
 - Each update to each word in the line leads to a broadcast
- 2. What about old data which other processors are no longer interested in?
 - We'll keep broadcasting updates indefinitely...
 - Do we really have to broadcast every store?
 - It would be nice to know that we have exclusive access to the cacheline so we don't have to broadcast updates...

A more cunning plan... invalidation

- Suppose instead of updating remote cache lines, we invalidate them all when a store occurs?
- After the first write to a cache line we know there are no remote copies – so subsequent writes don't lead to communication
- Is invalidate always better than update?
 - 🟓 Often
 - But not if the other processors really need the new data as soon as possible
- Note that to exploit this, we need a bit on each cache line indicating its sharing state
- (analogous to write-back vs write-through caches)

Update:

May reduce latency of subsequent remote reads

Update vs invalidate

lif any are made

Invalidate:

- May reduce network traffic
- e.g. if same CPU writes to the line before remote nodes take copies of it

Four cache line states:

Broadcast invalidations on bus unless cache line is exclusively "owned" (DIRTY)

• Read miss:

- If another cache has the line in SHARED-DIRTY or DIRTY,
 - it is supplied
 - changing state to SHARED-DIRTY
- Otherwise
 - the line comes from memory. The state of the
 - · line is set to VALID

- INVALID
- VALID : clean, potentially shared, unowned
- SHARED-DIRTY : modified, possibly shared, owned
- DIRTY : modified, only copy, owned

• Write hit:

• No action if line is DIRTY

The "Berkeley" Protocol

- If VALID or SHARED-DIRTY,
 - · an invalidation is sent, and
 - · the local state set to DIRTY

• Write miss:

- Line comes from owner (as with read miss).
- All other copies set to INVALID, and line in requesting cache is set to DIRTY

Advanced Computer Architecture Chapter 6.35



- 3. SHARED-DIRIT: modified, possibly shared, owne
- 4. DIRTY: modified, only copy, owned

Advanced Computer Architecture Chapter 6.36

miss"?

The job of the cache controller - snooping

- The protocol state transitions are implemented by the cache controller – which "snoops" all the bus traffic
- Transitions are triggered either by
 - the bus (Bus invalidate, Bus write miss, Bus read miss)
 - The CPU (Read hit, Read miss, Write hit, Write miss)
- For every bus transaction, it looks up the directory (cache line state) information for the specified address
 - If this processor holds the only valid data (DIRTY), it responds to a "Bus read miss" by providing the data to the requesting CPU
 - If the memory copy is out of date, one of the CPUs will have the cache line in the SHARED-DIRTY state (because it updated it last) – so must provide data to requesting CPU
 - State transition diagram doesn't show what happens when a cache line is displaced...

Advanced Computer Architecture Chapter 6.37

Berkeley protocol - summary

- Invalidate is usually better than update
- Cache line state "DIRTY" bit records whether remote copies exist
 - If so, remote copies are invalidated by broadcasting message on bus cache controllers snoop all traffic
- Where to get the up-to-date data from?
 - Broadcast read miss request on the bus
 - ➡ If this CPUs copy is DIRTY, it responds
 - ➡ If no cache copies exist, main memory responds
 - If several copies exist, the CPU which holds it in "SHARED-DIRTY" state responds
 - ▶ If a SHARED-DIRTY cache line is displaced, ... need a plan
- Mow well does it work?
 - See extensive analysis in Hennessy and Patterson

Advanced Computer Architecture Chapter 6.38



Snooping Cache Variations

Basic	Berkeley	Illinois	MESI
Protocol	Protocol	Protocol	Protocol
	Owned Exclusive	Private Dirty	Modified (private,!=Memory)
Exclusive 2	—Owned Shared	Private Clean	eXclusive (private,=Memory)
Shared	Shared	Shared	<u>Shared (shared,=Memory)</u>
Invalid	Invalid	Invalid	<u>I</u> nvalid

Owner can update via bus invalidate operation Owner must write back when replaced in cache

> If read sourced from memory, then Private Clean If read sourced from other cache, then Shared Can write in cache if held private clean or dirty

dvanced Computer Architecture Chapter 6.39

Implementation Complications

Write Races:

- Cannot update cache until bus is obtained
- Otherwise, another processor may get bus first, and then write the same cache block!
- Two step process:
 - Arbitrate for bus
- Place miss on bus and complete operation
- If miss occurs to block while waiting for bus,
- handle miss (invalidate may be needed) and then restart. Split transaction bus:
- Bus transaction is not atomic:
 - can have multiple outstanding transactions for a block
- Multiple misses can interleave, allowing two caches to grab block in the Exclusive state
- Must track and prevent multiple misses for one block

Must support interventions and invalidations

Multiple processors must be on bus, access to both addresses and data

Implementing Snooping Caches

- Add a few new commands to perform coherency, in addition to read and write
- Processors continuously snoop on address bus
 - If address matches tag, either invalidate or update
- Since every bus transaction checks cache tags, could interfere with CPU just to check:
 - solution 1: duplicate set of tags for L1 caches just to allow checks in parallel with CPU
 - solution 2: L2 cache already duplicate, provided L2 obeys inclusion with L1 cache
 - block size, associativity of L2 affects L1

Advanced Computer Architecture Chapter 6.42

Implementing Snooping Caches

- Bus serializes writes, getting bus ensures no one else can perform memory operation
- On a miss in a write-back cache, may have the desired copy and it's dirty, so must reply
- MAdd extra state bit to cache to determine shared or not
- Madd 4th state (MESI)

Large-Scale Shared-Memory Multiprocessors

Bus inevitably becomes a bottleneck when many processors are used

- Use a more general interconnection network
- So snooping does not work
- DRAM memory is also distributed
 - Each node allocates space from local DRAM
 - Copies of remote data are made in cache
- Major design issues:
 - How to find and represent the "directory" of each line?
 - How to find a copy of a line?
- As a case study, we will look at S3.MP (Sun's Scalable Shared memory Multi-Processor, a CC-NUMA (cachecoherent non-uniform memory access) architecture

Advanced Computer Architecture Chapter 6.43

Larger MPs

In Separate Memory per Processor

- In Local or Remote access via memory controller
- ▶ 1 Cache Coherency solution: non-cached pages
- Alternative: <u>directory</u> per cache that tracks state of every block in every cache
 - ➡ Which caches have a copies of block, dirty vs. clean, ...
- Info per memory block vs. per cache block?
 - PLUS: In memory => simpler protocol (centralized/one location)
 - MINUS: In memory => directory is f(memory size) vs. f(cache size)
- Prevent directory as bottleneck? distribute directory entries with memory, each keeping track of which Procs have copies of their blocks

Distributed Directory MPs



Advanced Computer Architecture Chapter 6.46

Directory Protocol

In Similar to Snoopy Protocol: Three states

- Shared: ≥ 1 processors have data, memory up-to-date
- Uncached (no processor has it; not valid in any cache)
- Exclusive: 1 processor (owner) has data; memory out-of-date
- In addition to cache state, must track which processors have data when in the shared state (usually bit vector, 1 if processor has copy)
- Keep it simple(r):
 - Writes to non-exclusive data => write miss
 - Processor blocks until access completes
 - Assume messages received and acted upon in order sent

Directory Protocol

- **Involution** No bus and don't want to broadcast:
 - interconnect no longer single arbitration point
 - all messages have explicit responses
- ▶ Terms: typically 3 processors involved
- Local node where a request originates
- Home node where the memory location of an address resides
- Remote node has a copy of a cache block, whether exclusive or shared
- Example messages on next slide:
 - P = processor number, A = address

vanced Computer Architecture Chapter 6.47

	Directo	ry Protocol	Messages
Message type	Source	Destination	Msg Content
Read miss	Local cache	Home directory	Ρ, Α
Processor i make P a i	P reads data at add read sharer and arr	lress A; ange to send data baci	k
Write miss	Local cache	Home directory	Ρ, Α
Processor make P the	P writes data at ad e exclusive owner an	dress A; ad arrange to send dat	a back
Invalidate	Home directory	Remote caches	Α
📫 Invalidate	a shared copy at ad	ddress A.	
Fetch	Home directory	Remote cache	Α
🔶 Fetch the	block at address A	and send it to its hom	e directory
Fetch/Invalidate	Home directory	Remote cache	Α
Fetch the the block	block at address A in the cache	and send it to its hom	e directory; invalidate
Data value reply	Home directory	Local cache	Data
📫 Return a a	lata value from the	home memory (read m	iss response)
Data write-back	Remote cache	Home directory	A, Data
🕈 Write-bac	k a data value for d	address A (invalidate r	esponse)

State Transition Diagram for an Individual Cache Block in a **Directory Based System**

- States identical to snoopy case; transactions very similar.
- In Transitions caused by read misses, write misses, invalidates, data fetch requests
- Menerates read miss & write miss msg to home directory.
- Write misses that were broadcast on the bus for snooping => explicit invalidate & data fetch requests.
- In Note: on a write, a cache block is bigger, so need to read the full cache block



State Transition Diagram for the Directory

- Same states & structure as the transition diagram for an individual cache
- ▶ 2 actions: update of directory state & send msgs to statisfy requests
- In Tracks all copies of memory block.
- Malso indicates an action that updates the sharing set, Sharers, as well as sending a message.



Example Directory Protocol

- Message sent to directory causes two actions:
 - Update the directory
 - More messages to satisfy request
- Block is in Uncached state: the copy in memory is the current value; only possible requests for that block are:
 - Read miss: requesting processor sent data from memory & requestor made <u>only</u> sharing node; state of block made Shared.
 - Write miss: requesting processor is sent the value & becomes the Sharing node. The block is made Exclusive to indicate that the only valid copy is cached. Sharers indicates the identity of the owner.
- Block is Shared => the memory value is up-to-date:
 - Read miss: requesting processor is sent back the data from memory & requesting processor is added to the sharing set.
 - Write miss: requesting processor is sent the value. All processors in the set Sharers are sent invalidate messages, & Sharers is set to identity of requesting processor. The state of the block is made Exclusive.

Advanced Computer Architecture Chapter 6.54

Example

Example Directory Protocol

Block is Exclusive: current value of the block is held in the cache of the processor identified by the set Sharers (the owner) => three possible directory requests:

Read miss: owner processor sent data fetch message, causing state of block in owner's cache to transition to Shared and causes owner to send data to directory, where it is written to memory & sent back to requesting processor.

Identity of requesting processor is added to set Sharers, which still contains the identity of the processor that was the owner (since it still has a readable copy). State is shared.

- Data write-back: owner processor is replacing the block and hence must write it back, making memory copy up-to-date (the home directory essentially becomes the owner), the block is now Uncached, and the Sharer set is empty.
- Write miss: block has a new owner. A message is sent to old owner causing the cache to send the value of the block to the directory from which it is sent to the requesting processor, which becomes the new owner. Sharers is set to identify of new owner, and state of block is made Exclusive.

Processor 1 Processor 2 Interconnect Directory Memory

	P1			P2			Bus				Direc	tory		Memor
step	State	Addr	Value	State	Addr	Value	Action	Proc.	Addr	Value	Addr	State	{Procs}	Value
P1: Write 10 to A1														
P1: Read A1														
P2: Read A1														
2: Write 20 to A1														
P2: Write 40 to A2														

A1 and A2 map to the same cache block

Advanced Computer Architecture Chapter 6.55

Example

Example

		Proce	sso	r 1	Pro	ces	sor	2 I	nter	con	nect		Dire	ctory	Me	mory
		P1			P2			Bus				Direc	tory		Memor	
	step	State	Addr	Value	State	Addr	Value	Action	Proc.	Addr	Value	Addr	State	{Procs}	Value	
	P1: Write 10 to A1							WrMs	P1	A1		<u>A1</u>	Ex	<u>{P1}</u>		
		Excl.	A1	10				DaRp	P1	A1	0					
	P1: Read A1															
	P2: Read A1															
Р	2: Write 20 to A1															
	P2: Write 40 to A2															
																J

A1 and A2 map to the same cache block

Advanced Computer Architecture Chapter 6,57

Example





Advanced Computer Architecture Chapter 6.59

		Proce	sso	r 1	Pro	ces	sor	2	nter	con	nect		Dire	ctory	Mem
		P1			P2			Bus				Direc	tory		Memor
	step	State	Addr	Value	State	Addr	Value	Action	Proc.	Addr	Value	Addr	State	{Procs}	Value
	P1: Write 10 to A1							WrMs	P1	A1		<u>A1</u>	Ex	<u>{P1}</u>	
		Excl.	A1	10				DaRp	P1	A1	0				
	P1: Read A1	Excl.	A1	10											
	P2: Read A1														
Р	2: Write 20 to A1														
	P2: Write 40 to A2														

A1 and A2 map to the same cache block

Advanced Computer Architecture Chapter 6.58



Processor 1 Processor 2 Interconnect Directory Memory P2 y Memor Bus step State Addr Value Action Proc. Addr Value Addr P1: Write 10 to A P1 Δ1 Δ1 Ex {P1} 14/14/10 Excl. A1 10 P1 A1 P1: Read A1 Excl. A1 10 P2: Read A1 RdMs P2 A1 Shar. A1 A1 10 P 10 Shar. A1 10 DaRp P2 A1 10 A1 P2: Write 20 to A1 10 10
 Excl.
 A1
 20
 WrMs
 P2
 A1

 Inval.
 P1
 A1
 A1 P2: Write 40 to A2 10

A1 and A2 map to the same cache block

Example

		Proce	SSO	r 1	Pro	ces	sor	2	nter	con	nect		Dire	ctory	Memory
1		P1			P2			Bus				Direc	tory		Memor
	step	State	Addr	Value	State	Addr	Value	Action	Proc.	Addr	Value	Addr	State	{Procs}	Value
	P1: Write 10 to A1							WrMs	P1	A1		<u>A1</u>	Ex	<u>{P1}</u>	
		Excl.	<u>A1</u>	10				<u>DaRp</u>	P1	A1	0				
	P1: Read A1	Excl.	A1	10											
	P2: Read A1				Shar.	A1		RdMs	P2	A1					
		Shar.	A1	10				<u>Ftch</u>	P1	A1	10	A	1		<u> 10</u>
					Shar.	A1	<u>10</u>	<u>DaRp</u>	P2	A1	10	A1	Shar.	P1,P2}	10
P	2: Write 20 to A1				Excl.	A1	20	WrMs	P2	A1					10
		<u>Inv.</u>						Inval.	P1	A1		A1	Excl.	<u>{P2}</u>	10
	P2: Write 40 to A2							WrMs	P2	A2		<u>A2</u>	Excl.	<u>{P2}</u>	0
								WrBk	P2	A1	20	A1	Unca.	Û	20
					Excl.	<u>A2</u>	<u>40</u>	<u>DaRp</u>	P2	A2	0	A2	Excl.	{P2}	0

A1 and A2 map to the same cache block

Implementing a Directory

We assume operations atomic, but they are not; reality is much harder; must avoid deadlock when run out of bufffers in network (see Appendix E)

Detimizations:

read miss or write miss in Exclusive: send data directly to requestor from owner vs. 1st to memory and then from memory to requestor

Synchronization

Why Synchronize? Need to know when it is safe for different processes to use shared data

Issues for Synchronization:

- Uninterruptable instruction to fetch and update memory (atomic operation);
- User level synchronization operation using this primitive;
- For large scale MPs, synchronization can be a bottleneck; techniques to reduce contention and latency of synchronization

Uninterruptable Instruction to Fetch and Update Memory

- Atomic exchange: interchange a value in a register for a value in memory
 - 0 => synchronization variable is free
 - 1 => synchronization variable is locked and unavailable
 - 🗯 Set register to 1 & swap
 - New value in register determines success in getting lock
 - 0 if you succeeded in setting the lock (you were first)
 - I if other processor had already claimed access
 - Key is that exchange operation is indivisible
- Method Test-and-set: tests a value and sets it if the value passes the test
- Fetch-and-increment: it returns the value of a memory location and atomically increments it
 - 0 => synchronization variable is free

Advanced Computer Architecture Chapter 6.63

Uninterruptable Instruction to Fetch and Update Memory

Mard to have read & write in 1 instruction: use 2 instead

Load linked (or load locked) + store conditional

Load linked returns the initial value

Store conditional returns 1 if it succeeds (no other store to same memory location since preceeding load) and 0 otherwise

Example doing atomic swap with LL & SC:

try:	mov	R3,R4 R2,0(R1)	; mov exchange value ; load linked
	sc	R3,0(R1)	; store conditional
	mov	R4,R2	; branch store fails (RS = 0) ; but load value in R4

▶ Example doing fetch & increment with LL & SC:

try: R2,0(R1)	; load linked
addi R2,R2,#1	; increment (OK if reg-reg)
sc R2,0(R1)	; store conditional
beqz R2,try	; branch store fails (R2 = 0)

User Level Synchronization— Operation Using this Primitive

- Spin locks: processor continuously tries to acquire, spinning around a loop trying to get the lock
 - li R2,#1 lockit: exch R2,0(R1) ;atomic exchange bnez R2,lockit ;already locked?
- What about MP with cache coherency?
 - → Want to spin on cache copy to avoid full memory latency
 - Likely to get cache hits for such variables
- Problem: exchange includes a write, which invalidates all other copies; this generates considerable bus traffic
- Solution: start by simply repeatedly reading the variable; when it changes, then try exchange ("test and test&set"):

try:	li	R2,#1	
lockit:	lw bpez	R3,0(R1)	;load var
	exch	R2,0(R1)	atomic exchange
	bnez	R2,try	;already locked?

Another MP Issue: Memory Consistency Models

What is consistency? When must a processor see the new value? e.g., seems that

P1:	A = 0;	P2:	B = 0;
L1:	A = 1; if (B == 0)	L2:	B = 1; if (A == 0)
		-	

- Impossible for both if statements L1 & L2 to be true?
 What if write invalidate is delayed & processor continues?
- Memory consistency models: what are the rules for such cases?
- Sequential consistency: result of any execution is the same as if the accesses of each processor were kept in order and the accesses among different processors were interleaved => assignments before ifs above
 - SC: delay all memory accesses until all invalidates done

dvanced Computer Architecture Chapter 6.67

Memory Consistency Model

- Weak consistency schemes offer faster execution than sequential consistency
- Several processors provide fence instructions to enforce sequential consistency when an instruction stream passes such a point. Expensive!
- Not really an issue for most programs; they are synchronized
 - A program is synchronized if all access to shared data are ordered by synchronization operations

write (x)

release (s) {unlock}

acquire (s) {lock}

read(x)

- Only those programs willing to be nondeterministic are not synchronized: "data race": outcome f(proc. speed)
- Several Relaxed Models for Memory Consistency since most programs are synchronized; characterized by their attitude towards: RAR, WAR, RAW, WAW to different addresses

Summary

- Caches contain all information on state of cached memory blocks
- Snooping and Directory Protocols similar; bus makes snooping easier because of broadcast (snooping => uniform memory access)
- Directory has extra data structure to keep track of state of all cache blocks
- Distributing directory => scalable shared address multiprocessor
 - => Cache coherent, Non uniform memory access



1. Read Request

Case study: Sun's S3MP

Protocol Basics

- S3.MP uses distributed singly-linked sharing lists, with static homes
- Each line has a "home" node, which stores the root of the directory
- Requests are sent to the home node
- Home either has a copy of the line, or knows a node which does

S3MP: Read Requests



HOME





- If the line is exclusive (i.e. dirty bit is set) no message is required
- Else send a write-request to the home
 - Home sends an invalidation message down the chain
 - Each copy is invalidated (other than that of the requester)
 - Final node in chain acknowledges the requester and the home
- Chain is locked for the duration of the invalidation

S3MP - Replacements



Finding your data

Writes

- How does a CPU find a valid copy of a specified ine address's data?
 - 1. Translate virtual address to physical
 - 2. Physical address includes bits which identify "home" node
 - 3. Home node is where DRAM for this address resides
 - 4. But current valid copy may not be there - may be in another CPU's cache
 - 5. Home node holds pointer to sharing chain, so always knows where valid copy can be found

ccNUMA summary

- S3MP's cache coherency protocol implements strong consistency
 - Many recent designs implement a weaker consistency model...
- 🕨 S3MP uses a singly-linked sharing chain
 - ➡ Widely-shared data long chains long invalidations, nasty replacements
 - "Widely shared data is rare"
- 🕨 In real life:
 - ▶ IEEE Scalable Coherent Interconnect (SCI): doubly-linked sharing list
 - ♦ SGI Origin 2000: bit vector sharing list
 - Real Origin 2000 systems in service with 256 CPUs
 - ▶ Sun E10000: hybrid multiple buses for invalidations, separate switched network for data transfers

Many E10000s in service, often with 64 CPUs

Beyond ccNUMA

IF COMA: cache-only memory architecture

- ➡ Data migrates into local DRAM of CPUs where it is being used
- + Handles very large working sets cleanly
- Replacement from DRAM is messy: have to make sure someone still has a copy
- Scope for interesting OS/architecture hybrid
- System slows down if total memory requirement approaches RAM available, since space for replication is reduced
- ▶ Examples: DDM, KSR-1/2, rumours from IBM...

Which cache should the cache controller control?

- L1 cache is already very busy with CPU traffic
- L2 cache also very busy...
- L3 cache doesn't always have the current value for a cache line
 - 1. Although L1 cache is normally write-through, L2 is normally writeback
 - 2. Some data may bypass L3 (and perhaps L2) cache (eg when stream-prefetched)
- In Power4, cache controller manages L2 cache all external invalidations/requests
- L3 cache improves access to DRAM for accesses both from CPU and from network

Advanced Computer Architecture Chapter 6.77

Advanced Computer Architecture Chapter 6.79

Summary and Conclusions

- Caches are essential to gain the maximum performance from modern microprocessors
- The performance of a cache is close to that of SRAM but at the cost of DRAM
- Caches can be used to form the basis of a parallel computer
- Bus-based multiprocessors do not scale well: max < 10 nodes
- Larger-scale shared-memory multiprocessors require more complicated networks and protocols
- CC-NUMA is becoming popular since systems can be built from commodity components (chips, boards, OSs) and use existing software
- 🕨 e.g. HP/Convex, Sequent, Data General, SGI, Sun, IBM