

Outline

- Shared-centralized vs. Distributed Memory
- •Challenges of Parallel Processing
- Centralized shared-memory architecture
- Snoop based Coherence
- •Directory based Coherence





1. Centralized Memory MultiProcessor Also called Symmetric MultiProcessors (SMPs) because single main memory has a symmetric relationship to all processors Large caches ⇒ single memory can satisfy memory demands of small number of processors Can scale to a few dozen processors by using a switch and by using many memory banks Although scaling beyond that is technically conceivable, it becomes less attractive as the number of processors sharing centralized memory increases

5





2. Distributed Shared Memory (DSM) MultiProcessor

- Also called Non Uniform Memory Access (NUMA) since the access time depends on the location of a data word in memory
- Pros:
 - $^\circ~$ Cost-effective way to scale memory bandwidth
 - If most accesses are to local memory
 - $^\circ~$ Reduces latency of local memory accesses
- Cons:
 - ° Communicating data between processors more complex
 - ° Must change software to take advantage of increased memory BW







Outline

- •Shared-centralized vs. Distributed Memory
- •Challenges of Parallel Processing
- Centralized shared-memory architecture
- Snoop based Coherence
- •Directory based Coherence

Challenges of Parallel Processing	5
 First challenge: percentage of program that is inherently sequential 	
 Suppose we want to achieve 80X speedup fro 100 processors. What fraction of original prog can be sequential? a.10% b.5% c.1% d.<1% 	m gram
e.Impossible	13







Challenges of Parallel Processing

Possible solutions to the two challenges:

- Application parallelism ⇒ primarily via new algorithms that have better parallel performance
- 2. Long remote latency impact \Rightarrow both by architect and by the programmer
- For example, reduce freq. of remote accesses by:
 - $^{\circ}\,$ Caching shared data (HW), or
 - $^\circ\,$ Restructuring the data layout to make more accesses local (SW)

17

Outline

- •Shared-centralized vs. Distributed Memory
- •Challenges of Parallel Processing
- Centralized shared-memory architecture
- Snoop based Coherence
- •Directory based Coherence







Defining Coherent Memory System

- 1. <u>Preserve Program Order</u>: A read by processor P to location X that follows a write by P to X, with no writes of X by another processor occurring between the write and the read by P, always returns the value written by P
- Coherent view of memory: Read by a processor to location X that follows a write by another processor to X returns the written value if the read and write are sufficiently separated in time and no other writes to X occur between the two accesses
- 3. Write serialization: 2 writes to same location by any 2 processors are seen in the same order by all processors ° For example, if values 1 and then 2 are written to a location, processors can never read the value of the location as 2 and then later read it as 1

21



Migration and Replication
 In a coherent multiprocessor, caches provide both migration and replication
 Migration and replication: key to performance of shared data
 Migration - data can be moved to a local cache and used there in a transparent fashion
Reduces both latency to access shared data that is allocated remotely and bandwidth demand on the shared memory
 Replication - for shared data being simultaneously read, since caches make a copy of data in local cache Reduces both latency of access and contention for read shared data
- Reduces both latency of access and contention for read shared data 23

Enforcing Coherence Two Classes of Cache Coherence Protocols SMPs use a HW protocol to maintain coherent caches

- 1. Snooping Every cache with a copy of data also has a copy of *sharing status* of block, but no centralized state is kept
 - All caches are accessible via some broadcast medium (a bus or switch)
 - All cache controllers monitor or snoop on the medium to determine whether or not they have a copy of a block that is requested on a bus or switch access
- 2. Directory based Sharing status of a block of physical memory is kept in just one location, the directory













Cache Resources for WB Snooping

 To track whether a cache block is shared, add extra state bit ("shared") associated with each cache block, like "valid" bit and "dirty" bit

- $^\circ\,$ Write to Shared block \Rightarrow Need to place invalidate on bus and mark cache block as private
- $^{\circ}\,$ No further invalidations will be sent for that block
- $^\circ\,$ This processor called ${\color{black} {owner}}$ of cache block
- $^{\circ}\,$ Owner then changes state from shared to unshared (or exclusive)

•Recall:

- $^\circ\,$ Valid bit: says whether or not this cache entry contains a valid address
- Dirty bit: status bit that indicates whether the block in cache is dirty (modified while in the cache) or clean (not modified). If clean, the block is not written back on a miss, since identical information to the cache is found in lower levels.



Write-Back Snoopy Protocol	
 Invalidation protocol, write-back cache [°] Snoops every address on bus 	
•Each memory block is in one state: ^o Clean in all caches and up-to-date in main memory (Shared) ^o OR Dirty in exactly one cache (Modified) ^o OR Not in any caches	
 Each cache block is in one state: Shared: block can be read OR Modified: cache has only copy, it's writeable, and dirty OR Invalid: block contains no data 	33













			C,	d	nþ	JE						
	P1			P2			Bus				Men	nory
step	State	Addr	Value	State	Addr	Valu	Actior	Proc	Addr	Value	Add	Valu
P1 Write 10 to A	1 <u>Mod.</u>	<u>A1</u>	<u>10</u>				<u>WrMs</u>	P1	A1			
P1: Read A1	Mod.	A1	10									
P2: Read A1												
P2: Write 20 to	A1											
P2: Write 40 to	A2											
	Δς	suma	ς Δ1 o	nd A2	man t	o sam			~k			



				C			•					
	P1			P2			Bus				Men	nory
step	State	Addr	Value	State	Addr	Valu	Actior	Proc	Addr	Value	Add	Val
P1 Write 10 to A1	Mod.	A1	10				WrMs	P1	A1			
P1: Read A1	Mod.	A1	10									
P2: Read A1				Shar.	A1		RdMs	P2	A1			
	Shar.	A1	10				WrBk	P1	A1	10	A1	10
				Shar.	A1	10	RdDa	P2	A1	10	A1	10
P2: Write 20 to A	Inv.			Mod.	A1	20	WrMs	P2	A1		A1	10
P2: Write 40 to A	2											







Extensions to the Basic Snoopy Coherence Protocol

2. MESIF

- Intel i7 uses a variant of a MESI protocol, called MESIF, which adds a state (Forward) to designate which sharing processor should respond to a request
 - In MESI, a cache line request that is received by multiple caches may either be satisfied from (slow) main memory, or *all* the sharing caches could respond, bombarding the requestor with redundant responses

3. MOESI

- $^\circ\,$ Adds the state <code>Owned</code> to the MESI protocol to indicate that the associated block is owned by that cache and out-of-date in memory
- In MSI and MESI protocols, when there is an attempt to share a block in the M state, the state is changed to S (in both the original and newly sharing cache), and the block must be written back to memory
- $^\circ\,$ In MOESI, the block can be changed from the Modified to Owned state in the original cache without writing it to memory
- Other caches, which are newly sharing the block, keep the block in the S state; the O state, which only the original cache holds, indicates that main memory copy is out of date and that the designated cache is the owner.
- $^{\circ}\,$ owner of the block must supply it on a miss
- ° AMD Opteron uses the MOESI protocol



Outline

- •Shared-centralized vs. Distributed Memory
- •Challenges of Parallel Processing
- Centralized shared-memory architecture
- Snoop based Coherence
- •Directory based Coherence





• ...

Directory Protocol 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) Modified: 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 for every block, 1 if processor has copy) No bus and don't want to broadcast: Interconnect no longer single arbitration point! Terms: typically 3 processors involved Local node where a request originates Home node where the memory location and directory of an address resides Remote node has a copy of a cache block, whether modified or shared

Directory Protocol Messages

Message type	Source	Destination	Message contents	Function of this message
Read miss	Local cache	Home directory	Р, А	Node P has a read miss at address A; request data and make P a read sharer.
Write miss	Local cache	Home directory	P, A	Node P has a write miss at address A; request data and make P the exclusive owner.
Invalidate	Local cache	Home directory	А	Request to send invalidates to all remote caches that are caching the block at address A.
Invalidate	Home directory	Remote cache	А	Invalidate a shared copy of data at address A.
Fetch	Home directory	Remote cache	А	Fetch the block at address A and send it to its home directory; change the state of A in the remote cache to shared.
Fetch/invalidate	Home directory	Remote cache	А	Fetch the block at address A and send it to its home directory; invalidate the block in the cache.
Data value reply	Home directory	Local cache	D	Return a data value from the home memory.
Data write-back	Remote cache	Home directory	A, D	Write-back a data value for address A.

State Transition Diagram for One Cache Block in Directory Based System

- States identical to snoopy case; transactions very similar
- •Transitions caused by read misses, write misses, invalidates, data fetch requests
- •Generates read miss & write miss messages to home directory
 - Over the observation of the bus for shooping => explicit invalidate & data fetch requests



State Transition Diagram for Directory

- •Same states & structure as the transition diagram for an individual cache
- •2 actions: update of directory state & send messages to satisfy requests
- •Tracks all copies of memory block
- •Also updates the sharing set, Sharers, when necessary, and sends back the *value* or *invalidate* messages 53







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. 	
^o 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 identity of new owner, and state of block is made Exclusive. 56	
6	



















Feature	AMD Opteron 8439	IBM Power 7	Intel Xenon 7560	Sun T2
Transistors	904 M	1200 M	2300 M	500 M
Power (nominal)	137 W	140 W	130 W	95 W
Max. cores/chip	6	8	8	8
Multithreading	No	SMT	SMT	Fine-grained
Threads/core	1	4	2	8
Instruction issue/clock	3 from one thread	6 from one thread	4 from one thread	2 from 2 threads
Clock rate	2.8 GHz	4.1 GHz	2.7 GHz	1.6 GHz
Outermost cache	L3; 6 MB; shared	L3; 32 MB (using embedded DRAM); shared or private/core	L3; 24 MB; shared	L2; 4 MB; shared
Inclusion	No, although L2 is superset of L1	Yes, L3 superset	Yes, L3 superset	Yes
Multicore coherence protocol	MOESI	Extended MESI with behavioral and locality hints (13-state protocol)	MESIF	MOESI
Multicore coherence implementation	Snooping	Directory at L3	Directory at L3	Directory at L2
Extended coherence upport	Up to 8 processor chips can be connected with HyperTransport in a ring, using directory or snooping. System is NUMA.	Up to 32 processor chips can be connected with the SMP links. Dynamic distributed directory structure. Memory access is symmetric outside of an	Up to 8 processor cores can be implemented via Quickpath Interconnect. Support for directories with external logic.	Implemented via four coherence links per processor that can be used to snoop. Up to two chips directly connect, and up to four connect using

