

**Lecture 18:  
Multiprocessors 2:  
Snooping v. Directory Coherency,  
Memory Consistency Models**

**Professor David A. Patterson  
Computer Science 252  
Spring 1998**

DAP Spr:'98 ©UCB 1

## Review: Parallel Framework



- **Layers:**

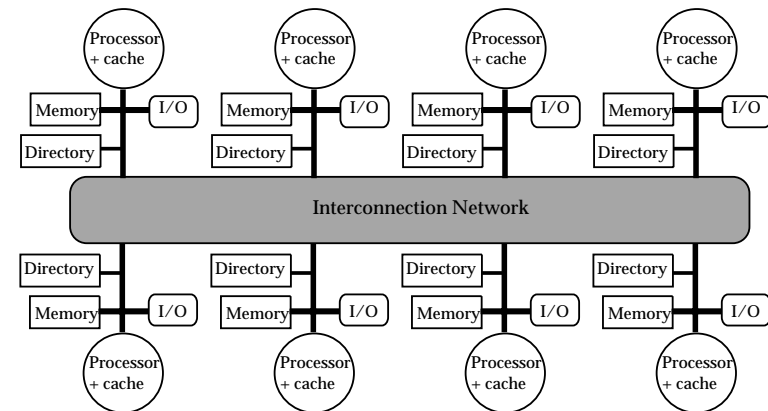
- **Programming Model:**

- » **Multiprogramming** : lots of jobs, no communication
- » **Shared address space**: communicate via memory
- » **Message passing**: send and receive messages
- » **Data Parallel**: several agents operate on several data sets simultaneously and then exchange information globally and simultaneously (shared or message passing)

- **Communication Abstraction:**

- » **Shared address space**: e.g., load, store, atomic swap
- » **Message passing**: e.g., send, receive library calls
- » **Debate over this topic (ease of programming, scaling)**  
=> many hardware designs 1:1 programming model

## Distributed Directory MPs



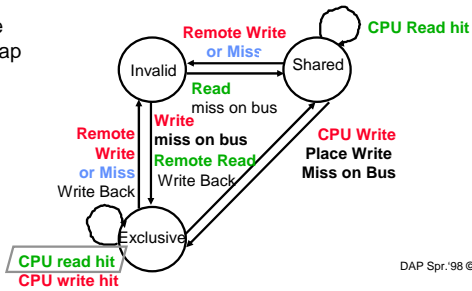
DAP Spr.'98 ©UCB 4



## Example: Step 2

step	P1			P2			Bus			Memory		
	State	Addr	Value	State	Addr	Value	Action	Proc.	Addr	Value	Addr	Value
P1: Write 10 to A1	Excl.	A1	10				WrMs	P1	A1			
P1: Read A1	Excl.	A1	10									
P2: Read A1												
P2: Write 20 to A1												
P2: Write 40 to A2												

Assumes initial cache state is invalid and A1 and A2 map to same cache block, but A1 ≠ A2.

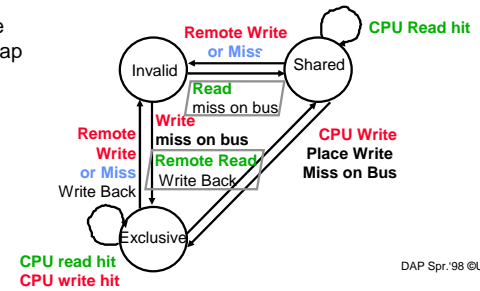


DAP Spr.'98 ©UCB 7

## Example: Step 3

step	P1			P2			Bus			Memory		
	State	Addr	Value	State	Addr	Value	Action	Proc.	Addr	Value	Addr	Value
P1: Write 10 to A1	Excl.	A1	10				WrMs	P1	A1			
P1: Read A1	Excl.	A1	10									
P2: Read A1				Shar.	A1		RdMs	P2	A1			
				Shar.	A1	10	WrBk	P1	A1	10	A1	10
P2: Write 20 to A1							RdDa	P2	A1	10		10
P2: Write 40 to A2												10

Assumes initial cache state is invalid and A1 and A2 map to same cache block, but A1 ≠ A2.



DAP Spr.'98 ©UCB 8



## Snooping Coherency 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** DAP Spr:98 ©UCB 11

## Implementing Snooping Caches

- **Multiple processors must be on bus, access to both addresses and data**
- **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** to allow checks in parallel with CPU
  - solution 2: L2 cache already duplicate and underutilized, **provided L2 obeys inclusion** with L1 cache
    - » block size, associativity of L2 affects L1

DAP Spr:98 ©UCB 12

## 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 its dirty, so must reply
- Add extra state bit to cache to determine shared or not
- Add 4th state (MESI)

DAP Spr:98 ©UCB 13

## Larger MPs

- Separate Memory per Processor
- Local or Remote access via memory controller
- 1 Cache Coherency solution: non-cached pages
- Alternative: **directory** 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  $f(\text{memory size})$  vs.  $f(\text{cache size})$
- Prevent directory as bottleneck?  
distribute directory entries with memory, each keeping track of which Procs have copies of their blocks

DAP Spr:98 ©UCB 14

## 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)
  - **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

DAP Spr:98 ©UCB 15

## Directory Protocol

- 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

DAP Spr:98 ©UCB 16



## Directory Protocol Messages

Message type	Source	Destination	Msg Content
<b>Read miss</b>	Local cache	Home directory	P, A
– Processor P reads data at address A; make P a read sharer and arrange to send data back			
<b>Write miss</b>	Local cache	Home directory	P, A
– Processor P writes data at address A; make P the exclusive owner and arrange to send data back			
<b>Invalidate</b>	Home directory	Remote caches	A
– Invalidate a shared copy at address A.			
<b>Fetch</b>	Home directory	Remote cache	A
– Fetch the block at address A and send it to its home directory			
<b>Fetch/Invalidate</b>	Home directory	Remote cache	A
– Fetch the block at address A and send it to its home directory; invalidate the block in the cache			
<b>Data value reply</b>	Home directory	Local cache	Data
– Return a data value from the home memory (read miss response)			
<b>Data write-back</b>	Remote cache	Home directory	A, Data
– Write-back a data value for address A (invalidate response)			

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

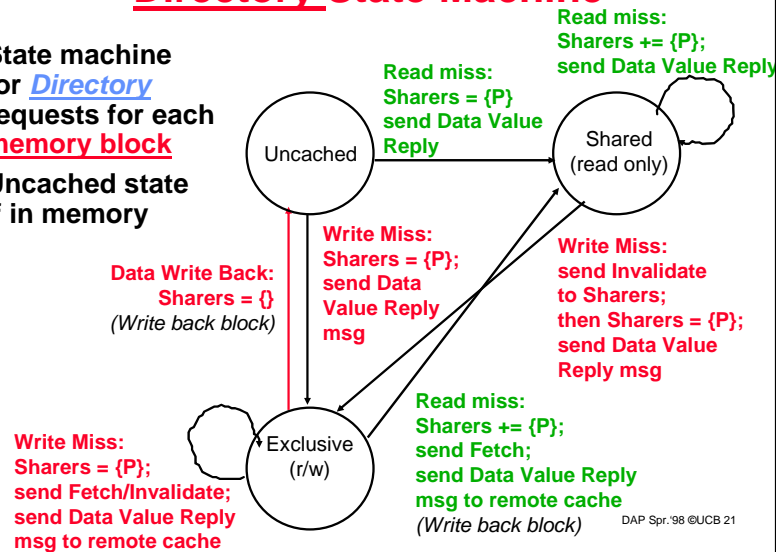
- States identical to snoopy case; transactions very similar.
- Transitions caused by read misses, write misses, invalidates, data fetch requests
- Generates read miss & write miss msg to home directory.
- Write misses that were broadcast on the bus for snooping => explicit invalidate & data fetch requests.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

DAP Spr. '98 ©UCB 18



## Directory State Machine

- State machine for Directory requests for each memory block
- Uncached state if in memory



## 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 only 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.

DAP Spr: '98 ©UCB 22

## 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 identity of new owner, and state of block is made Exclusive.

DAP Spr. '98 ©UCB'23

## Example

Processor 1 Processor 2 Interconnect Directory Memory

step	P1			P2			Bus			Directory			Memory
	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													
P2: Write 20 to A1													
P2: Write 40 to A2													

A1 and A2 map to the same cache block

DAP Spr. '98 ©UCB 24

# Example

	Processor 1			Processor 2			Interconnect			Directory			Memory	
step	P1 State	P1 Addr	P1 Value	P2 State	P2 Addr	P2 Value	Bus Action	Proc.	Addr	Value	Addr	State	{Procs}	Value
P1: Write 10 to A1	Excl.	A1	10				WrMs	P1	A1		A1	Ex	{P1}	
P1: Read A1							DaRp	P1	A1	0				
P2: Read A1														
P2: Write 20 to A1														
P2: Write 40 to A2														

A1 and A2 map to the same cache block

DAP Spr:98 ©UCB 25

# Example

	Processor 1			Processor 2			Interconnect			Directory			Memory	
step	P1 State	P1 Addr	P1 Value	P2 State	P2 Addr	P2 Value	Bus Action	Proc.	Addr	Value	Addr	State	{Procs}	Value
P1: Write 10 to A1	Excl.	A1	10				WrMs	P1	A1		A1	Ex	{P1}	
P1: Read A1							DaRp	P1	A1	0				
P2: Read A1														
P2: Write 20 to A1														
P2: Write 40 to A2														

A1 and A2 map to the same cache block

DAP Spr:98 ©UCB 26

## Example

Processor 1 Processor 2 Interconnect Directory Memory

step	P1			P2			Bus			Directory			Memory	
	State	Addr	Value	State	Addr	Value	Action	Proc.	Addr	Value	Addr	State		{Procs}
P1: Write 10 to A1							WrMs	P1	A1		A1	Ex	{P1}	
P1: Read A1	Excl.	A1	10				DaRp	P1	A1	0				
P2: Read A1				Shar.	A1		RdMs	P2	A1					
	Shar.	A1	10				Fich	P1	A1	10			A1	10
				Shar.	A1	10	DaRp	P2	A1	10	A1	Shar.	{P1,P2}	10
P2: Write 20 to A1														10
P2: Write 40 to A2														10

Write Back

A1 and A2 map to the same cache block

DAP Spr.'98 ©UCB 27

## Example

Processor 1 Processor 2 Interconnect Directory Memory

step	P1			P2			Bus			Directory			Memory	
	State	Addr	Value	State	Addr	Value	Action	Proc.	Addr	Value	Addr	State		{Procs}
P1: Write 10 to A1							WrMs	P1	A1		A1	Ex	{P1}	
P1: Read A1	Excl.	A1	10				DaRp	P1	A1	0				
P2: Read A1				Shar.	A1		RdMs	P2	A1					
	Shar.	A1	10				Fich	P1	A1	10			A1	10
				Shar.	A1	10	DaRp	P2	A1	10	A1	Shar.	{P1,P2}	10
P2: Write 20 to A1														10
				Excl.	A1	20	WrMs	P2	A1					10
P2: Write 40 to A2							Inval.	P1	A1		A1	Excl.	{P2}	10

A1 and A2 map to the same cache block

DAP Spr.'98 ©UCB 28

## Example

Processor 1 Processor 2 Interconnect Directory Memory

step	P1			P2			Bus			Directory			Memor	
	State	Addr	Value	State	Addr	Value	Action	Proc.	Addr	Value	Addr	State	{Procs}	Value
P1: Write 10 to A1							WrMs	P1	A1		A1	Ex	{P1}	
P1: Read A1	Excl.	A1	10				DaRp	P1	A1	0				
P2: Read A1				Shar.	A1		RdMs	P2	A1					
	Shar.	A1	10				Fich	P1	A1	10			A1	10
				Shar.	A1	10	DaRp	P2	A1	10	A1	Shar.	{P1,P2}	10
P2: Write 20 to A1				Excl.	A1	20	WrMs	P2	A1					10
	Inv.						Inval.	P1	A1		A1	Excl.	{P2}	10
P2: Write 40 to A2							WrMs	P2	A2		A2	Excl.	{P2}	0
							WrBk	P2	A1	20	A1	Unca.	∅	20
				Excl.	A2	40	DaRp	P2	A2	0	A2	Excl.	{P2}	0

A1 and A2 map to the same cache block

DAP Spr. '98 ©UCB 29

## Implementing a Directory

- We assume operations atomic, but they are not; reality is much harder; must avoid deadlock when run out of buffers in network (see Appendix E)
- Optimizations:
  - read miss or write miss in Exclusive: send data directly to requestor from owner vs. 1st to memory and then from memory to requestor

DAP Spr. '98 ©UCB 30

## Synchronization

- **Why Synchronize?** Need to know when it is safe for different processes to use shared data
- **Issues for Synchronization:**
  - Uninterruptible 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

DAP Spr: '98 ©UCB 31

## Uninterruptible 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)
    - 1 if other processor had already claimed access
  - Key is that exchange operation is indivisible
- **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

DAP Spr: '98 ©UCB 32



## Uninterruptible Instruction to Fetch and Update Memory

- Hard 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 preceding load) and 0 otherwise

- **Example doing atomic swap with LL & SC:**

```
try:  mov    R3,R4      ; mov exchange value
      ll     R2,0(R1)   ; load linked
      sc     R3,0(R1)   ; store conditional
      beqz   R3,try     ; branch store fails (R3 = 0)
      mov    R4,R2     ; put load value in R4
```

- **Example doing fetch & increment with LL & SC:**

```
try:  ll     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)
```

DAP Spr: '98 @UCB 33

## User Level Synchronization— Operation Using this Primitive

- **Spin locks:** processor continuously tries to acquire, spinning around a loop trying to get the lock

```
lockit:  li     R2,#1
         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     R3,0(R1)   ;load var
         bnez  R3,lockit  ;not free=>spin
         exch  R2,0(R1)   ;atomic exchange
         bnez  R2,try     ;already locked?
```

DAP Spr: '98 @UCB 34

## 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;
.....	.....
A = 1;	B = 1;
L1: if (B == 0) ...	L2: 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

DAP Spr. '98 ©UCB 35

## Memory Consistency Model

- Schemes faster execution to sequential consistency
- 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

DAP Spr. '98 ©UCB 36

## Review

- 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

DAP Spr.'98 ©UCB 37