
EECS 252 Graduate Computer Architecture

Lec 13 – Snooping Cache and Directory Based Multiprocessors

**David Patterson
Electrical Engineering and Computer Sciences
University of California, Berkeley**

**<http://www.eecs.berkeley.edu/~pattarn>
<http://vlsi.cs.berkeley.edu/cs252-s06>**

Outline

- **Review**
- **Coherence**
- **Write Consistency**
- **Administrivia**
- **Snooping**
- **Building Blocks**
- **Snooping protocols and examples**
- **Coherence traffic and Performance on MP**
- **Directory-based protocols and examples (if get this far)**
- **Conclusion**

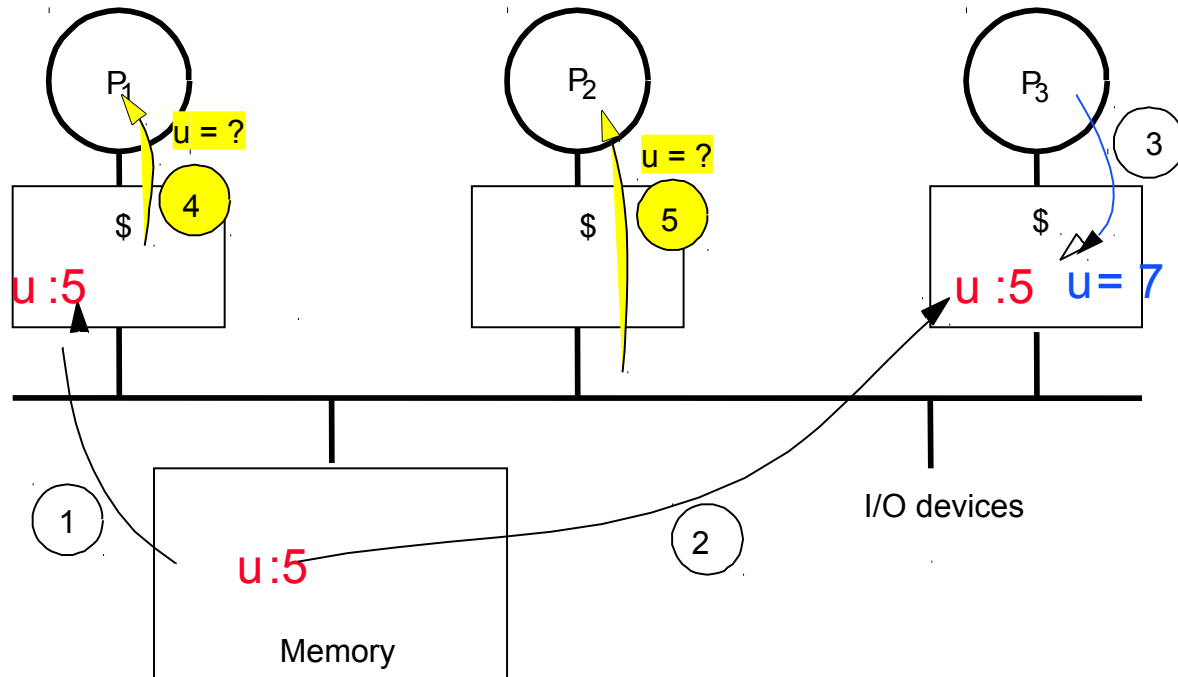
Challenges of Parallel Processing

- 1. Application parallelism \Rightarrow 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 frequency of remote accesses either by**
 - Caching shared data (HW)**
 - Restructuring the data layout to make more accesses local (SW)**
 - Today's lecture on HW to help latency via caches**

Symmetric Shared-Memory Architectures

- From multiple boards on a shared bus to multiple processors inside a single chip
- Caches both
 - Private data are used by a single processor
 - Shared data are used by multiple processors
- Caching shared data
 - ⇒ reduces latency to shared data, memory bandwidth for shared data, and interconnect bandwidth
 - ⇒ cache coherence problem

Example Cache Coherence Problem

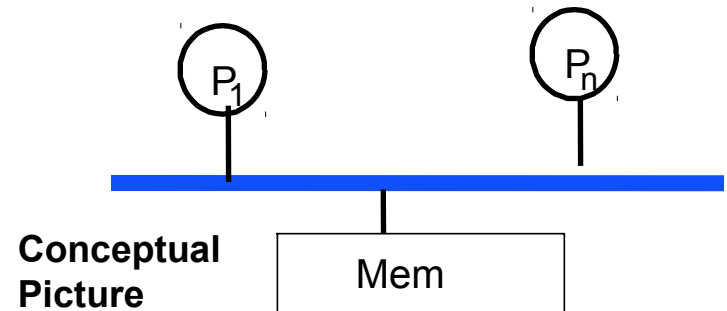


- Processors see different values for **u** after event 3
- With write back caches, value written back to memory depends on happenstance of which cache flushes or writes back value when
 - » Processes accessing main memory may see very stale value
- Unacceptable for programming, and its frequent!

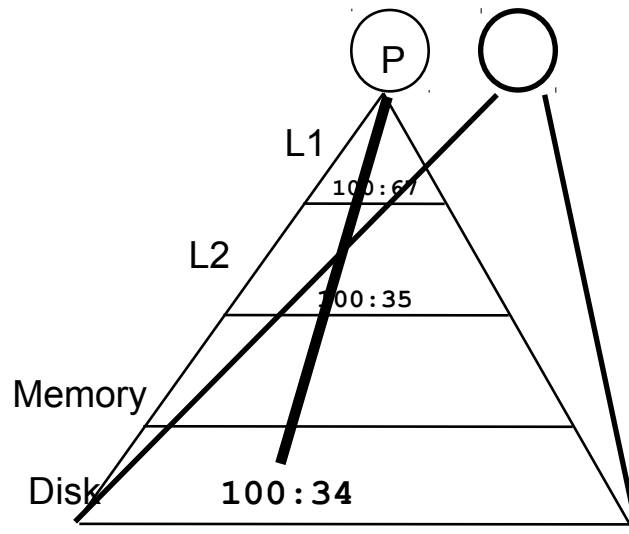
Example

P_1	P_2
/*Assume initial value of A and flag is 0*/	
A = 1;	while (flag == 0); /*spin idly*/
flag = 1;	print A;

- **Intuition not guaranteed by coherence**
- **expect memory to respect order between accesses to *different* locations issued by a given process**
 - to preserve orders among accesses to same location by different processes
- **Coherence is not enough!**
 - pertains only to single location



Intuitive Memory Model



- Reading an address should **return the last value written** to that address
 - Easy in uniprocessors, except for I/O
- **Too vague and simplistic; 2 issues**
 1. **Coherence** defines **values** returned by a read
 2. **Consistency** determines **when** a written value will be returned by a read
- **Coherence** defines behavior to same location, **Consistency** defines behavior to other locations

Defining Coherent Memory System

1. **Preserve Program Order**: 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
2. **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
 - If not, a processor could keep value 1 since saw as last write
 - For example, if the 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

Write Consistency

- For now assume
 1. A write does not complete (and allow the next write to occur) until all processors have seen the effect of that write
 2. The processor does not change the order of any write with respect to any other memory access
- ⇒ if a processor writes location A followed by location B, any processor that sees the new value of B must also see the new value of A
- These restrictions allow the processor to reorder reads, but forces the processor to finish writes in program order

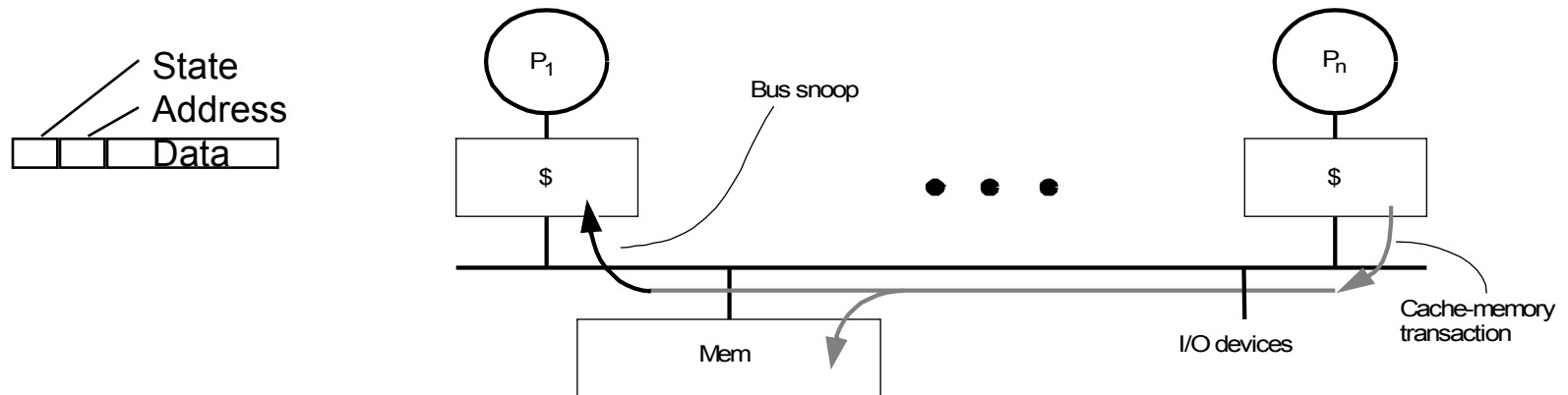
Basic Schemes for Enforcing Coherence

- **Program on multiple processors will normally have copies of the same data in several caches**
 - Unlike I/O, where its rare
- **Rather than trying to avoid sharing in SW, SMPs use a HW protocol to maintain coherent caches**
 - 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 reading shared data simultaneously, since caches make a copy of data in local cache**
 - Reduces both latency of access and contention for read shared data

2 Classes of Cache Coherence Protocols

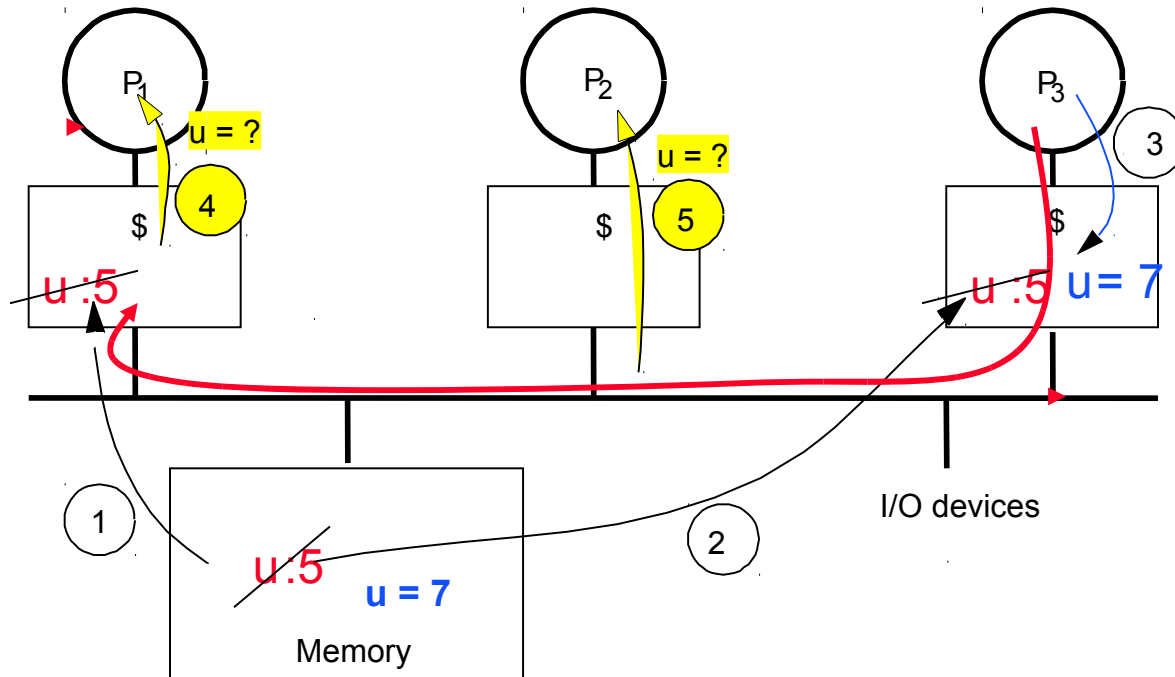
1. **Directory based** — Sharing status of a block of physical memory is kept in just one location, the **directory**
2. **Snooping** — Every cache with a copy of data also has a copy of sharing status of block, but no centralized state is kept
 - All caches are accessible via some broadcast medium (a 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

Snoopy Cache-Coherence Protocols



- Cache Controller “**snoops**” all transactions on the shared medium (bus or switch)
 - relevant transaction if for a block it contains
 - take action to ensure coherence
 - » invalidate, update, or supply value
 - depends on state of the block and the protocol
- Either get exclusive access before write via write invalidate or update all copies on write

Example: Write-thru Invalidate



- **Must invalidate before step 3**
- **Write update uses more broadcast medium BW**
⇒ all recent MPUs use write invalidate

Architectural Building Blocks

- **Cache block state transition diagram**
 - FSM specifying how disposition of block changes
 - » invalid, valid, exclusive
- **Broadcast Medium Transactions (e.g., bus)**
 - Fundamental system design abstraction
 - Logically single set of wires connect several devices
 - Protocol: arbitration, command/addr, data
 - ⇒ Every device observes every transaction
- **Broadcast medium enforces serialization of read or write accesses ⇒ Write serialization**
 - 1st processor to get medium invalidates others copies
 - Implies cannot complete write until it obtains bus
 - All coherence schemes require serializing accesses to same cache block
- **Also need to find up-to-date copy of cache block**

Locate up-to-date copy of data

- **Write-through: get up-to-date copy from memory**
 - Write through simpler if enough memory BW
- **Write-back harder**
 - Most recent copy can be in a cache
- **Can use same snooping mechanism**
 1. **Snoop every address placed on the bus**
 2. **If a processor has dirty copy of requested cache block, it provides it in response to a read request and aborts the memory access**
 - Complexity from retrieving cache block from cache, which can take longer than retrieving it from memory
- **Write-back needs lower memory bandwidth**
 - ⇒ **Support larger numbers of faster processors**
 - ⇒ **Most multiprocessors use write-back**

Cache Resources for WB Snooping

- To track whether a cache block is shared, add extra state bit associated with each cache block, like valid bit and dirty bit
 - Write to Shared block \Rightarrow Need to place invalidate on bus and mark cache block as private (if an option)
 - No further invalidations will be sent for that block
 - This processor called owner of cache block
 - Owner then changes state from shared to unshared (or exclusive)

Cache behavior in response to bus

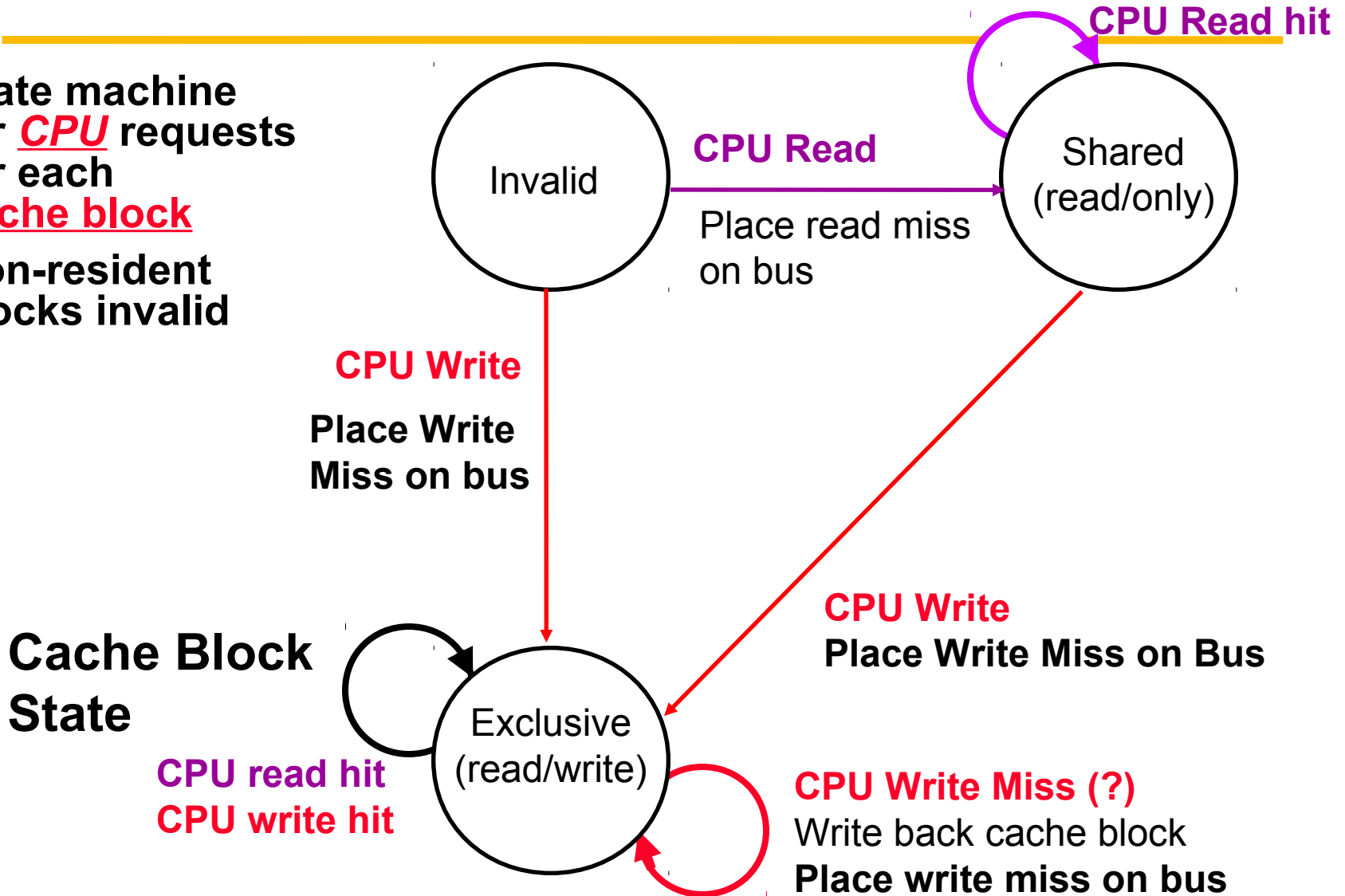
- **Every bus transaction must check the cache-address tags**
 - could potentially interfere with processor cache accesses
 - **A way to reduce interference is to duplicate tags**
 - One set for caches access, one set for bus accesses
 - **Another way to reduce interference is to use L2 tags**
 - Since L2 less heavily used than L1
- ⇒ **Every entry in L1 cache must be present in the L2 cache, called the inclusion property**
- If Snoop gets a hit in L2 cache, then it must arbitrate for the L1 cache to update the state and possibly retrieve the data, which usually requires a stall of the processor

Example Write Back Snoopy Protocol

- Invalidation protocol, write-back cache
 - Snoops every address on bus
 - If it has a dirty copy of requested block, provides that block in response to the read request and aborts the memory access
- Each **memory** block is in one state:
 - Clean in all caches and up-to-date in memory (**Shared**)
 - OR Dirty in exactly one cache (**Exclusive**)
 - OR Not in any caches
- Each **cache** block is in one state (track these):
 - **Shared** : block can be read
 - OR **Exclusive** : cache has only copy, its writeable, and dirty
 - OR **Invalid** : block contains no data (in uniprocessor cache too)
- Read misses: cause all caches to snoop bus
- Writes to clean blocks are treated as misses

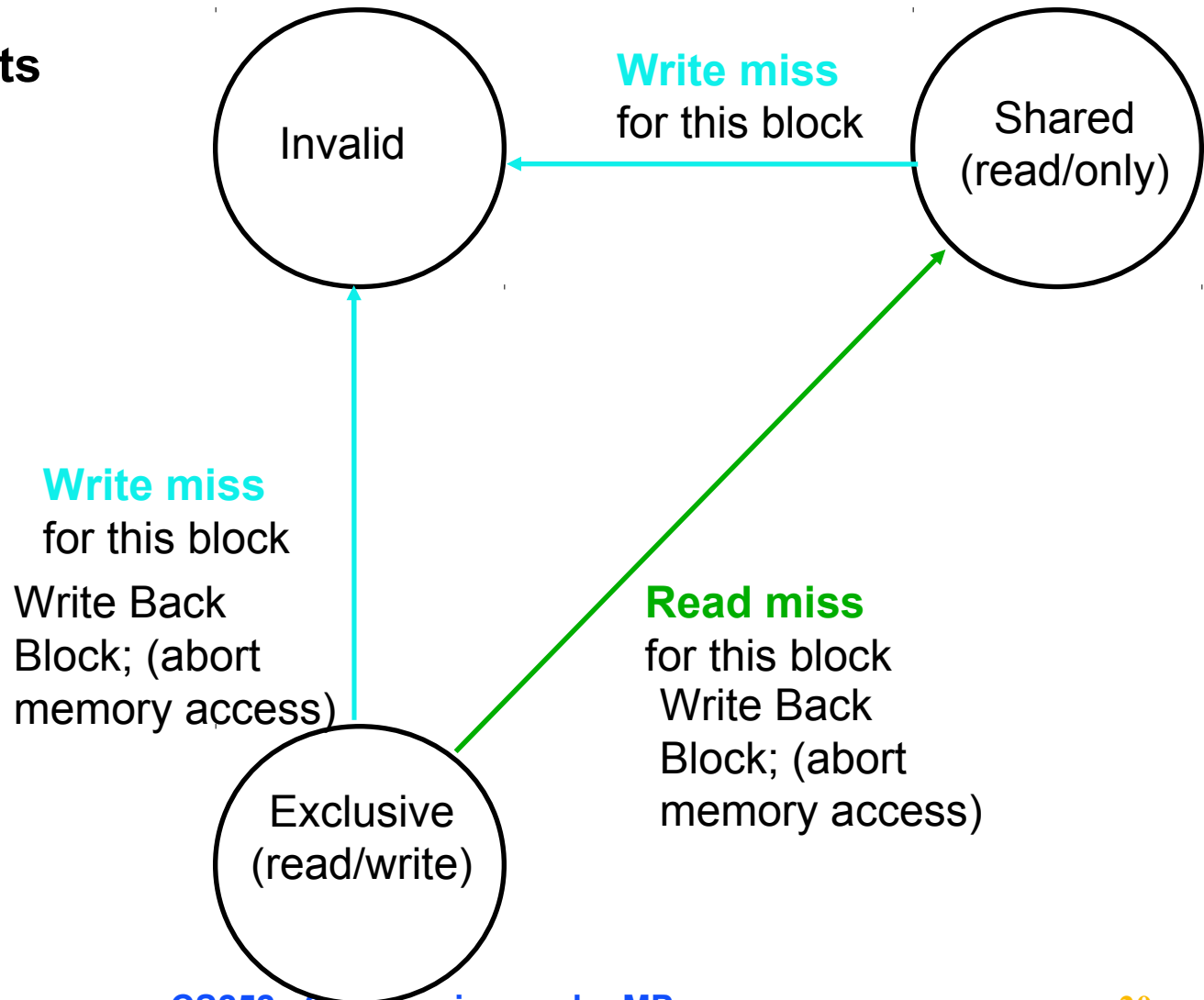
Write-Back State Machine - CPU

- State machine for **CPU** requests for each **cache block**
- Non-resident blocks invalid



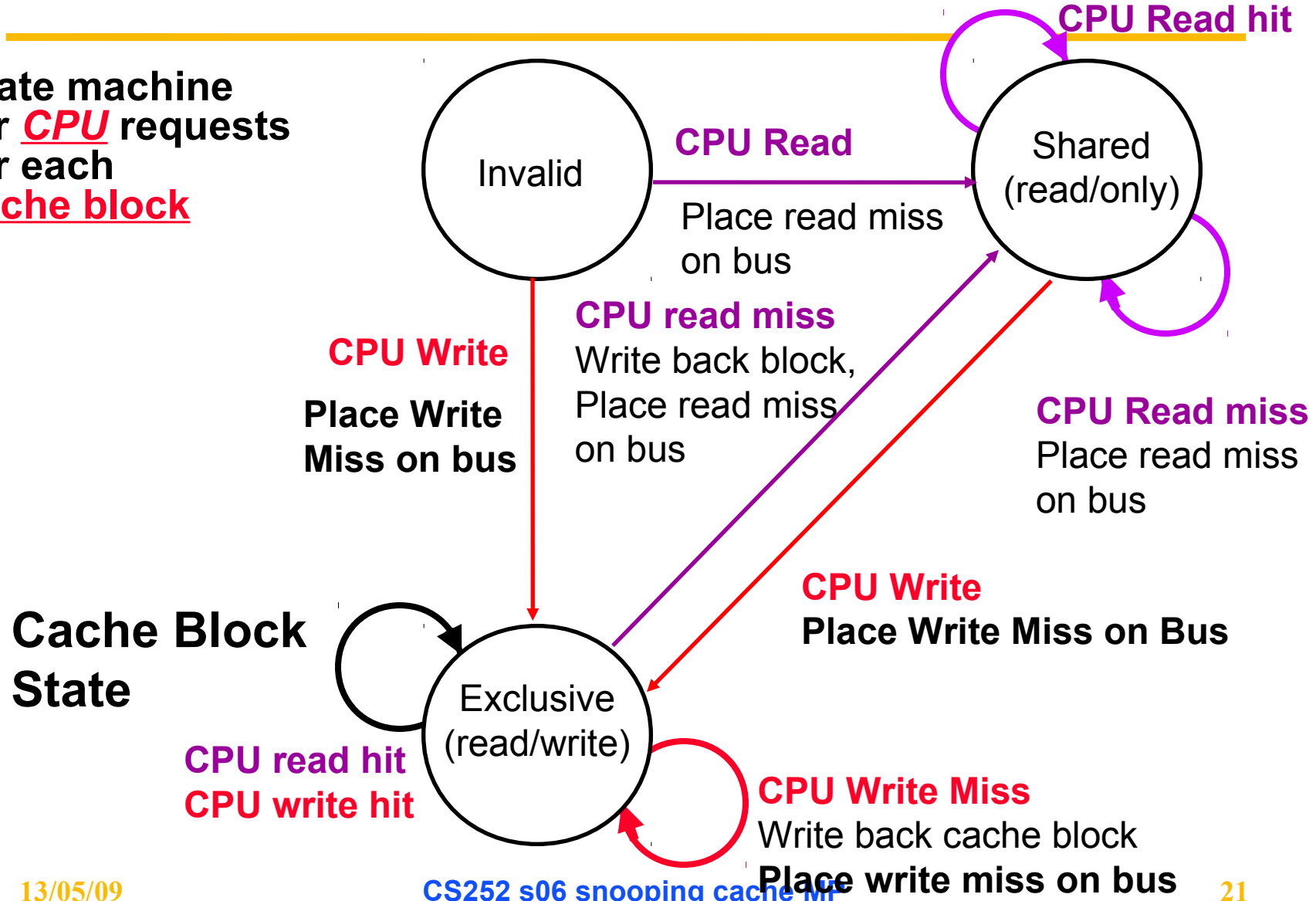
Write-Back State Machine- Bus request

- State machine for **bus** requests for each **cache block**



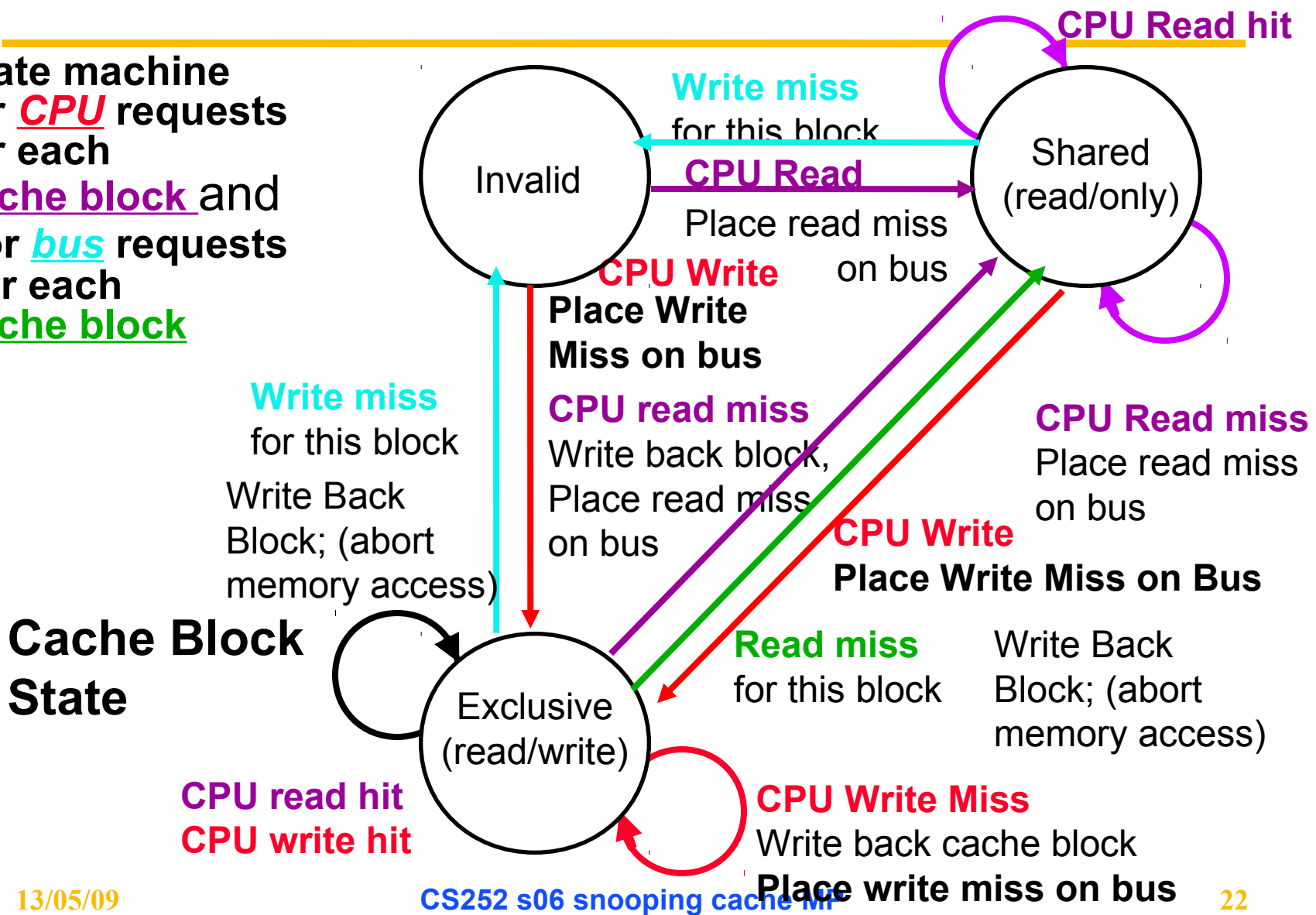
Block-replacement

- State machine for **CPU** requests for each **cache block**



Write-back State Machine-III

- State machine for **CPU** requests for each **cache block** and for **bus** requests for each **cache block**



Example

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

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

Example

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

Assumes A1 and A2 map to same cache block

Example

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

Assumes A1 and A2 map to same cache block

Example

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

Assumes A1 and A2 map to same cache block

Example

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

Assumes A1 and A2 map to same cache block

Example

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

Assumes A1 and A2 map to same cache block,
but A1 != A2

Limitations in Symmetric Shared-Memory Multiprocessors and Snooping Protocols

- **Single memory accommodate all CPUs**
⇒ **Multiple memory banks**
- **Bus-based multiprocessor, bus must support both coherence traffic & normal memory traffic**
⇒ **Multiple buses or interconnection networks (cross bar or small point-to-point)**
- **Opteron**
 - **Memory connected directly to each dual-core chip**
 - **Point-to-point connections for up to 4 chips**
 - **Remote memory and local memory latency are similar, allowing OS Opteron as UMA computer**

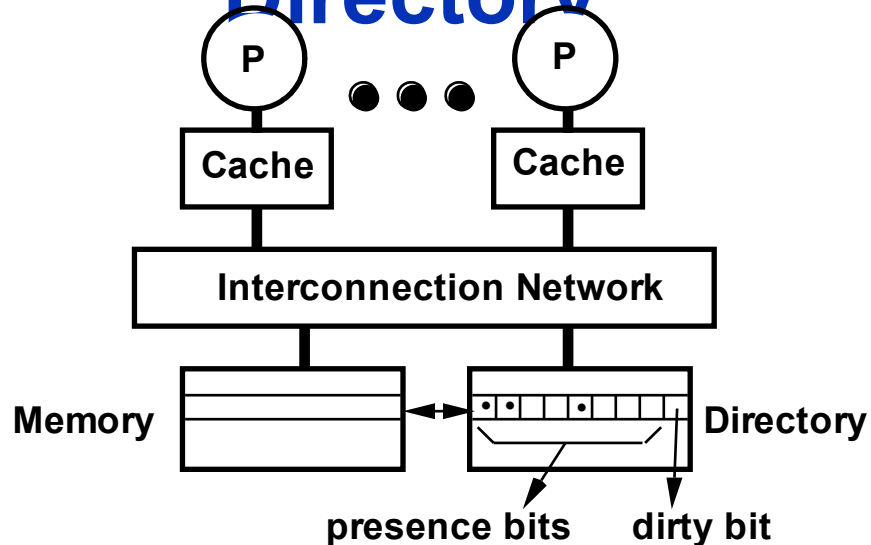
Bus-based Coherence

- **All of (a), (b), (c) done through broadcast on bus**
 - faulting processor sends out a “search”
 - others respond to the search probe and take necessary action
- **Could do it in scalable network too**
 - broadcast to all processors, and let them respond
- **Conceptually simple, but broadcast doesn't scale with p**
 - on bus, bus bandwidth doesn't scale
 - on scalable network, every fault leads to at least p network transactions
- **Scalable coherence:**
 - can have same cache states and state transition diagram
 - different mechanisms to manage protocol

Scalable Approach: Directories

- **Every memory block has associated directory information**
 - keeps track of copies of cached blocks and their states
 - on a miss, find directory entry, look it up, and communicate only with the nodes that have copies if necessary
 - in scalable networks, communication with directory and copies is through network transactions
- **Many alternatives for organizing directory information**

Basic Operation of Directory



- k processors.
- With each cache-block in memory: k presence-bits, 1 dirty-bit
- With each cache-block in cache: 1 valid bit, and 1 dirty (owner) bit

- Read from main memory by processor i :
- If dirty-bit OFF then { read from main memory; turn $p[i]$ ON; }
- if dirty-bit ON then { recall line from dirty proc (cache state to shared); update memory; turn dirty-bit OFF; turn $p[i]$ ON; supply recalled data to i ;}
 - Write to main memory by processor i :

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

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

Directory Protocol Messages (Fig 4.22)

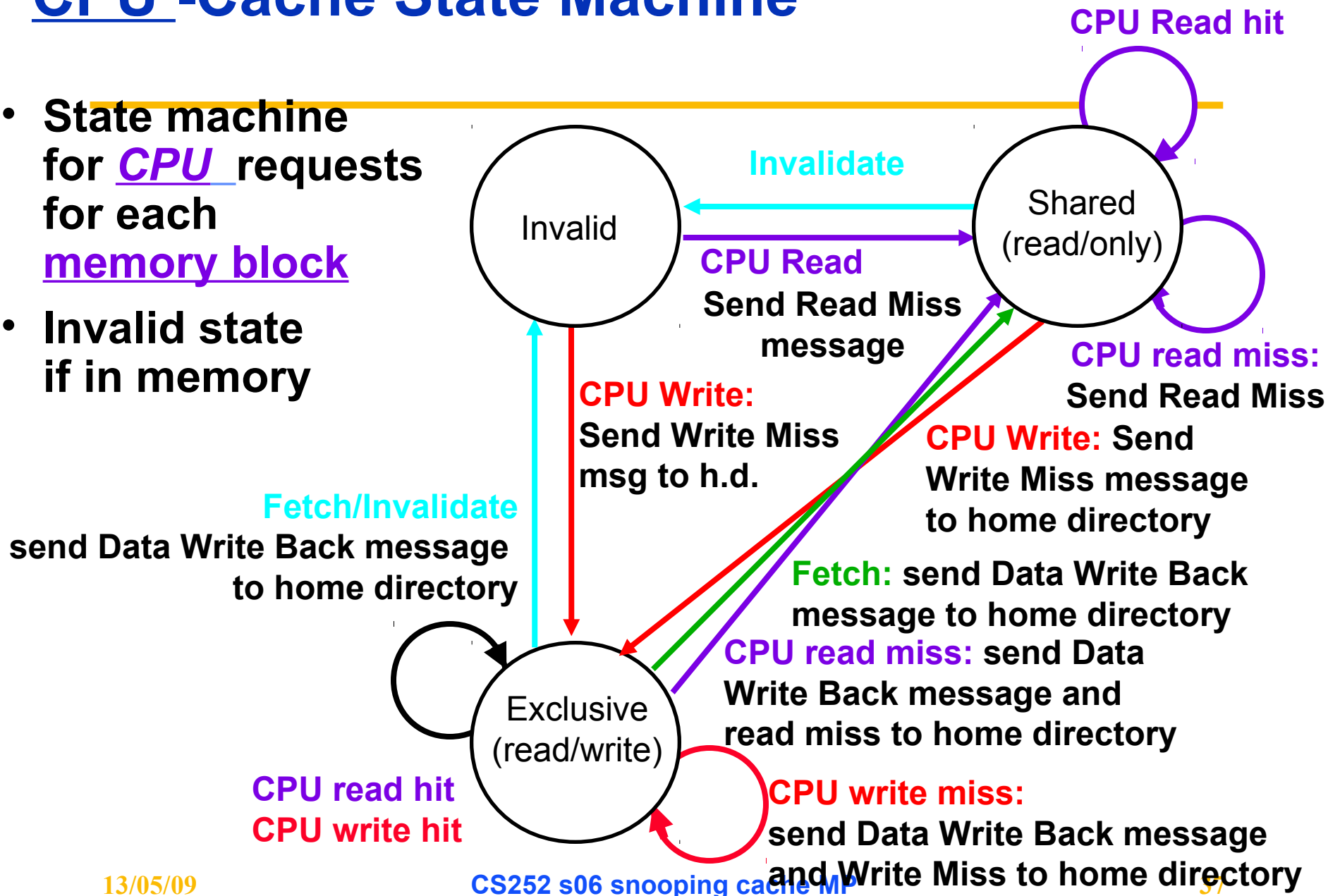
<i>Message type</i>	<i>Source</i>	<i>Destination</i>	<i>Msg Content</i>
Read miss	Local cache	Home directory	P, A
<i>– Processor P reads data at address A; make P a read sharer and request data</i>			
Write miss	Local cache	Home directory	P, A
<i>– Processor P has a write miss at address A; make P the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

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 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**

CPU -Cache State Machine

- State machine for CPU requests for each memory block
- Invalid state if in memory

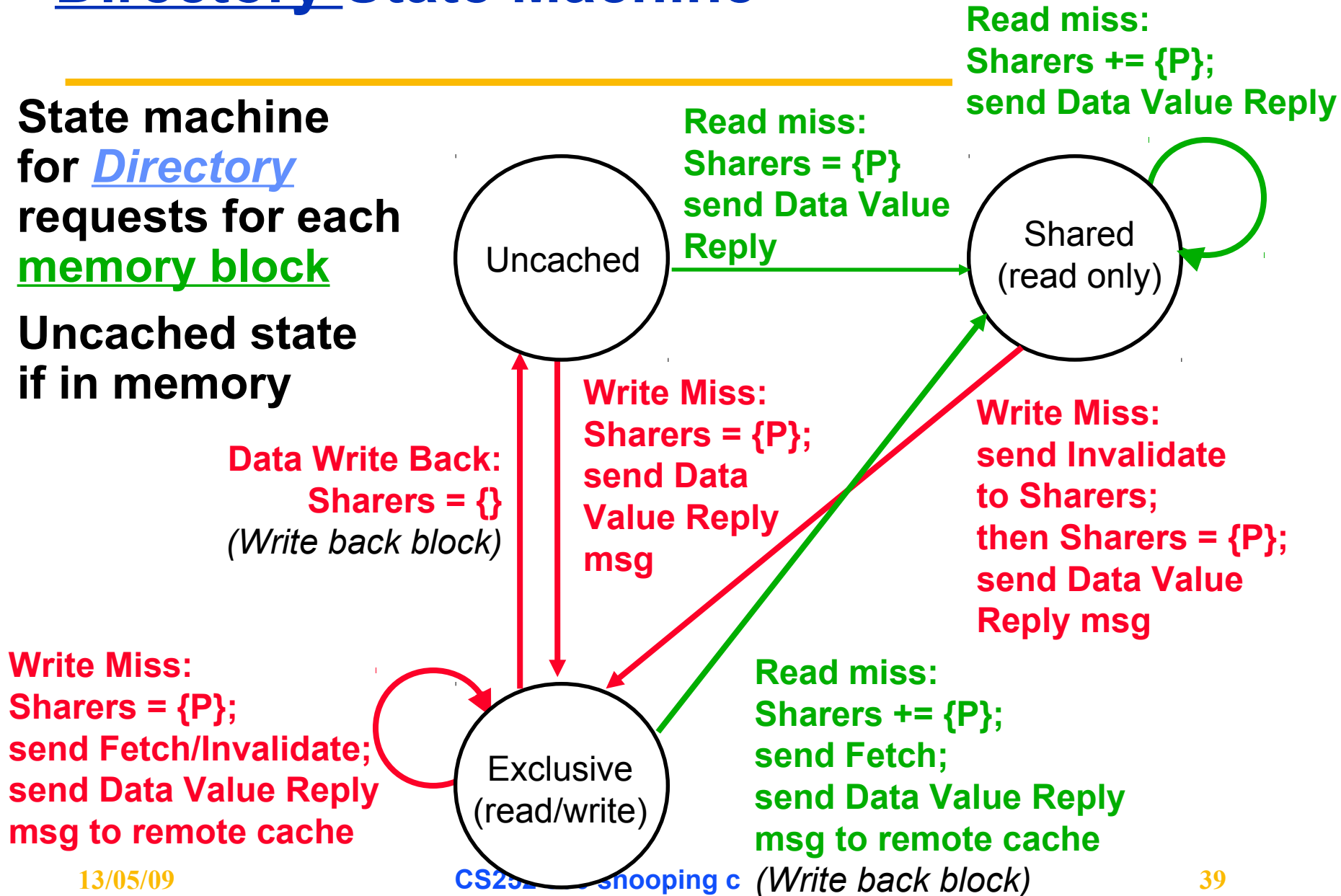


State Transition Diagram for Directory

- 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 indicates an action that updates the sharing set, **Sharers**, as well as sending a message

Directory State Machine

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



And in Conclusion ...

- **Caches contain all information on state of cached memory blocks**
- **Snooping cache over shared medium for smaller MP by invalidating other cached copies on write**
- **Sharing cached data \Rightarrow Coherence (values returned by a read), Consistency (when a written value will be returned by a read)**
- **Snooping and Directory Protocols similar; bus makes snooping easier because of broadcast (snooping \Rightarrow uniform memory access)**
- **Directory has extra data structure to keep track of state of all cache blocks**
- **Distributing directory \Rightarrow scalable shared address multiprocessor
 \Rightarrow Cache coherent, Non uniform memory access**