

# **EEL 6764 Principles of Computer Architecture**

# **Multiprocessors and**

# **Thread-Level Parallelism**

Dr Hao Zheng  
Dept of Comp Sci & Eng  
U of South Florida

# What are Multiprocessors?

- Tightly coupled processors
  - Controlled by a single OS
  - With shared memory space
  - Communication done in HW
- Clusters = processors connected by network
  - Comm. among different processors coordinated by OSs
- Support MIMD execution.
- Single Multicore chips, and systems with multiple chips
- Multithreading – thread executions interleaved on a single processor

# Why Multiprocessors?

- Diminishing return from exploiting ILP with rising cost of power and chip area
- Easier to replicate cores
- Rise of web applications and cloud computing where natural parallelism found in web/data-intensive applications
  - Ex. vector operations  $\mathbf{A} * \mathbf{B} + \mathbf{C} * \mathbf{D}$

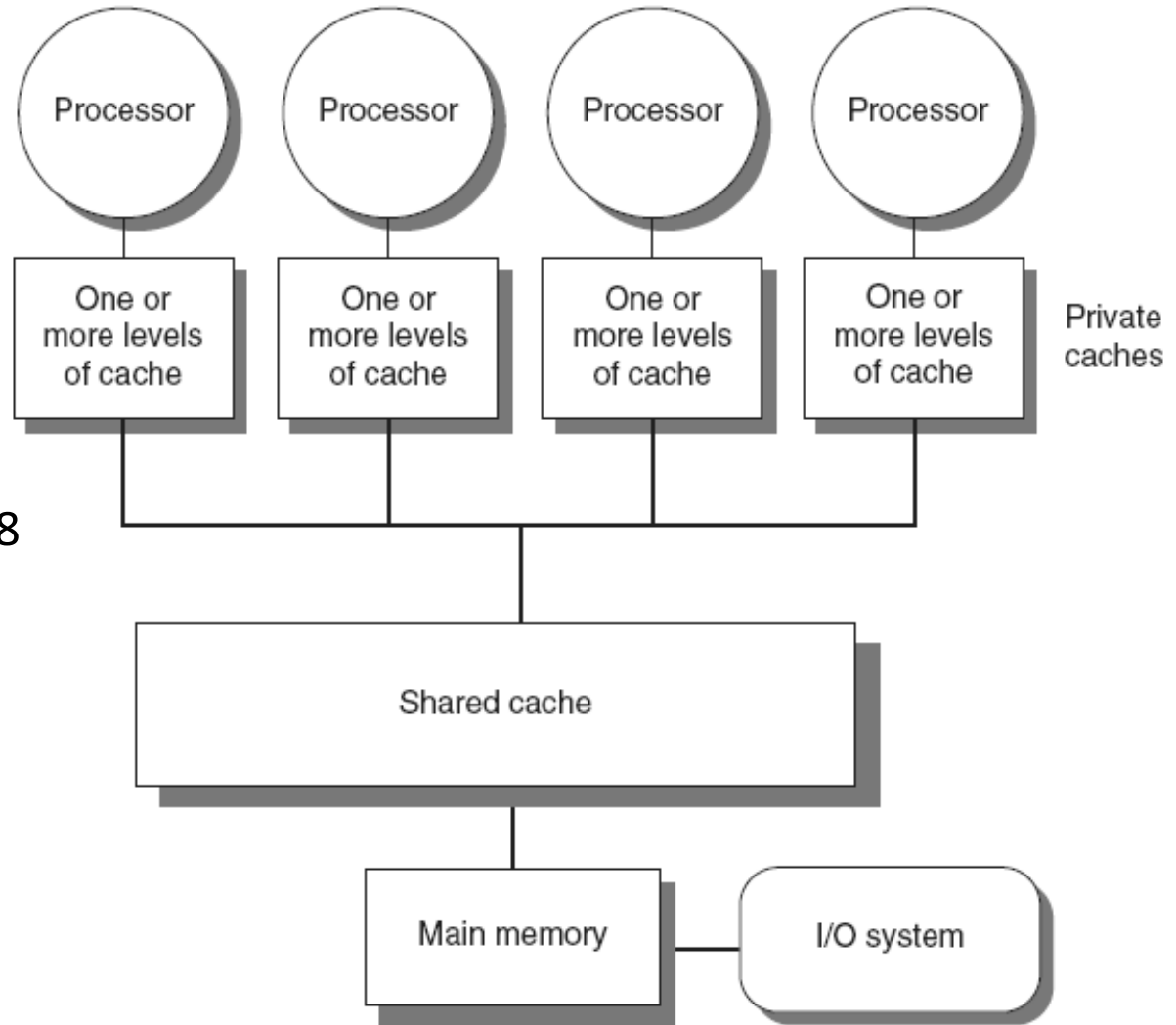
# Thread-Level Parallelism (TLP)

- Thread-Level parallelism
  - multiple running threads – multiple program counters
  - exploited through MIMD model
  - Targeted for tightly-coupled shared-memory multiprocessors
- Types of parallelism
  - Tightly coupled – threads collaborating for a single task
  - Loosely coupled – multiple programs running independently

# Exploiting TLP

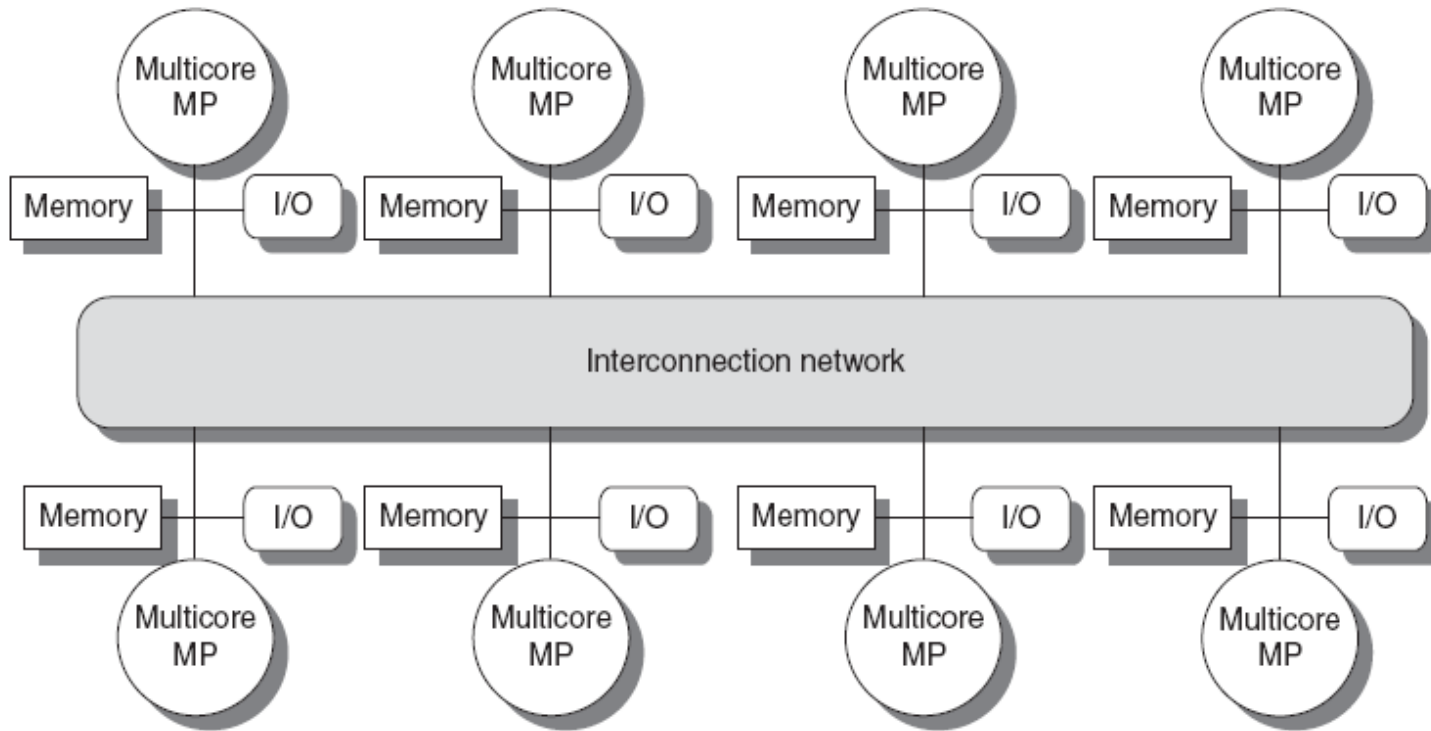
- Applications should contain abundant concurrency
  - For  $n$  processors, need  $n+$  independent threads
- Identified by programmers or OS
- Amount of computation assigned to each thread = grain size
- Thread grain size must be large
  - To reduce overheads associated with thread execution

# Symmetric multiprocessors (SMP)



- Small number of cores ( $\leq 8$  cores)
- Share single memory with uniform memory latency
- Also called UMA MP.

# Distributed Shared Memory (DSM)



- Memory distributed among processors
- Non-uniform memory access/latency (NUMA)
- Processors connected via direct (switched) and non-direct (multi-hop) interconnection networks
- Long latency to access remote memory.

# Challenges – Limited Parallelism

To achieve a speedup of 80 with 100 processors, what fraction of the original code can be sequential?

---

$$Speedup = \frac{1}{\frac{Fraction_{enhanced}}{Speedup_{enhanced}} + (1 - Fraction_{enhanced})}$$

$$80 = \frac{1}{\frac{Fraction_{parallel}}{100} + (1 - Fraction_{parallel})}$$

$$Fraction_{parallel} = 0.9975$$

Achieve linear speedup requires entire program to be parallel!



# Challenges – Remote Comm. Cost

Assume remote mem access latency is 100ns. Processors run at 4 GHz, and stall on remote accesses. Ideal CPI = 0.5. Compute speedup of code without comm versus 0.2% instructions require remote mem accesses.

---

$$\begin{aligned}\text{CPI} &= \text{Ideal CPI} + \text{stalls of remote comm} \\ &= 0.5 + 0.2\% * 400 = 1.3\end{aligned}$$

$$\text{Speedup} = 1.3/0.5 = 2.6$$

Caching is a technique to reduce latency

# **Centralized Shared-Memory Architecture**

# Caching

- Reduce latency of main memory and remote access
  - Also reduce contentions among memory accesses from different processors
- There are **private** & **shared** caches
- **Private data**: used by a single processor
- **Shared data**: used by multiple processors
  - Replicated in multiple private caches

# Cache Coherence Problem

- Processors may see different values through their private caches

Time	Event	Cache contents for processor A	Cache contents for processor B	Memory contents for location X
0				1

# Cache Coherence Problem

- Processors may see different values through their private caches

Time	Event	Cache contents for processor A	Cache contents for processor B	Memory contents for location X
0				1
1	Processor A reads X	1		1

# Cache Coherence Problem

- Processors may see different values through their private caches

Time	Event	Cache contents for processor A	Cache contents for processor B	Memory contents for location X
0				1
1	Processor A reads X	1		1
2	Processor B reads X	1	1	1

# Cache Coherence Problem

- Processors may see different values through their private caches

Time	Event	Cache contents for processor A	Cache contents for processor B	Memory contents for location X
0				1
1	Processor A reads X	1		1
2	Processor B reads X	1	1	1
3	Processor A stores 0 into X	0	1	0

# Cache Coherence

- **Coherence** – what values returned for reads
  - A read by a processor **A** to a location **X** that follows a write to **X** by **A** returns the value written by **A** if no other processors write in between
  - A read by processor **A** to location **X** after a write to **X** by processor **B** returns the written value if the read and write are sufficiently separated, and no other writes occur in between
  - Writes to the same location are serialized
- **Consistency** – when a written value seen by a read
  - Concerns reads & writes to different memory locations from multiple processors



# Enforcing Coherence

- Coherent caches provide:
  - *Migration*: movement of data – reduce latency
  - *Replication*: multiple copies of data – reduce latency & memory bandwidth demand
- Cache coherence protocols
  - *Snooping*
    - Every cache tracks sharing status of each cache block
    - Mem requests are broadcast on a bus to all caches
    - Writes serialized naturally by a single bus
  - *Directory based*
    - Sharing status of each cache block kept in one location

# Snoopy Coherence Protocols

## → Write invalidate

- On write, invalidate all other copies
- Use bus itself to serialize
  - Write cannot complete until bus access is obtained

## → Write update

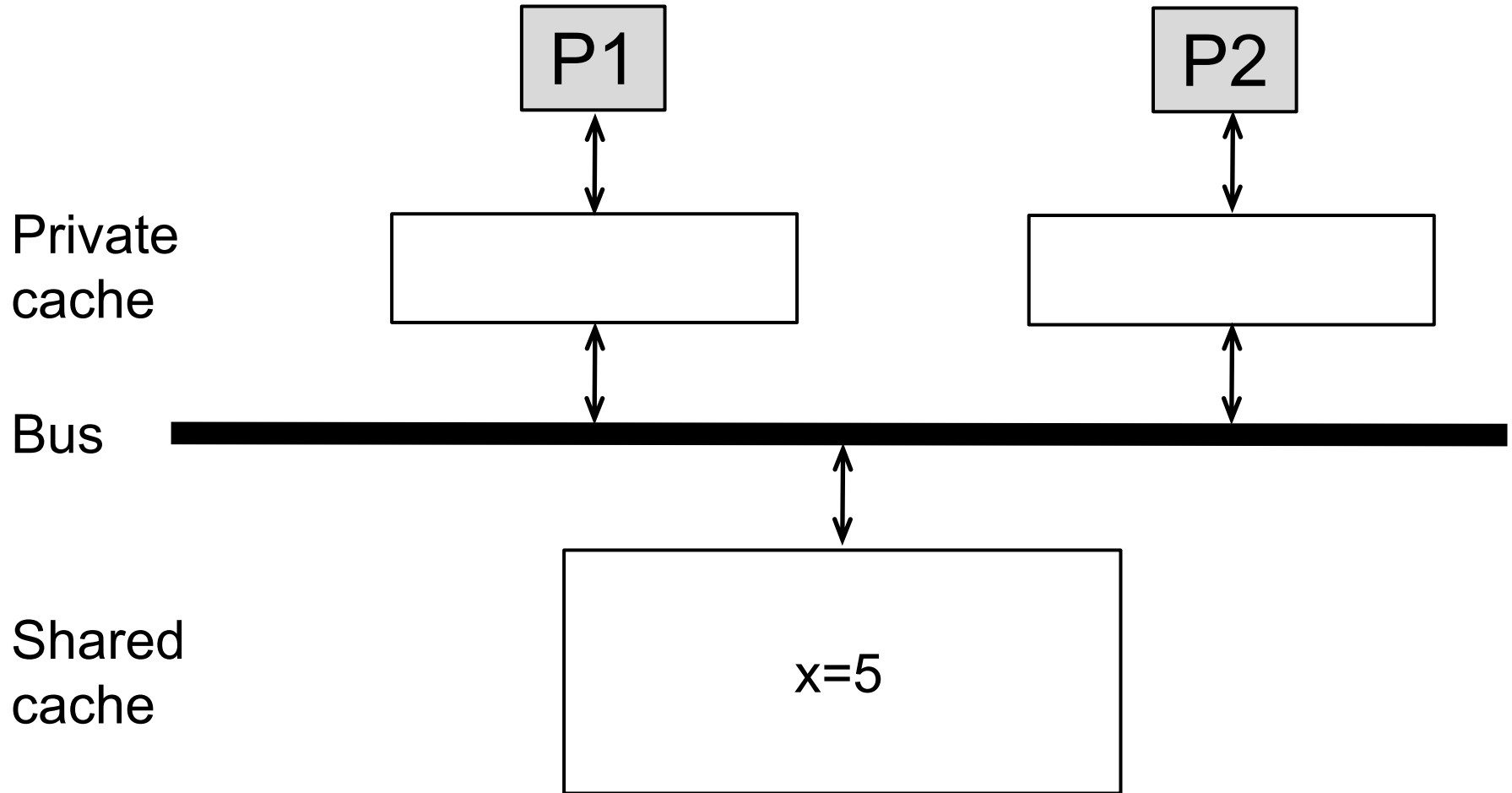
- On write, update all copies

→ **Invariant:** blocks in cache are always coherent

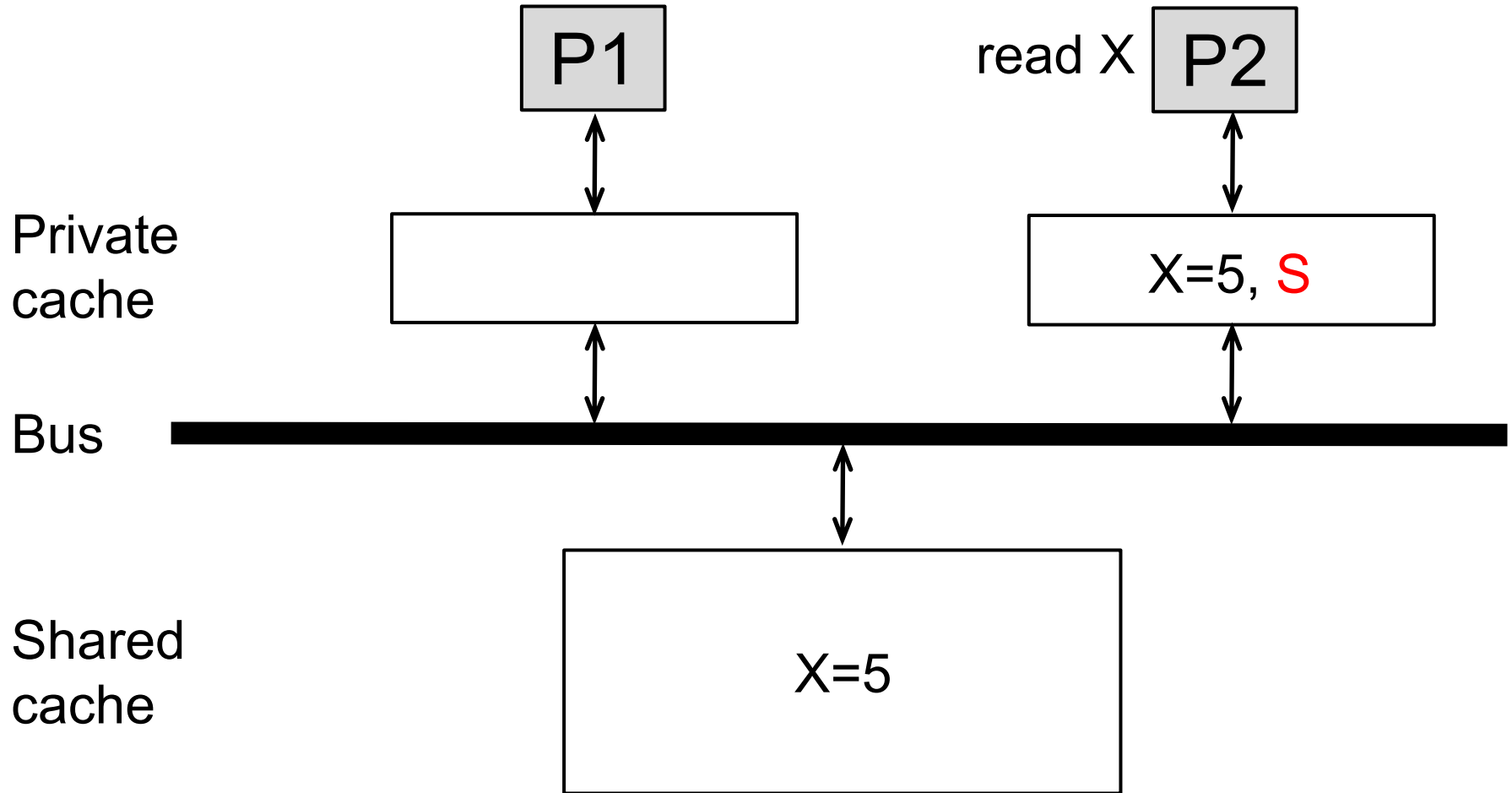
# Snoopy Coherence Protocols

- Each cache block is in one of following states
  - Invalid (I)
  - Shared (S)
  - Modified (M) – implies **exclusion** or **not shared**
- Locating an item when a read miss occurs
  - In write-back cache, the updated value must be sent to the requesting processor
- Cache lines marked as shared or exclusive/modified
  - Only writes to shared lines need an invalidate broadcast
    - After this, the line is marked as exclusive

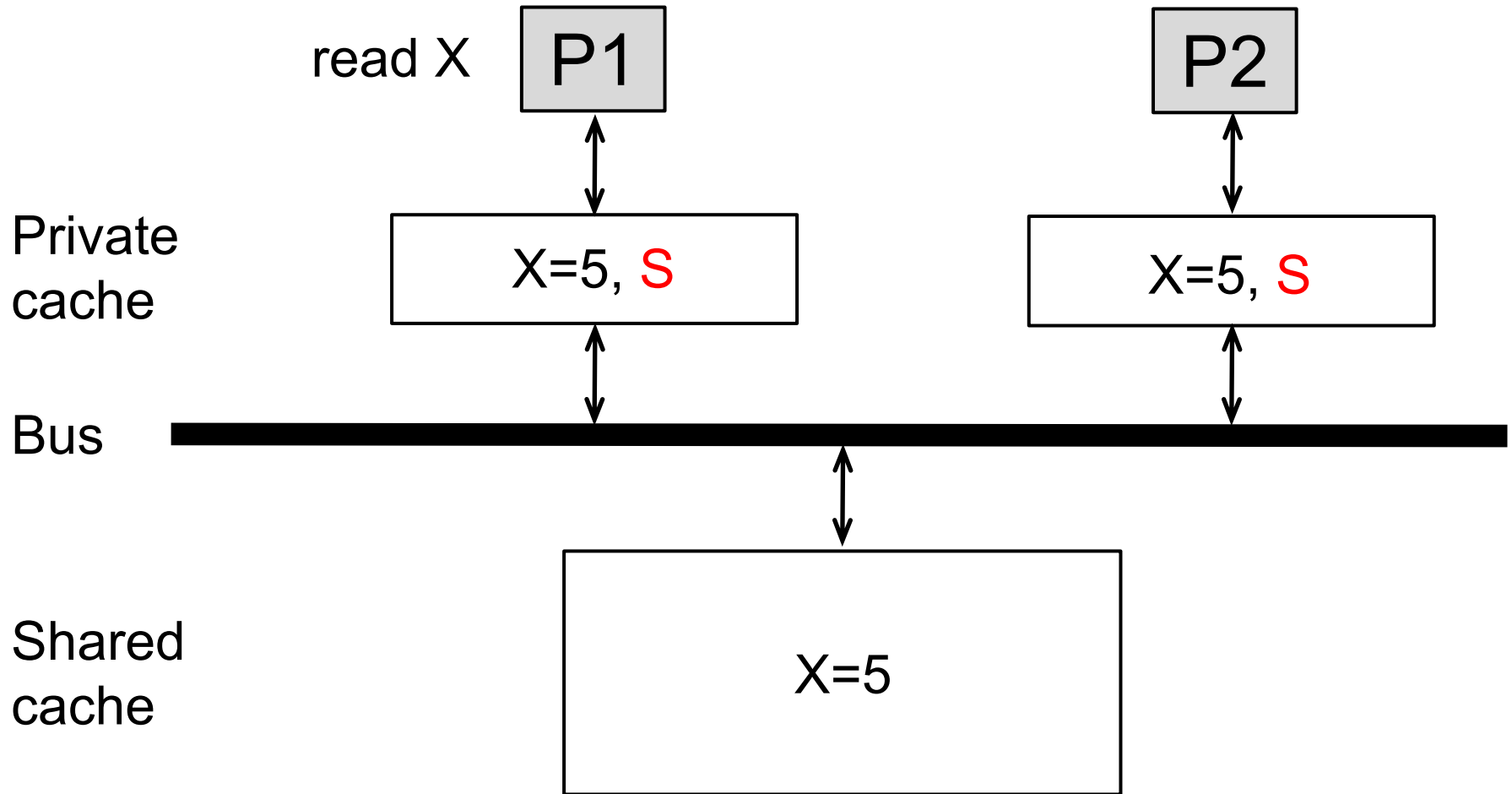
# Snoopy Coherence Protocols



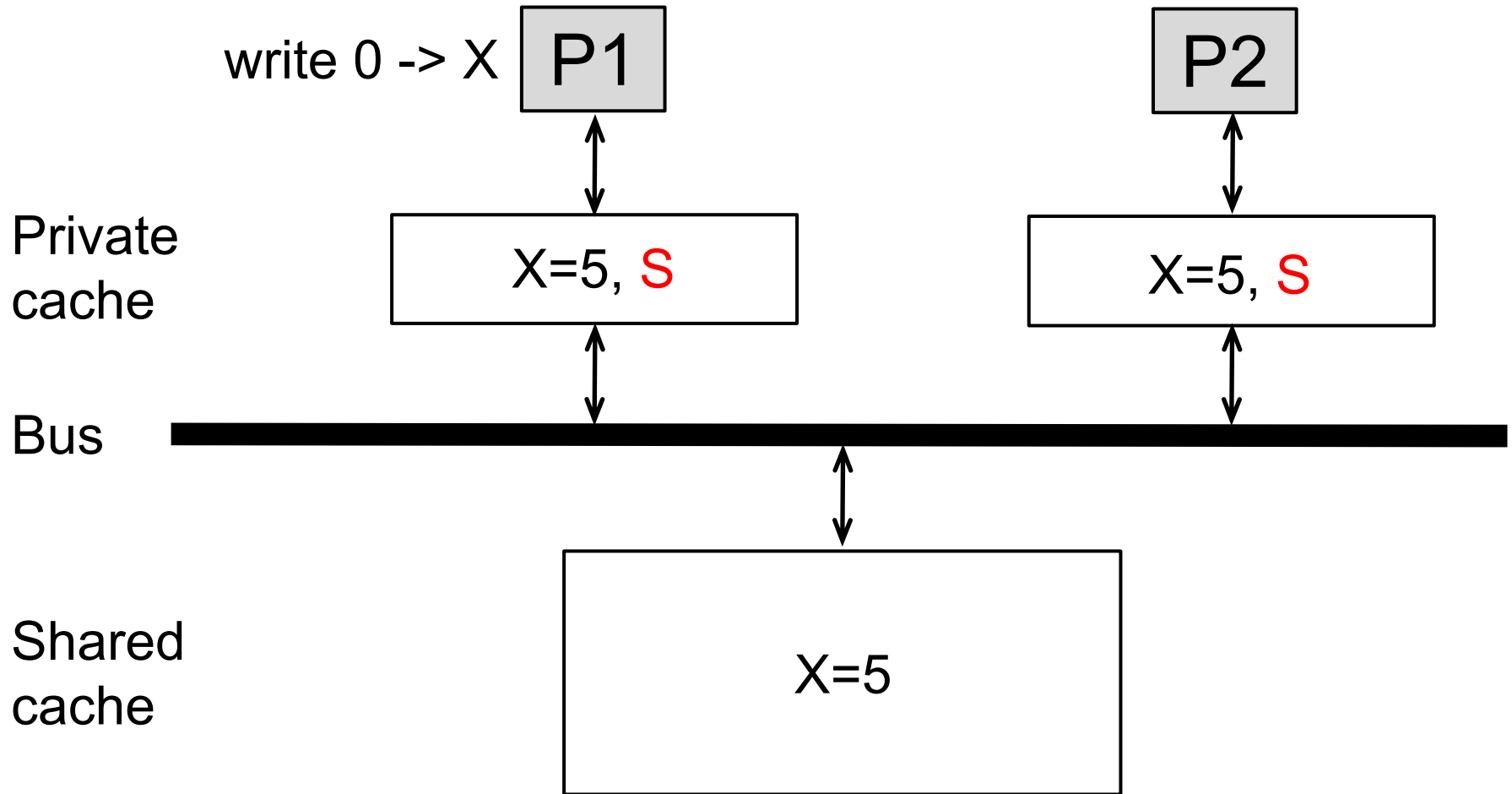
# Snoopy Coherence Protocols



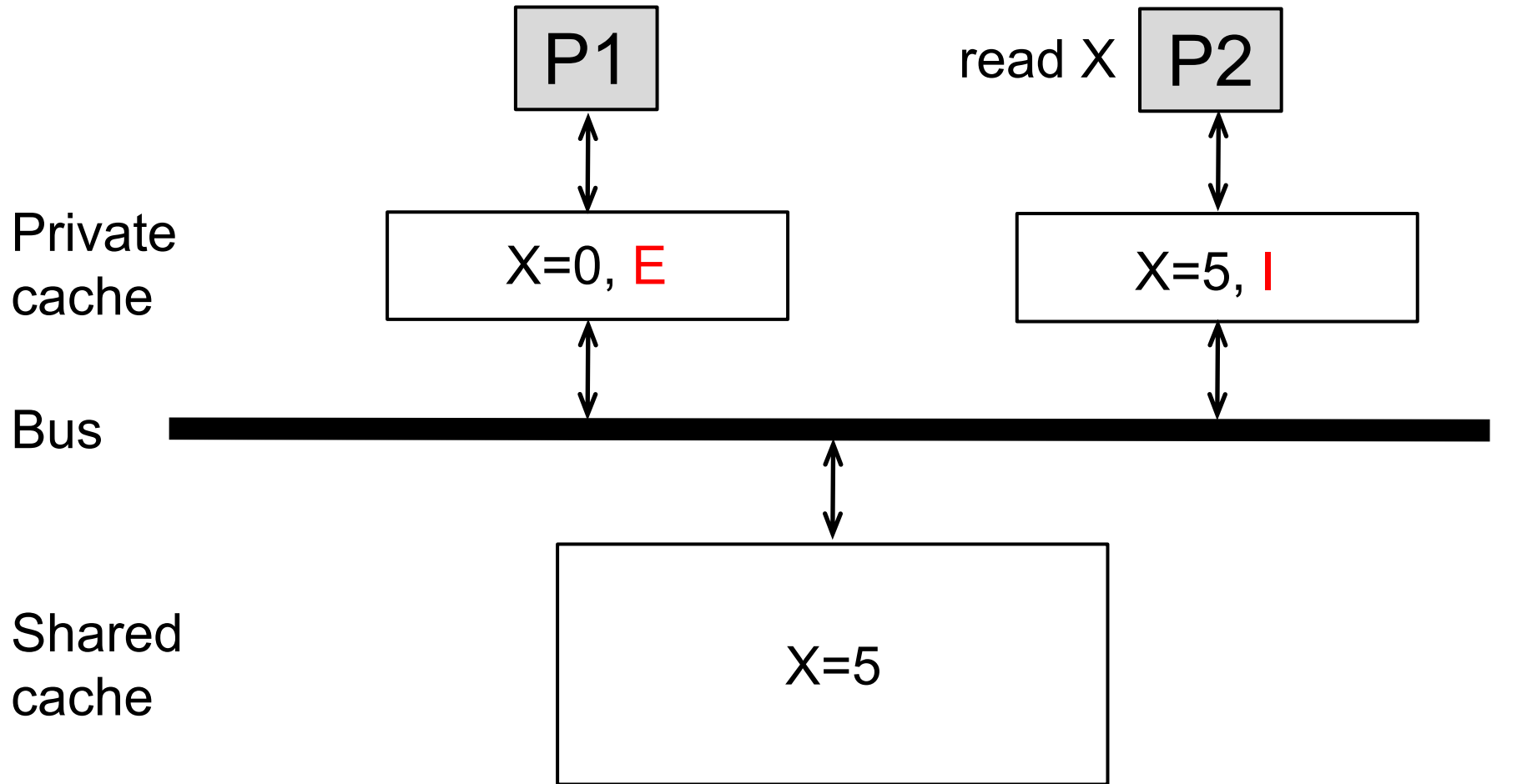
# Snoopy Coherence Protocols



# Snoopy Coherence Protocols



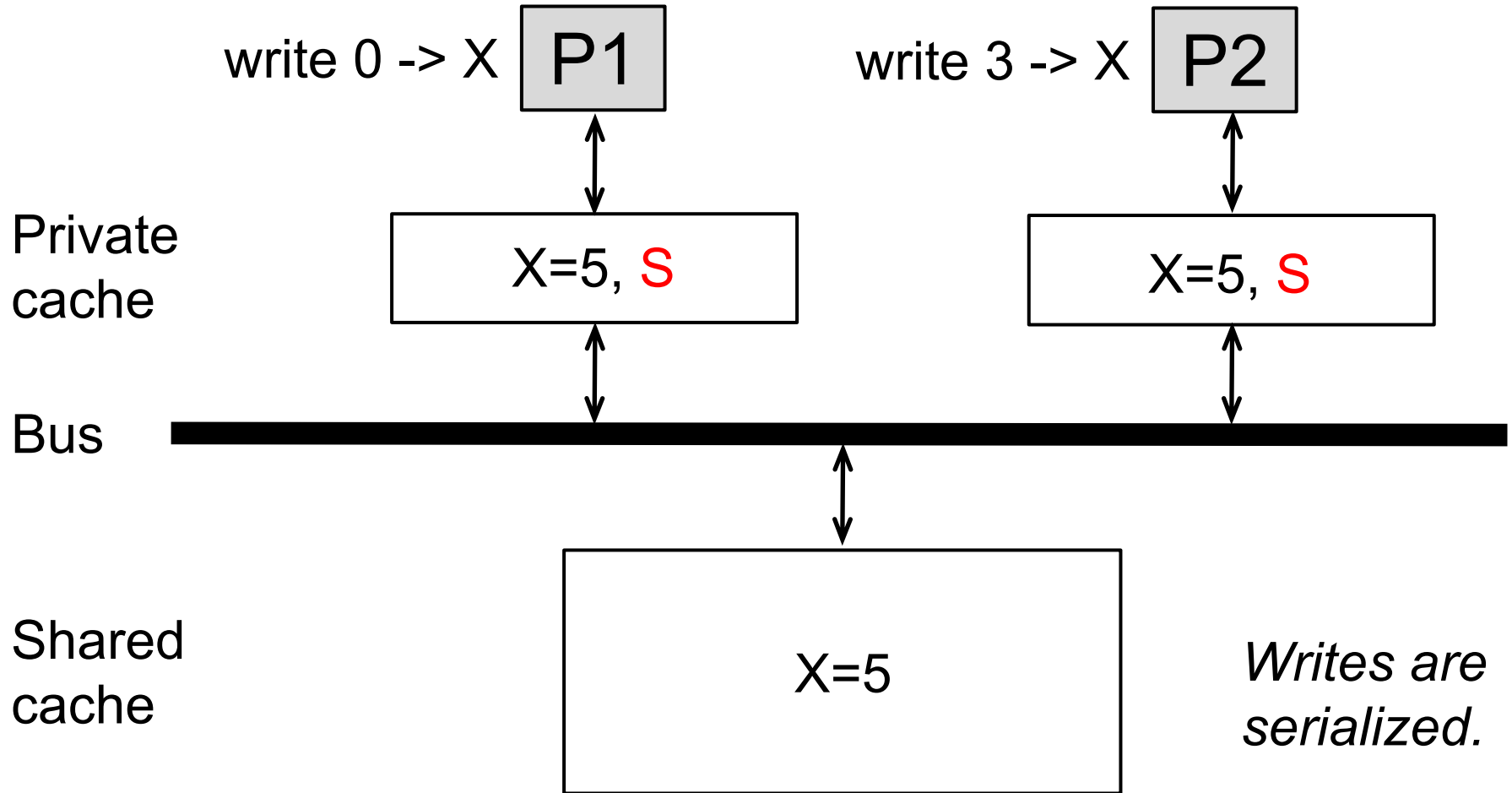
# Snoopy Coherence Protocols



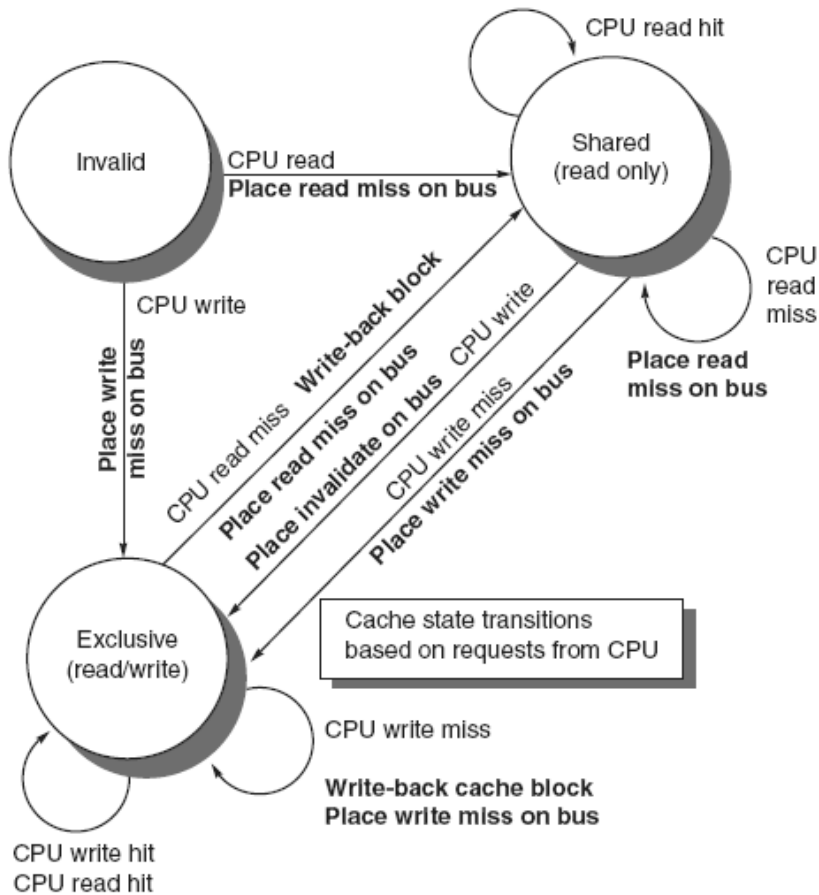
*where does P2 find data for X?*



# Snoopy Coherence Protocols



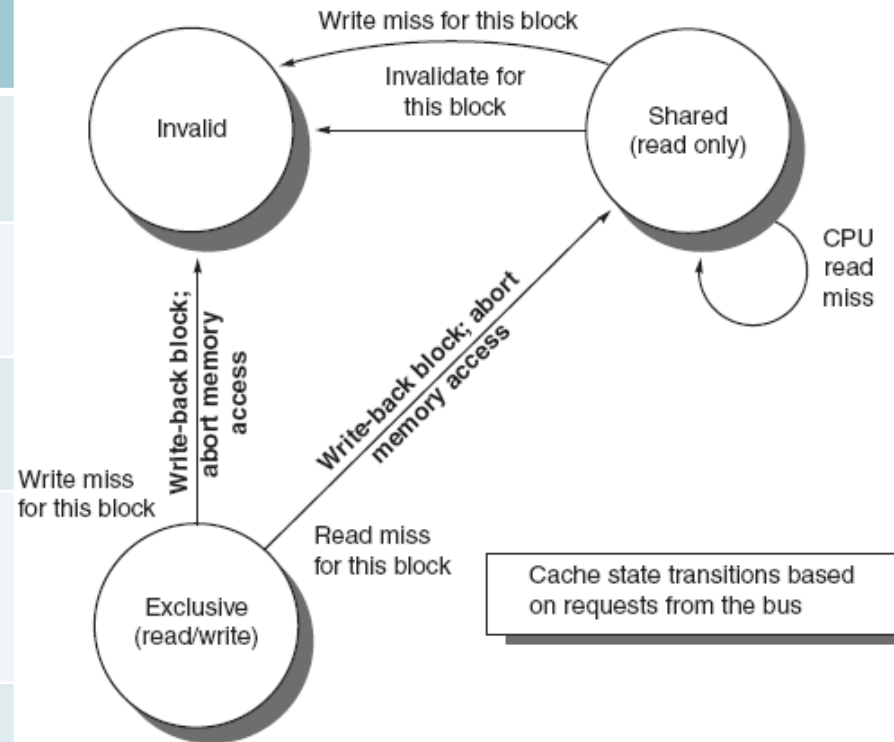
# Snoopy Coherence Protocols



State	CPU request	Action
Invalid	read	read miss on bus
Invalid	write	write miss on bus
Shared	read hit	no action
Shared	read miss	read miss on bus
Shared	write hit	invalidate on bus
Shared	write miss	write miss on bus
Exclusive	read/write hit	no action
Exclusive	write miss	write back write miss on bus
Exclusive	read miss	write back read miss on bus

# Snoopy Coherence Protocols

State	Bus request	Action
Shared	write miss	invalidate
Shared	invalidate	invalidate
Shared	read miss	no action
Exclusive	write miss	write back
Exclusive	read miss	write back put data on bus



# Snoopy Coherence Protocols

Request	Source	State of addressed cache block	Type of cache action	Function and explanation
Read hit	Processor	Shared or modified	Normal hit	Read data in local cache.
Read miss	Processor	Invalid	Normal miss	Place read miss on bus.
Read miss	Processor	Shared	Replacement	Address conflict miss: place read miss on bus.
Read miss	Processor	Modified	Replacement	Address conflict miss: write-back block, then place read miss on bus.
Write hit	Processor	Modified	Normal hit	Write data in local cache.
Write hit	Processor	Shared	Coherence	Place invalidate on bus. These operations are often called upgrade or <i>ownership</i> misses, since they do not fetch the data but only change the state.
Write miss	Processor	Invalid	Normal miss	Place write miss on bus.
Write miss	Processor	Shared	Replacement	Address conflict miss: place write miss on bus.
Write miss	Processor	Modified	Replacement	Address conflict miss: write-back block, then place write miss on bus.
Read miss	Bus	Shared	No action	Allow shared cache or memory to service read miss.
Read miss	Bus	Modified	Coherence	Attempt to share data: place cache block on bus and change state to shared.
Invalidate	Bus	Shared	Coherence	Attempt to write shared block; invalidate the block.
Write miss	Bus	Shared	Coherence	Attempt to write shared block; invalidate the cache block.
Write miss	Bus	Modified	Coherence	Attempt to write block that is exclusive elsewhere; write-back the cache block and make its state invalid in the local cache.

# Snoopy Coherence Protocols

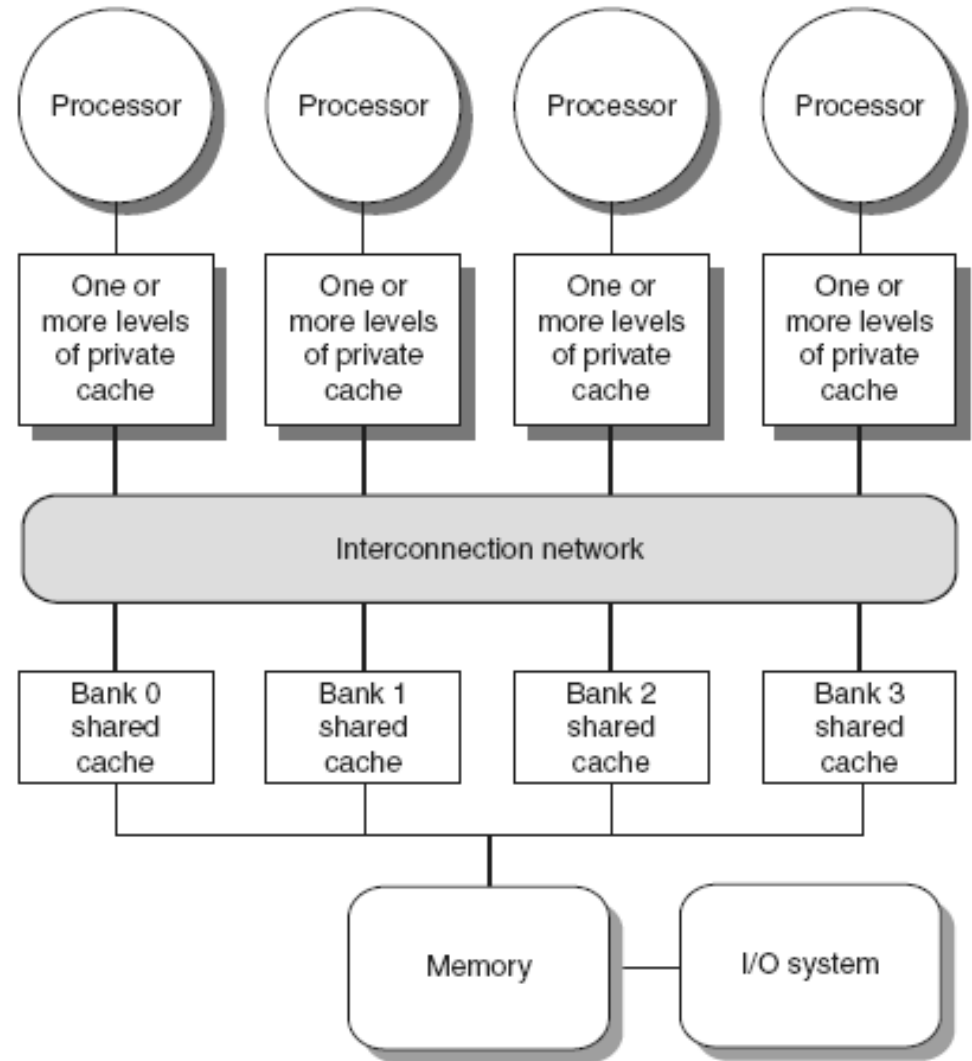
- Complications for the basic MSI protocol:
  - Operations are not atomic
    - E.g. detect miss, acquire bus, receive a response
    - Creates possibility of **deadlock** and **races**
    - One solution: processor that sends invalidate can hold bus until other processors receive the invalidate
- Extensions – optimize performance
  - Add exclusive (E) state to indicate clean block in only one cache (MESI protocol)
    - No need to invalidate on a write to a block in E state.
  - MOESI – add Owned (O) state
    - Dirty block is in local caches, but not in the shared cache

# Limitations

- Bus or shared cache accesses become bottleneck
  - limit #core to 4 – 8
- Accesses to private caches also a bottleneck
  - handle accesses from core and bus due to snooping
- Solution - increase memory bandwidth
  - see next slide

# Coherence Protocols: Extensions

- ➔ Shared memory bus and snooping bandwidth is bottleneck for scaling symmetric multiprocessors
  - ➔ Use crossbars or point-to-point networks with banked memory

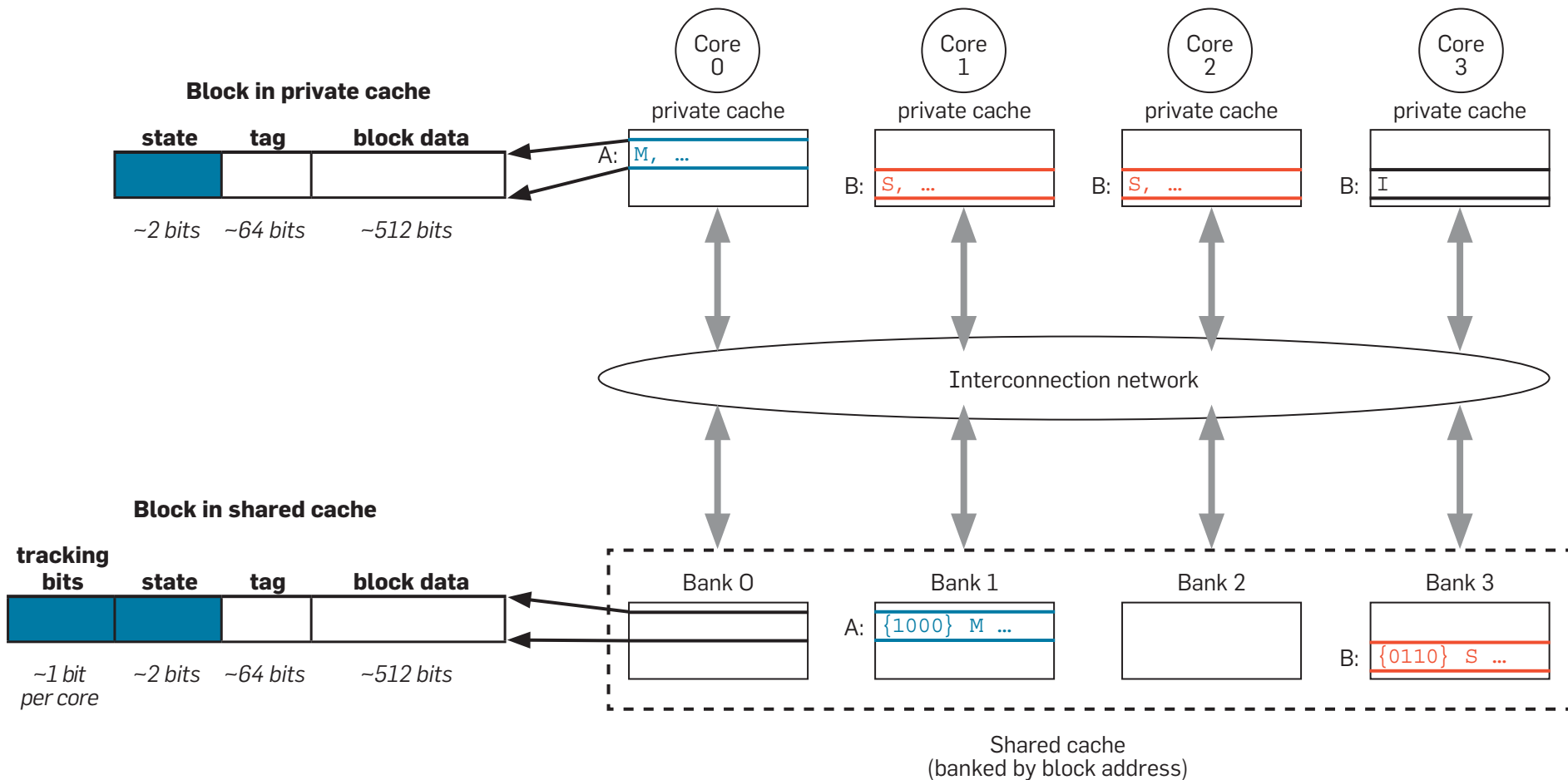


# **Distributed Shared-Memory Architecture**



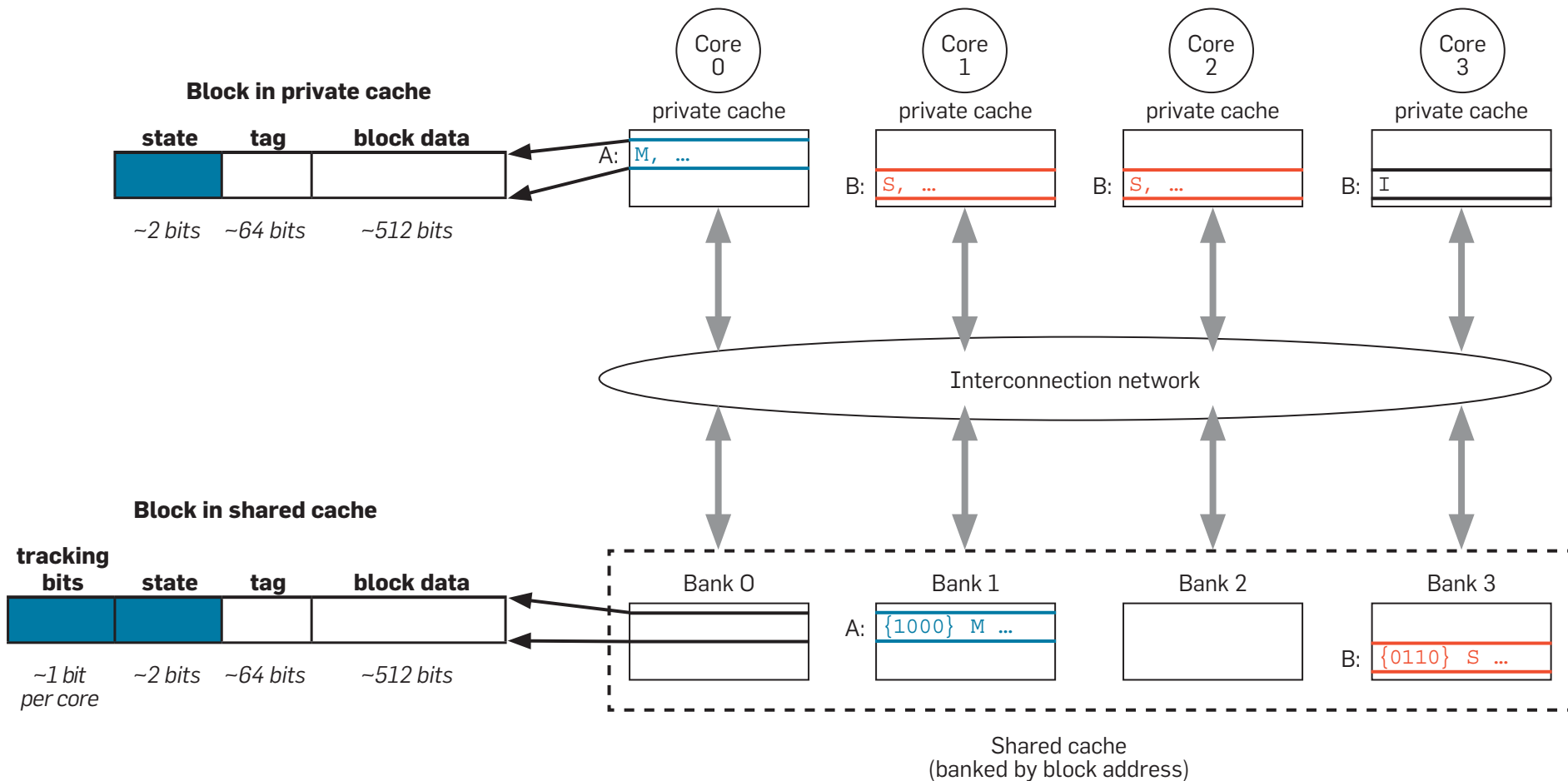
# Directory Protocols

- Directory keeps track of every block
  - Sharing caches and dirty status of each block



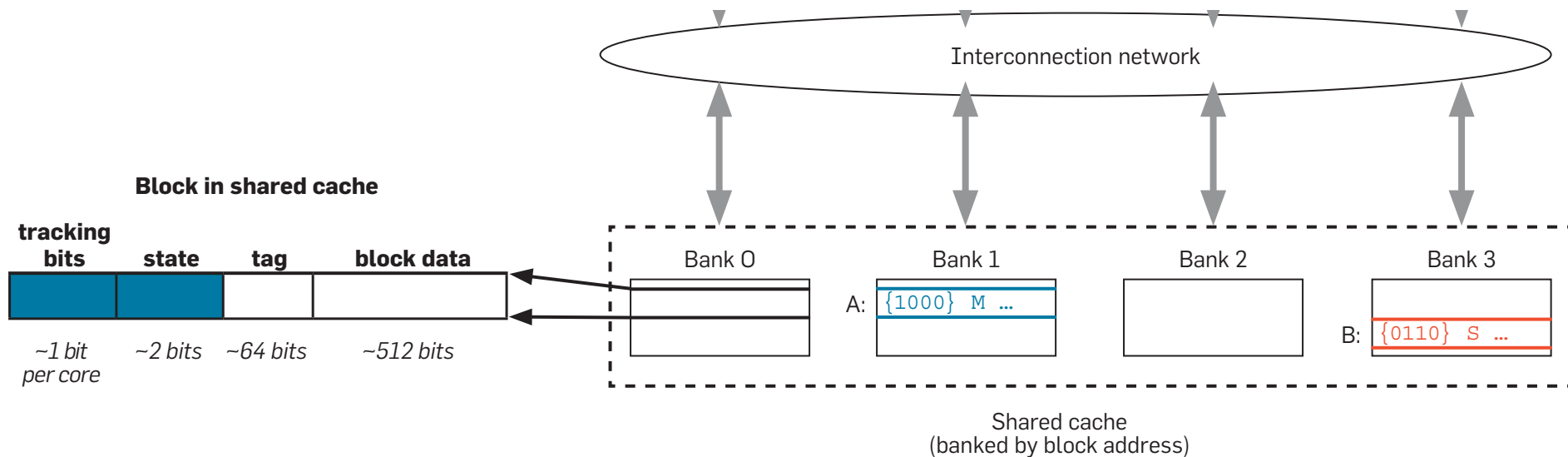
# Directory Protocols

- Implement in shared L3 cache
  - Status bit vector, its size = # cores for each block in L3

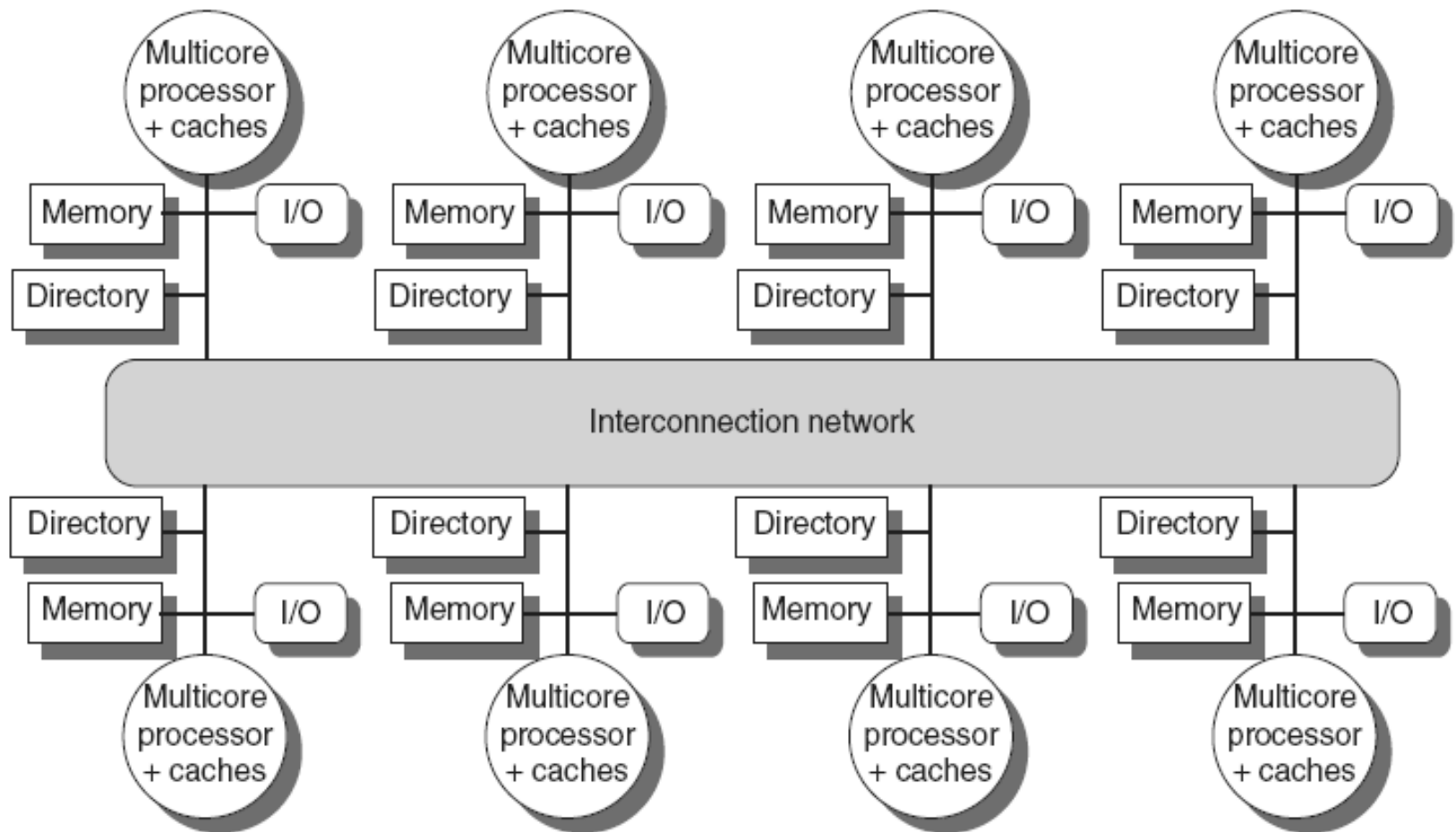


# Directory Protocols

- Not scalable beyond L3 cache – centralized cache is bottleneck
- Must be implemented in a distributed fashion
  - Each distributed memory has a directory
  - size = # memory blocks X # nodes



# Directory Protocols



Local directory only stores coherence information of cache blocks in local memory

# Directory Protocols

- Directory maintains block states and sends invalidation messages
  - Tracks states of all local memory blocks (simplest sol.)
- For each block, maintain state:
  - **Shared**
    - One or more nodes have the block cached, value in memory is up-to-date
  - **Uncached – invalid**
  - **Modified – Exclusive**
    - Exactly one node has a copy of the cache block, value in memory is out-of-date
    - This node is the owner

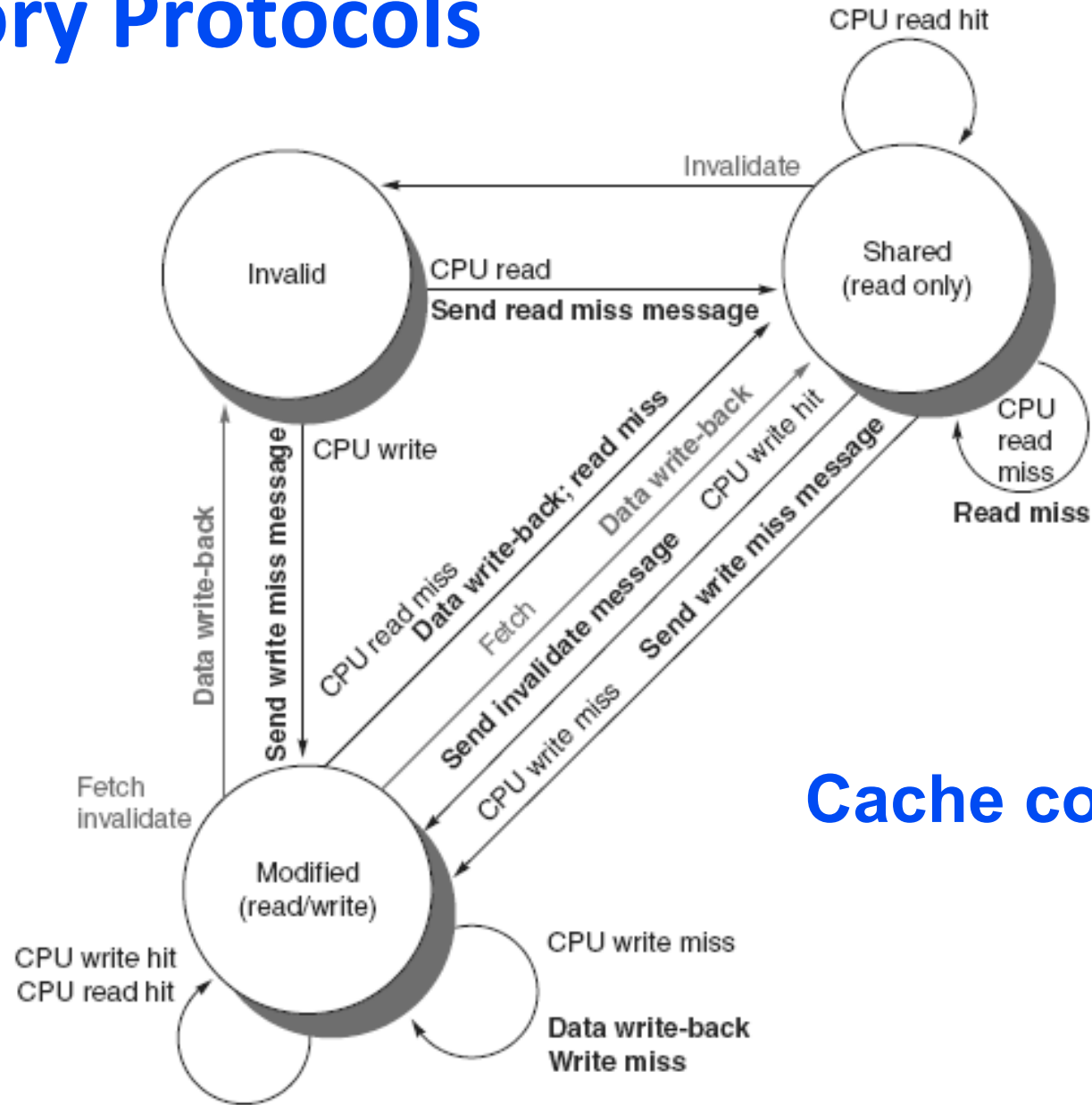
# Messages

Message type	Source	Destination	Message contents	Function of this message
Read miss	Local cache	Home directory	P, A	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	A	Request to send invalidates to all remote caches that are caching the block at address A.
Invalidate	Home directory	Remote cache	A	Invalidate a shared copy of data at address A.
Fetch	Home directory	Remote cache	A	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	A	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.

Local node: source of requests

Home node: destination of the requests

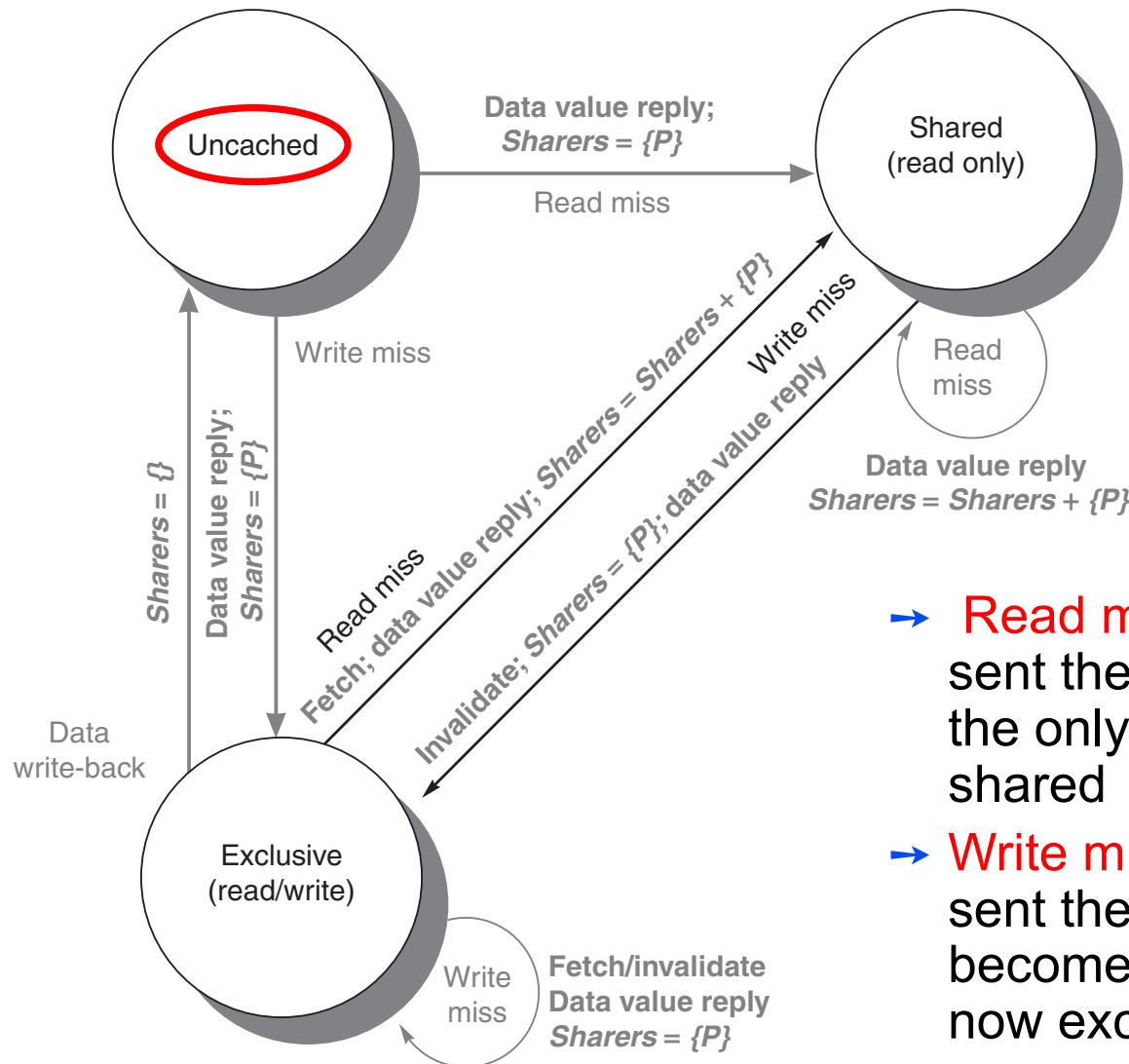
# Directory Protocols



# Cache controller

# Directory Protocols

## Directory controller

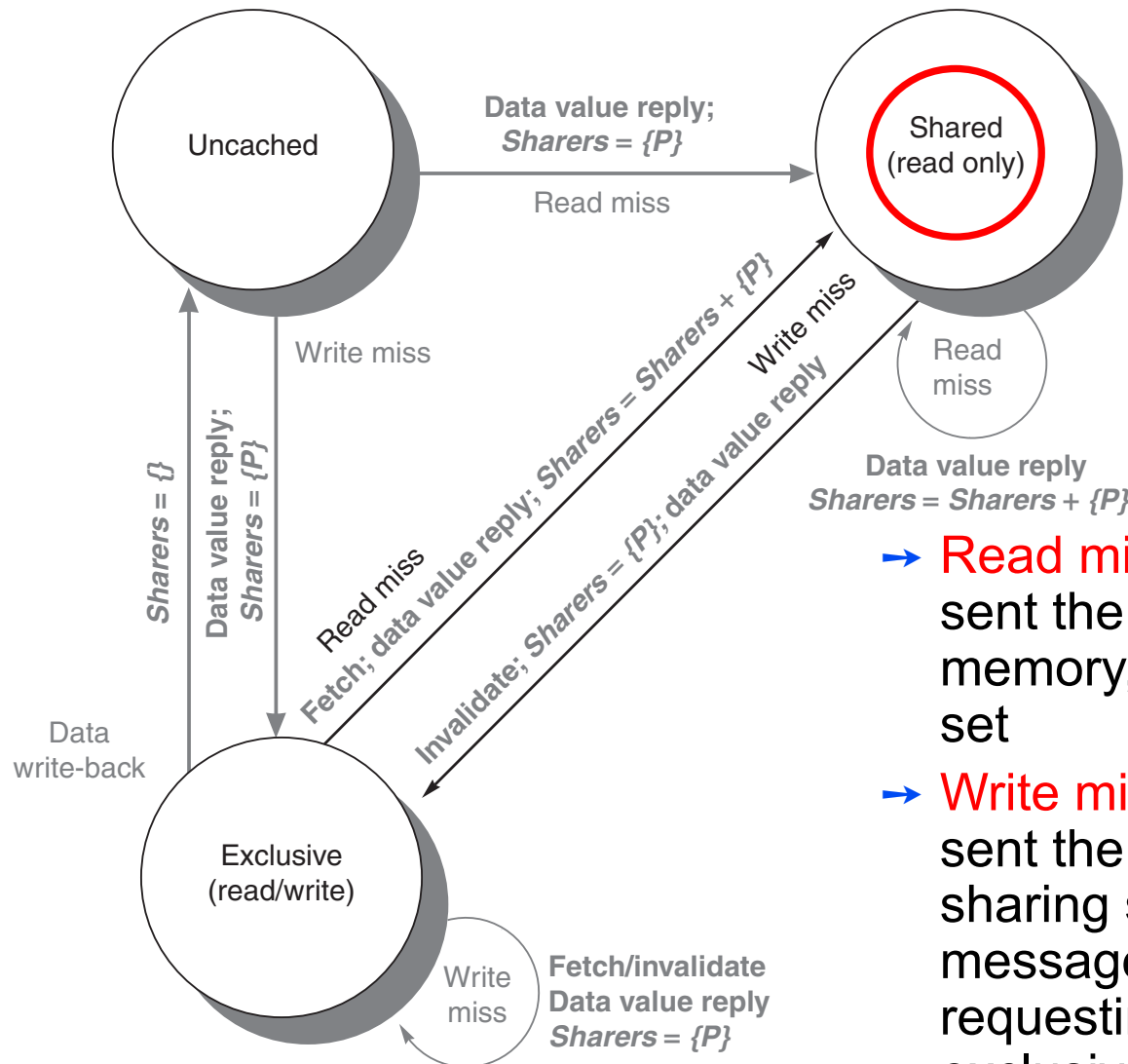


- **Read miss** – Requesting node is sent the requested data and is made the only sharing node, block is now shared
- **Write miss** – The requesting node is sent the requested data and becomes the sharing node, block is now exclusive



# Directory Protocols

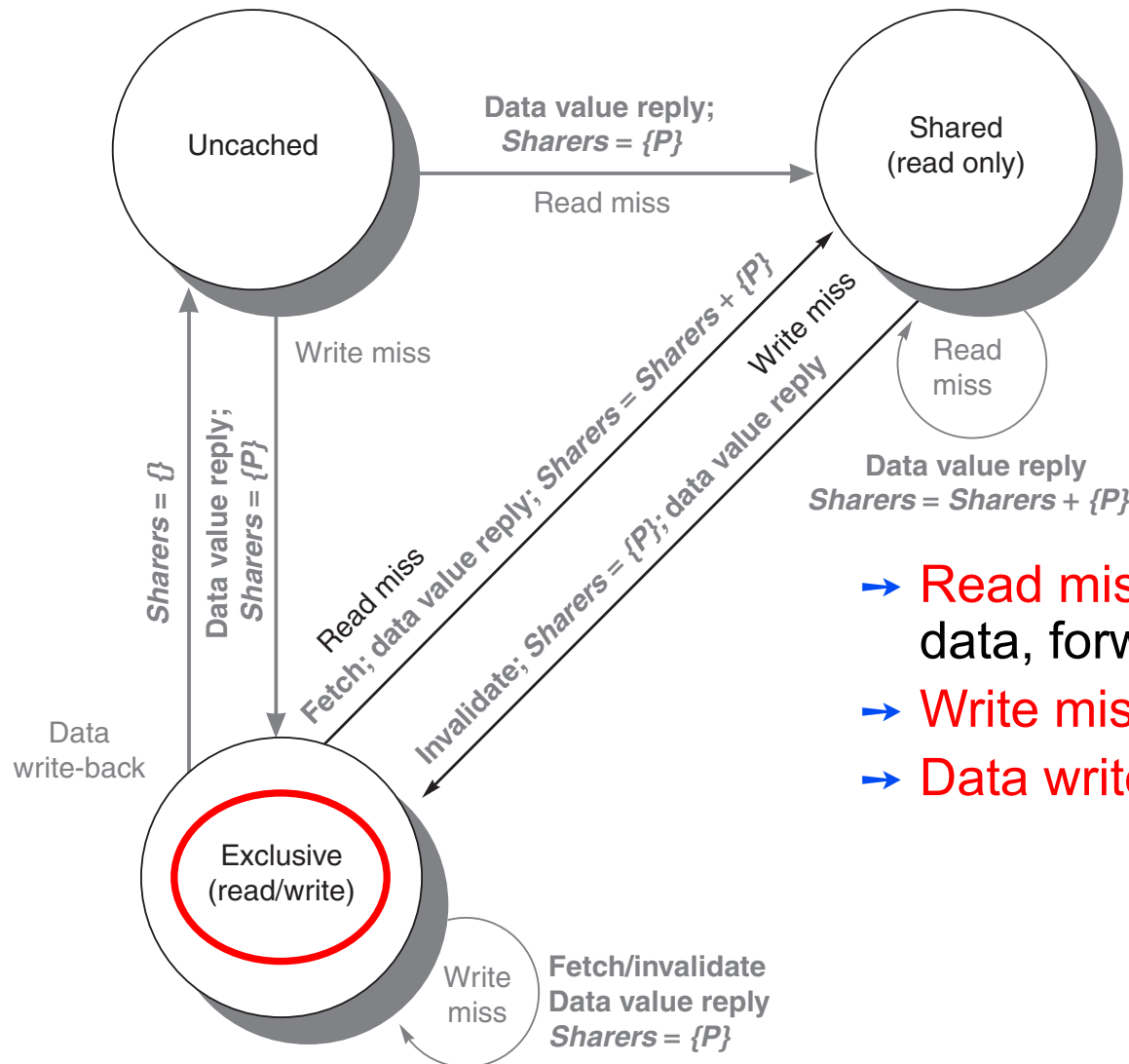
## Directory controller



- **Read miss** – The requesting node is sent the requested data from memory, node is added to sharing set
- **Write miss** – The requesting node is sent the value, all nodes in the sharing set are sent invalidate messages, sharing set only contains requesting node, block is now exclusive

# Directory Protocols

## Directory controller



- **Read miss** – ask owner to send data, forward to the requesting node
- **Write miss** – invalidate owner's copy
- **Data write back**

# Directory Protocols

## → For uncached (Invalid) block:

### → Read miss

- Requesting node is sent the requested data and is made the only sharing node, block is now shared

### → Write miss

- The requesting node is sent the requested data and becomes the sharing node, block is now exclusive

## → For shared block:

### → Read miss

- The requesting node is sent the requested data from memory, node is added to sharing set

### → Write miss

- The requesting node is sent the value, all nodes in the sharing set are sent invalidate messages, sharing set only contains requesting node, block is now exclusive

# Directory Protocols

→ For **exclusive** block:

→ **Read miss**

→ The owner is sent a data fetch message, block becomes shared, owner sends data to the directory, data written back to memory, sharers set contains old owner and requestor

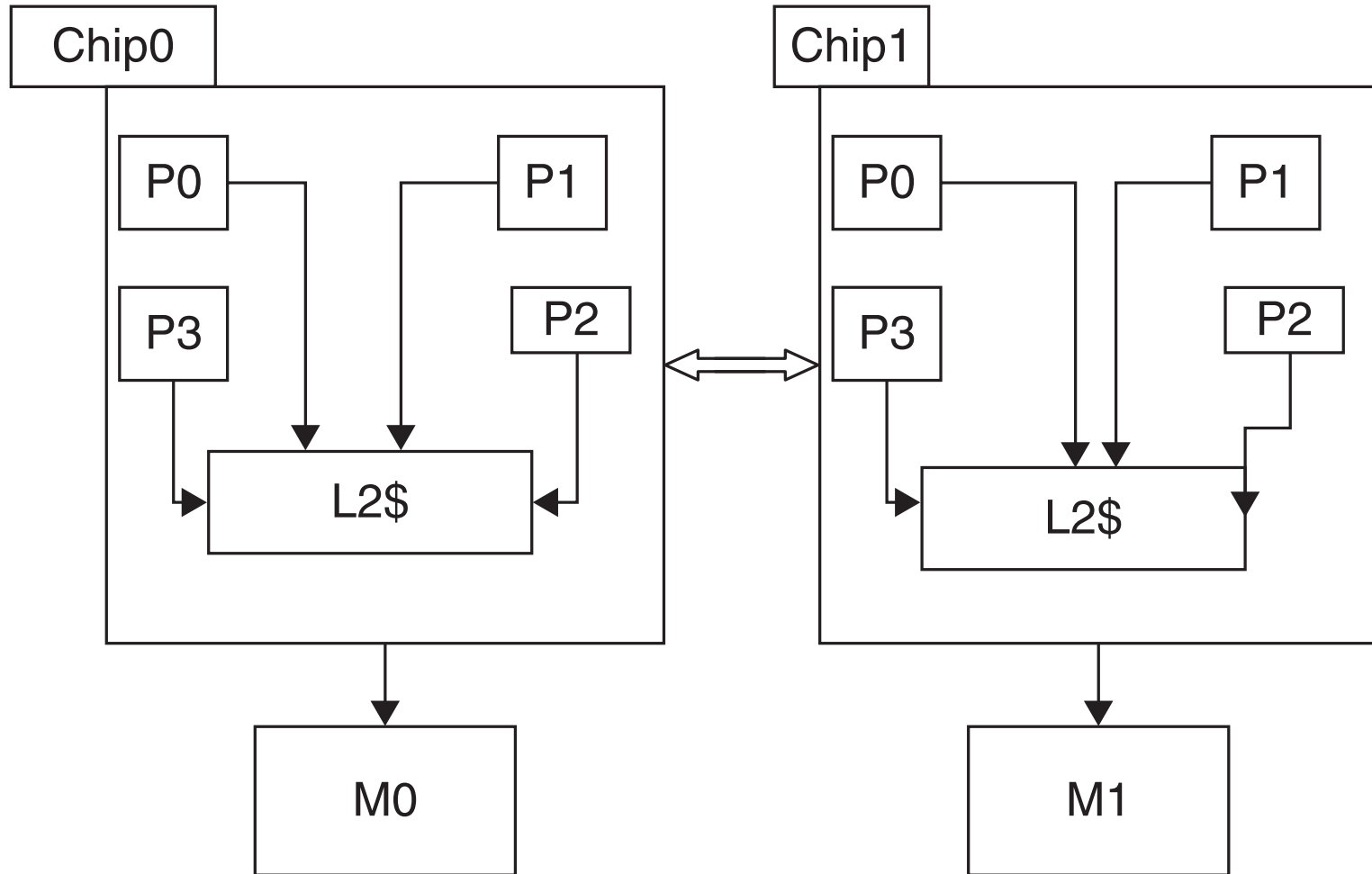
→ **Data write back**

→ Block becomes uncached, sharer set is empty

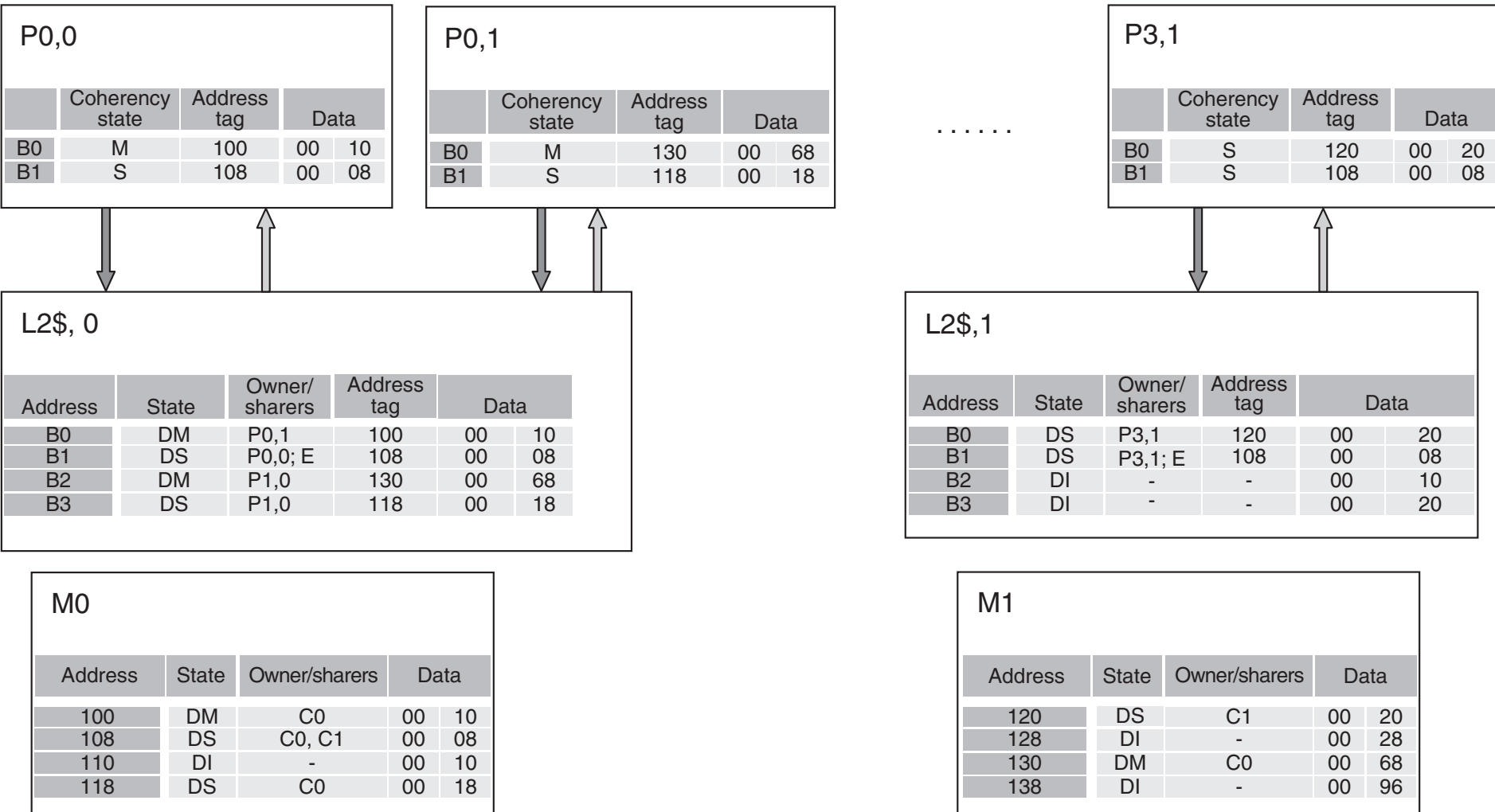
→ **Write miss**

→ Message is sent to old owner to invalidate and send the value to the directory, requestor becomes new owner, block remains exclusive

# Directory Protocols – Example



# P0,0: read 110



**P0,0: write 128 <-- 78**

P0,0

	Coherency state	Address tag	Data	
B0	M	100	00	10
B1	S	108	00	08

P0,1

	Coherency state	Address tag	Data	
B0	M	130	00	68
B1	S	118	00	18

P3,1

	Coherency state	Address tag	Data	
B0	S	120	00	20
B1	S	108	00	08

.....

L2\$, 0

Address	State	Owner/sharers	Address tag	Data	
B0	DM	P0,1	100	00	10
B1	DS	P0,0; E	108	00	08
B2	DM	P1,0	130	00	68
B3	DS	P1,0	118	00	18

L2\$,1

Address	State	Owner/sharers	Address tag	Data	
B0	DS	P3,1	120	00	20
B1	DS	P3,1; E	108	00	08
B2	DI	-	-	00	10
B3	DI	-	-	00	20

M0

Address	State	Owner/sharers	Data	
100	DM	C0	00	10
108	DS	C0, C1	00	08
110	DI	-	00	10
118	DS	C0	00	18

M1

Address	State	Owner/sharers	Data	
120	DS	C1	00	20
128	DI	-	00	28
130	DM	C0	00	68
138	DI	-	00	96

# Backup



# Models of Memory Consistency

Processor 1:

A=0

...

A=1

if (B==0) ...

Processor 2:

B=0

...

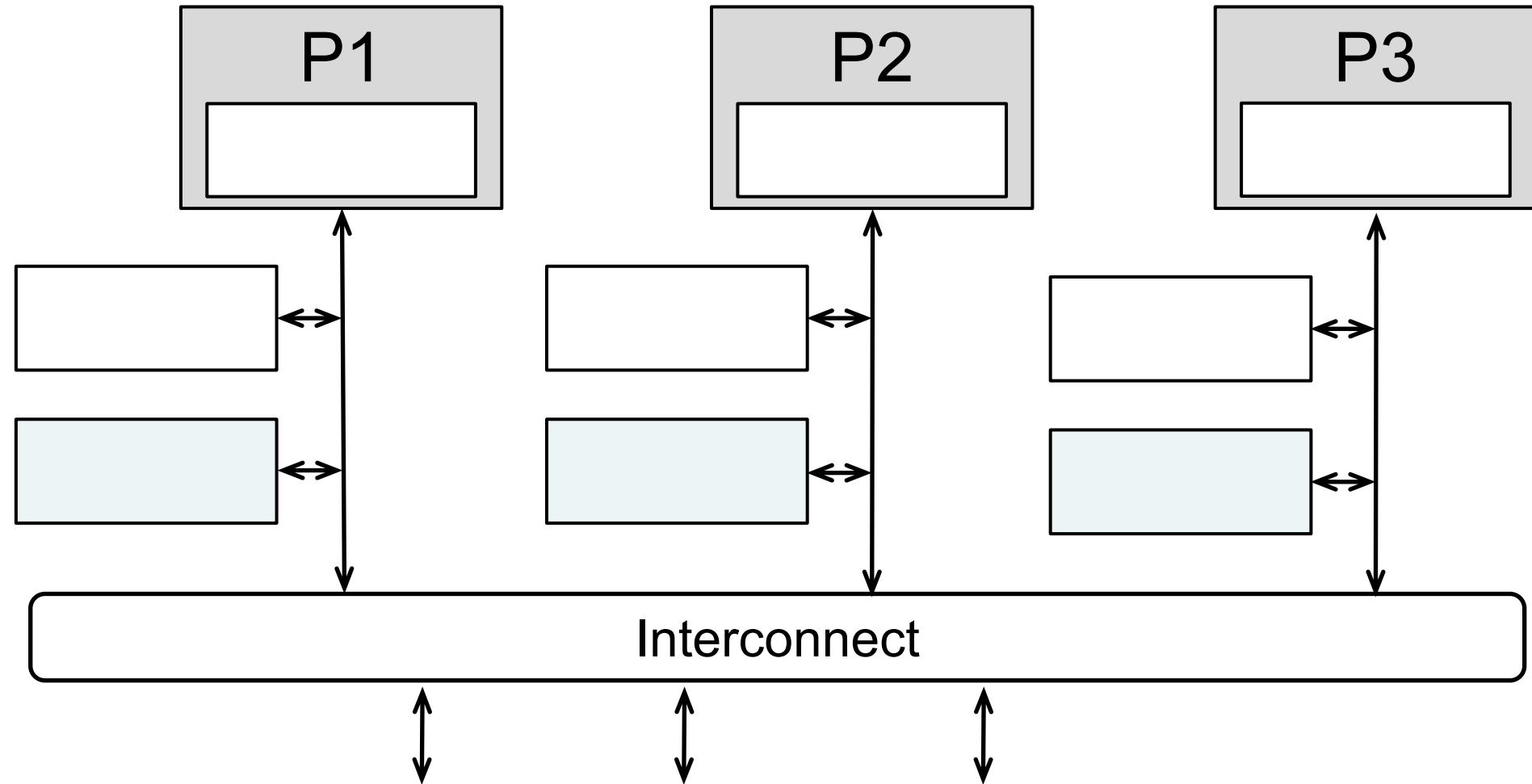
B=1

if (A==0) ...

- Possible for both **if** to be evaluated as true?
  - If so, what could be the cause?
- **Sequential consistency:**
  - Result of execution should be the same as long as:
    - Accesses on each processor were kept in order

A different ordering of processors were arbitrarily

# Directory Protocols – Example



# Snoopy Coherence Protocols

## → Write invalidate

→ On write, invalidate all other copies

→ Use bus itself to serialize

→ Write cannot complete until bus access is obtained

Processor activity	Bus activity	Contents of processor A's cache	Contents of processor B's cache	Contents of memory location X
				0
Processor A reads X	Cache miss for X	0		0
Processor B reads X	Cache miss for X	0	0	0
Processor A writes a 1 to X	Invalidation for X	1		0
Processor B reads X	Cache miss for X	1	1	1

## → Write update

# Performance

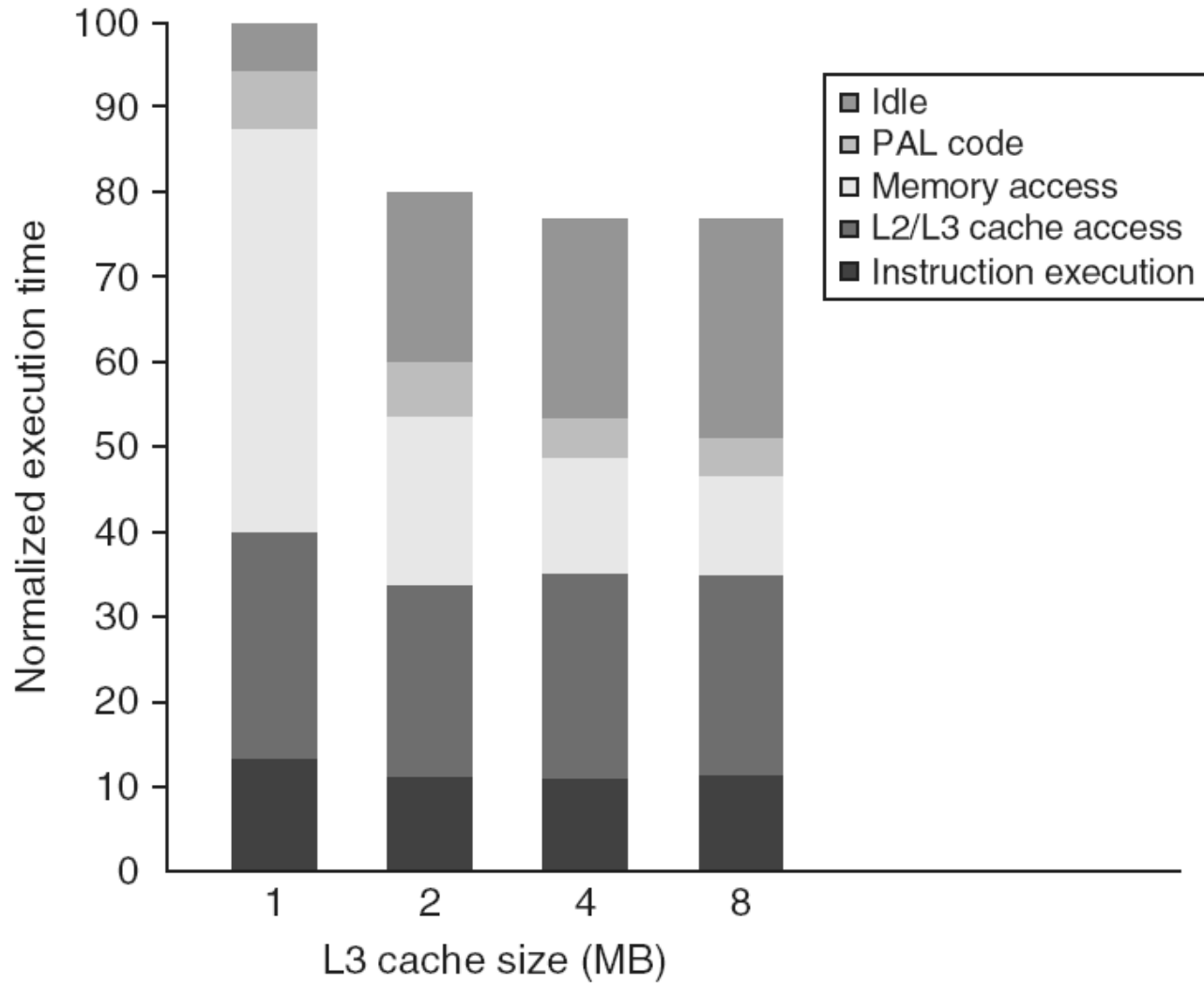
- Coherence influences cache miss rate
  - Coherence misses
    - True sharing misses
      - Write to shared block (transmission of invalidation)
      - Read an invalidated block
    - False sharing misses
      - Read an unmodified word in an invalidated block

# Coherence Protocols

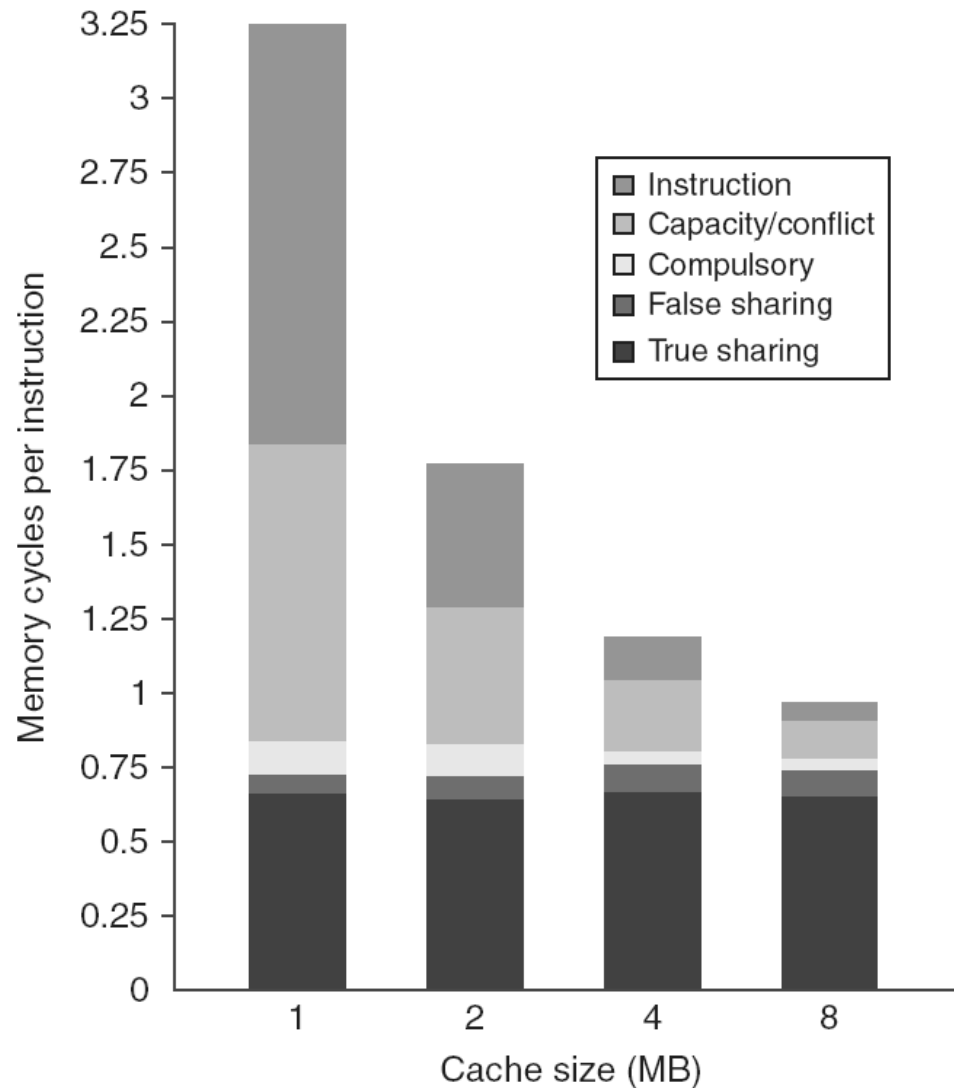
## → AMD Opteron:

- Memory directly connected to each multicore chip in NUMA-like organization
- Implement coherence protocol using point-to-point links
- Use explicit acknowledgements to order operations

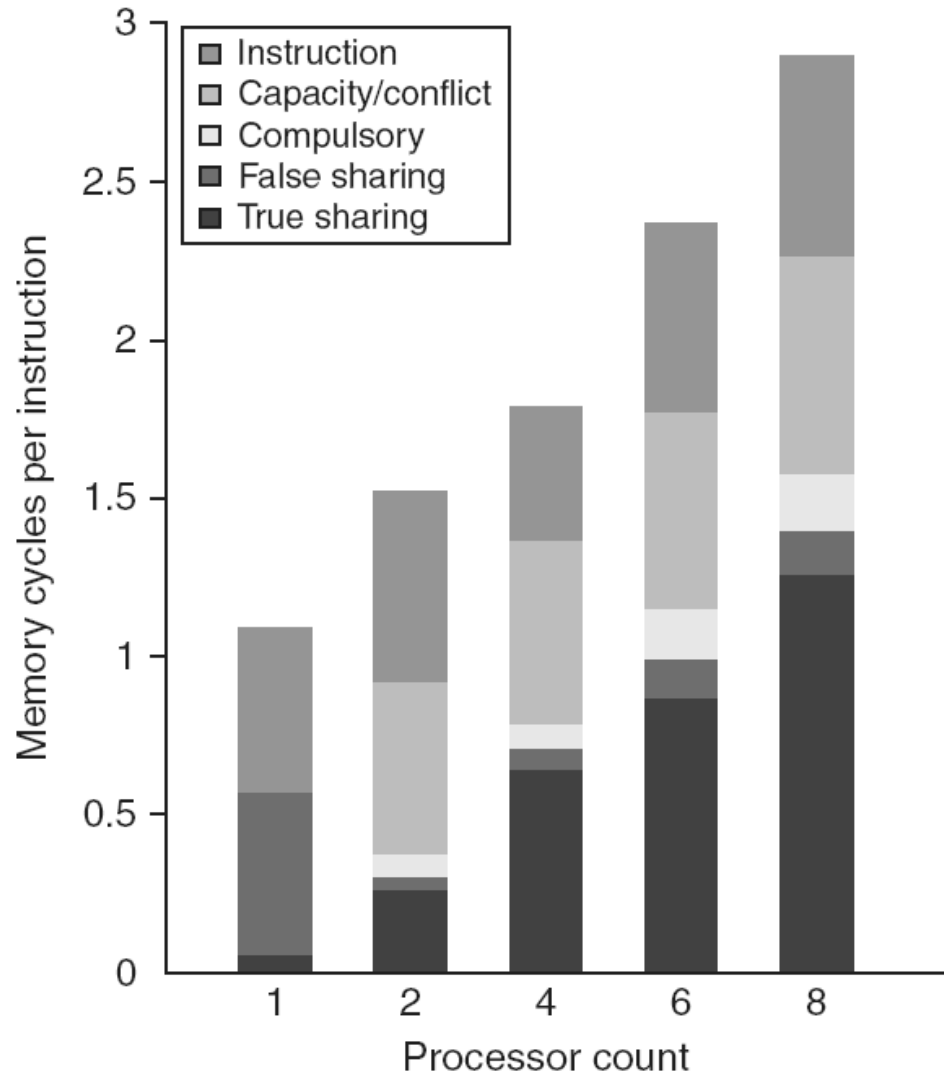
# Performance Study: Commercial Workload



# Performance Study: Commercial Workload

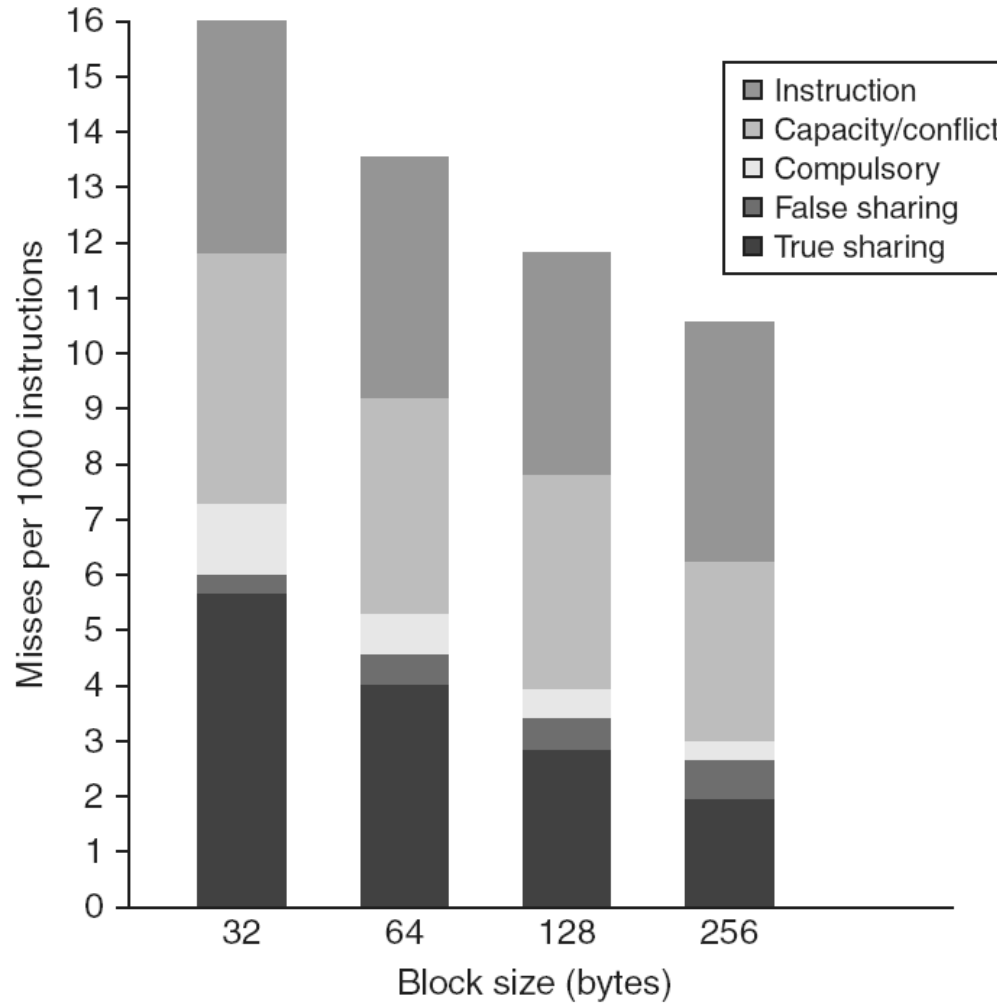


# Performance Study: Commercial Workload





# Performance Study: Commercial Workload



# Synchronization

- Basic building blocks:
  - Atomic exchange
    - Swaps register with memory location
  - Test-and-set
    - Sets under condition
  - Fetch-and-increment
    - Reads original value from memory and increments it in memory
  - Requires memory read and write in uninterruptable instruction
- load linked/store conditional
  - If the contents of the memory location specified by the load linked are changed before the store conditional to the same address, the store conditional fails

# Implementing Locks

## → Spin lock

### → If no coherence:

```
DADDUI R2,R0,#1
```

```
lockit:    EXCH      R2,0(R1);atomic exchange
           BNEZ      R2,lockit;already locked?
```

### → If coherence:

```
lockit:    LD        R2,0(R1);load of lock
           BNEZ      R2,lockit;not available-spin
           DADDUI R2,R0,#1      ;load locked value
           EXCH      R2,0(R1);swap
           BNEZ      R2,lockit;branch if lock wasn't 0
```

# Implementing Locks

→ Advantage of this scheme: reduces memory traffic

Step	P0	P1	P2	Coherence state of lock at end of step	Bus/directory activity
1	Has lock	Begins spin, testing if lock = 0	Begins spin, testing if lock = 0	Shared	Cache misses for P1 and P2 satisfied in either order. Lock state becomes shared.
2	Set lock to 0	(Invalidate received)	(Invalidate received)	Exclusive (P0)	Write invalidate of lock variable from P0.
3		Cache miss	Cache miss	Shared	Bus/directory services P2 cache miss; write-back from P0; state shared.
4		(Waits while bus/directory busy)	Lock = 0 test succeeds	Shared	Cache miss for P2 satisfied
5		Lock = 0	Executes swap, gets cache miss	Shared	Cache miss for P1 satisfied
6		Executes swap, gets cache miss	Completes swap: returns 0 and sets lock = 1	Exclusive (P2)	Bus/directory services P2 cache miss; generates invalidate; lock is exclusive.
7		Swap completes and returns 1, and sets lock = 1	Enter critical section	Exclusive (P1)	Bus/directory services P1 cache miss; sends invalidate and generates write-back from P2.
8		Spins, testing if lock = 0			None

# Implementing Locks

- To implement, delay completion of all memory accesses until all invalidations caused by the access are completed
  - Reduces performance!
- Alternatives:
  - Program-enforced synchronization to force write on processor to occur before read on the other processor
    - Requires synchronization object for A and another for B
      - “Unlock” after write
      - “Lock” after read

# Relaxed Consistency Models

## → Rules:

→  $X \rightarrow Y$

→ Operation X must complete before operation Y is done

→ Sequential consistency requires:

→  $R \rightarrow W, R \rightarrow R, W \rightarrow R, W \rightarrow W$

→ Relax  $W \rightarrow R$

→ “Total store ordering”

→ Relax  $W \rightarrow W$

→ “Partial store order”

→ Relax  $R \rightarrow W$  and  $R \rightarrow R$

# Relaxed Consistency Models

- Consistency model is multiprocessor specific
- Programmers will often implement explicit synchronization
- Speculation gives much of the performance advantage of relaxed models with sequential consistency
  - Basic idea: if an invalidation arrives for a result that has not been committed, use speculation recovery