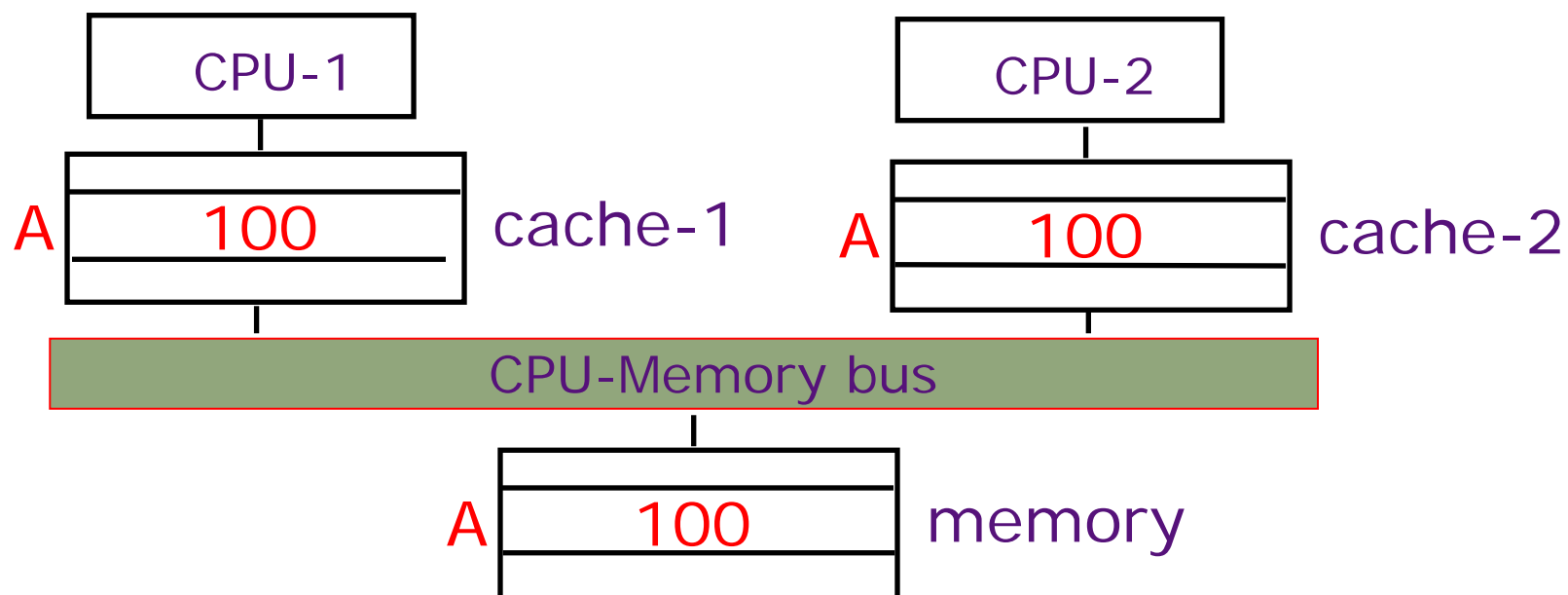# Sequential Consistency and Cache Coherence Protocols

*Arvind*
Computer Science and Artificial Intelligence Lab
M.I.T.

*Based on the material prepared by
Arvind and Krste Asanovic*

# Memory Consistency in SMPs


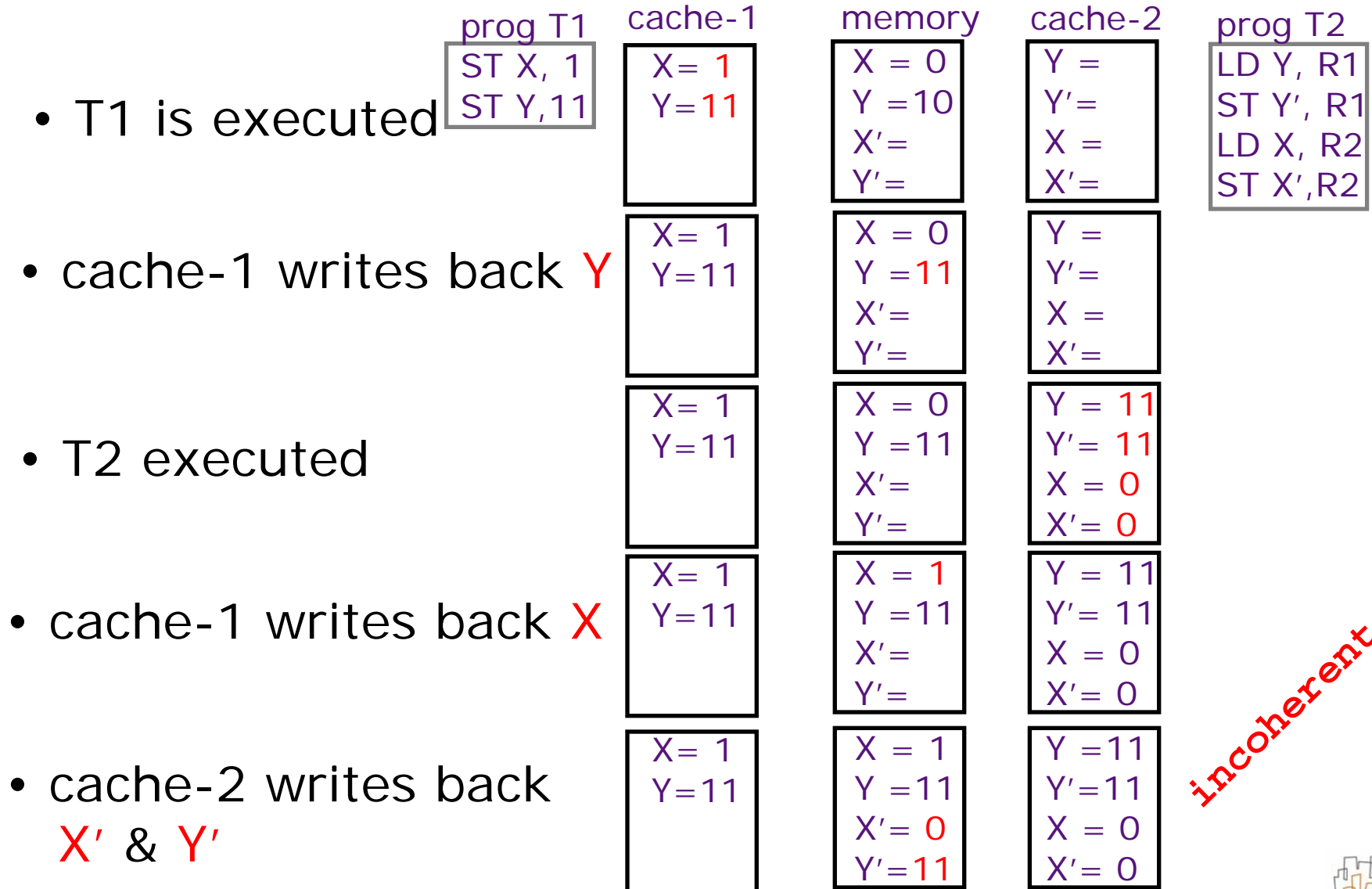
Suppose CPU-1 updates A to 200.

*write-back:*   memory and cache-2 have stale values
*write-through:*   cache-2 has a stale value

*Do these stale values matter?*
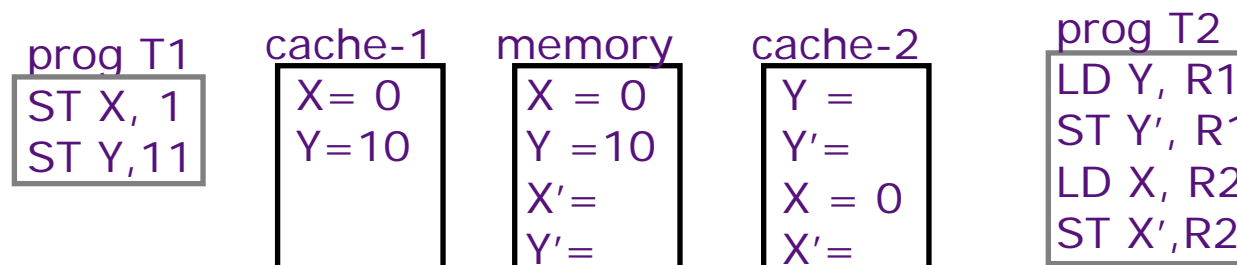*What is the view of shared memory for programming?*
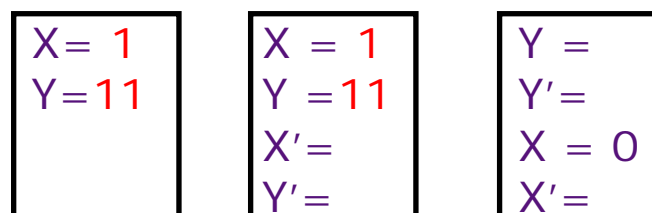
# Write-back Caches & SC

| | prog T1 | cache-1 | memory | cache-2 | prog T2 |
|---|---|---|---|---|---|
| • T1 is executed | ST X, 1<br>ST Y,11 | X= 1<br>Y=11 | X = 0<br>Y =10<br>X'=<br>Y'= | Y =<br>Y'=<br>X =<br>X'= | LD Y, R1<br>ST Y', R1<br>LD X, R2<br>ST X',R2 |
| • cache-1 writes back Y | | X= 1<br>Y=11 | X = 0<br>Y =11<br>X'=<br>Y'= | Y =<br>Y'=<br>X =<br>X'= | |
| • T2 executed | | X= 1<br>Y=11 | X = 0<br>Y =11<br>X'=<br>Y'= | Y = 11<br>Y'= 11<br>X = 0<br>X'= 0 | |
| • cache-1 writes back X | | X= 1<br>Y=11 | X = 1<br>Y =11<br>X'=<br>Y'= | Y = 11<br>Y'= 11<br>X = 0<br>X'= 0 | |
| • cache-2 writes back<br>  X' & Y' | | X= 1<br>Y=11 | X = 1<br>Y =11<br>X'= 0<br>Y'=11 | Y =11<br>Y'=11<br>X = 0<br>X'= 0 | |

*incoherent*

# Write-through Caches & SC

| prog T1 | cache-1 | memory | cache-2 | prog T2 |
|---------|---------|--------|---------|---------|
| ST X, 1 | X= 0 | X = 0 | Y = | LD Y, R1 |
| ST Y,11 | Y=10 | Y =10 | Y'= | ST Y', R1 |
| | | X'= | X = 0 | LD X, R2 |
| | | Y'= | X'= | ST X',R2 |

- T1 executed

| cache-1 | memory | cache-2 |
|---------|--------|---------|
| X= 1 | X = 1 | Y = |
| Y=11 | Y =11 | Y'= |
| | X'= | X = 0 |
| | Y'= | X'= |

- T2 executed

| cache-1 | memory | cache-2 |
|---------|--------|---------|
| X= 1 | X = 1 | Y = 11 |
| Y=11 | Y =11 | Y'= 11 |
| | X'= 0 | X = 0 |
| | Y'=11 | X'= 0 |

*Write-through caches don't preserve sequential consistency either*

CSAIL

# Maintaining Sequential Consistency

SC is sufficient for correct producer-consumer and mutual exclusion code (e.g., Dekker)

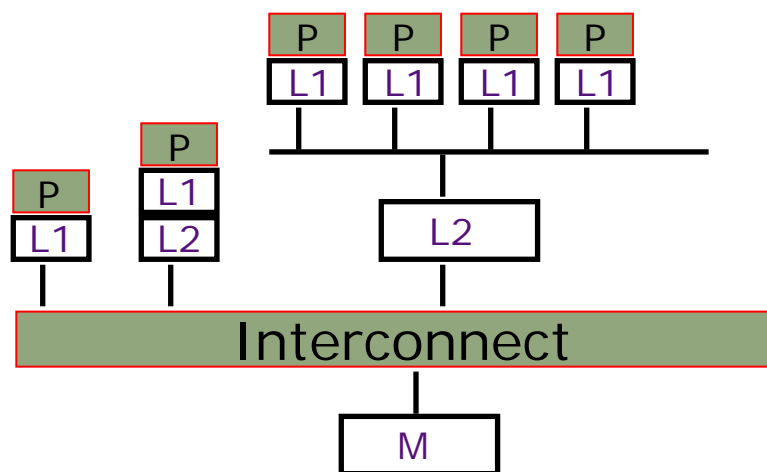Multiple copies of a location in various caches can cause SC to break down.

Hardware support is required such that
- only one processor at a time has write permission for a location
- no processor can load a stale copy of the location after a write

$$\Rightarrow \quad \textit{cache coherence protocols}$$

# A System with Multiple Caches



- Modern systems often have hierarchical caches
- Each cache has exactly one parent but can have zero or more children
- Only a parent and its children can communicate directly
- *Inclusion property* is maintained between a parent and its children, i.e.,

$$a \in L_i \quad \Rightarrow \quad a \in L_{i+1}$$
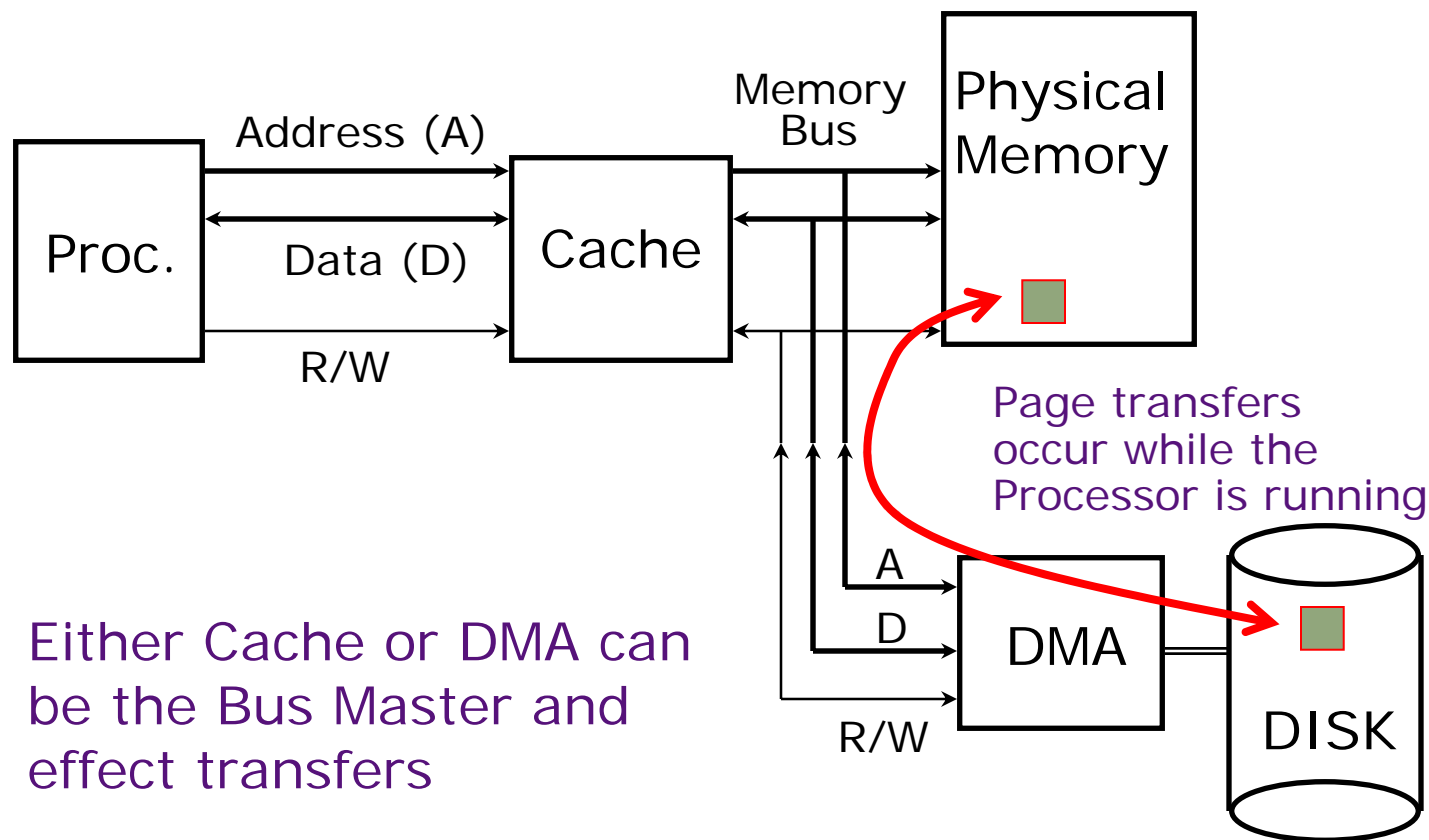
# Cache Coherence Protocols for SC

*write request:*
> the address is *invalidated* (*updated*) in all other caches *before* (*after*) the write is performed

*read request:*
> if a dirty copy is found in some cache, a write-back is performed before the memory is read

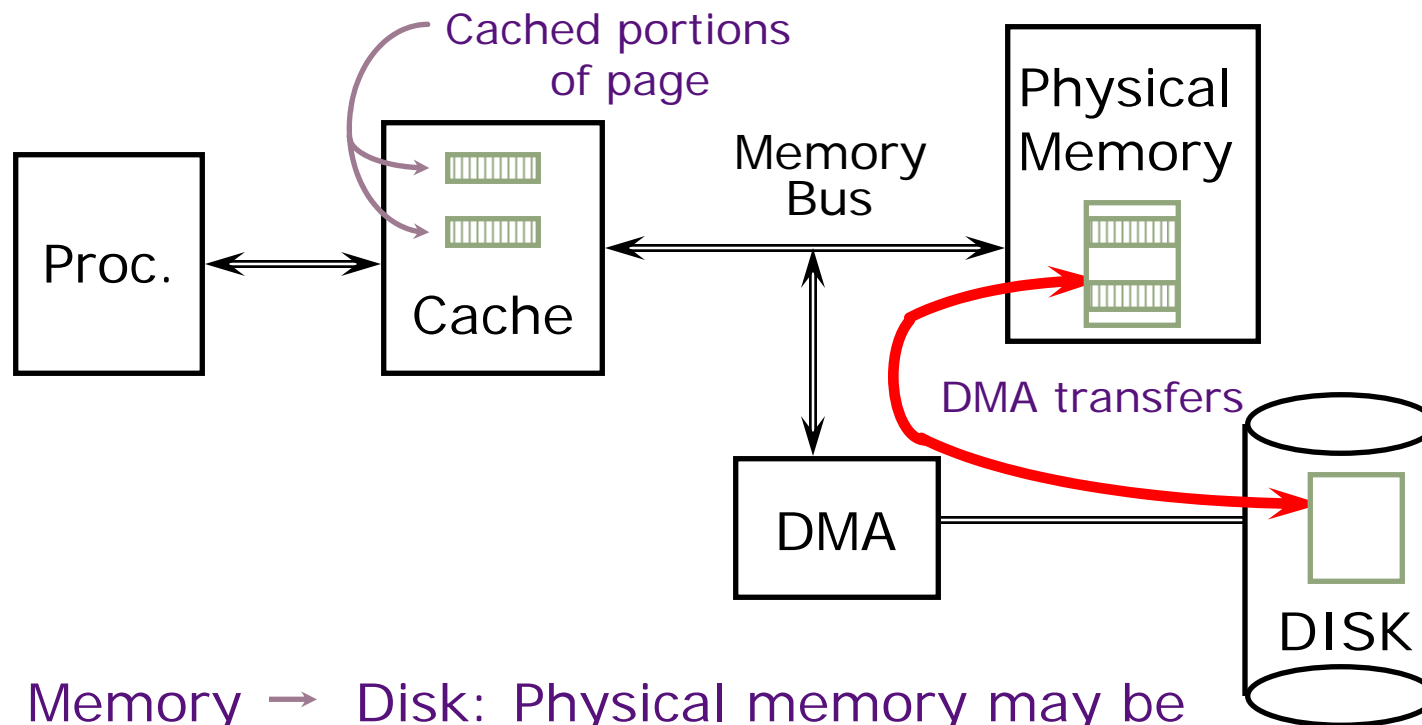*We will focus on Invalidation protocols as opposed to Update protocols*

# Warmup: Parallel I/O

Memory
Bus

Physical
Memory

Address (A)

Proc.

Data (D)

Cache

R/W

A

D

DMA

R/W

DISK

Page transfers
occur while the
Processor is running

Either Cache or DMA can
be the Bus Master and
effect transfers

DMA stands for Direct Memory Access

# Problems with Parallel I/O



Cached portions
of page

Proc.

Cache

Memory
Bus
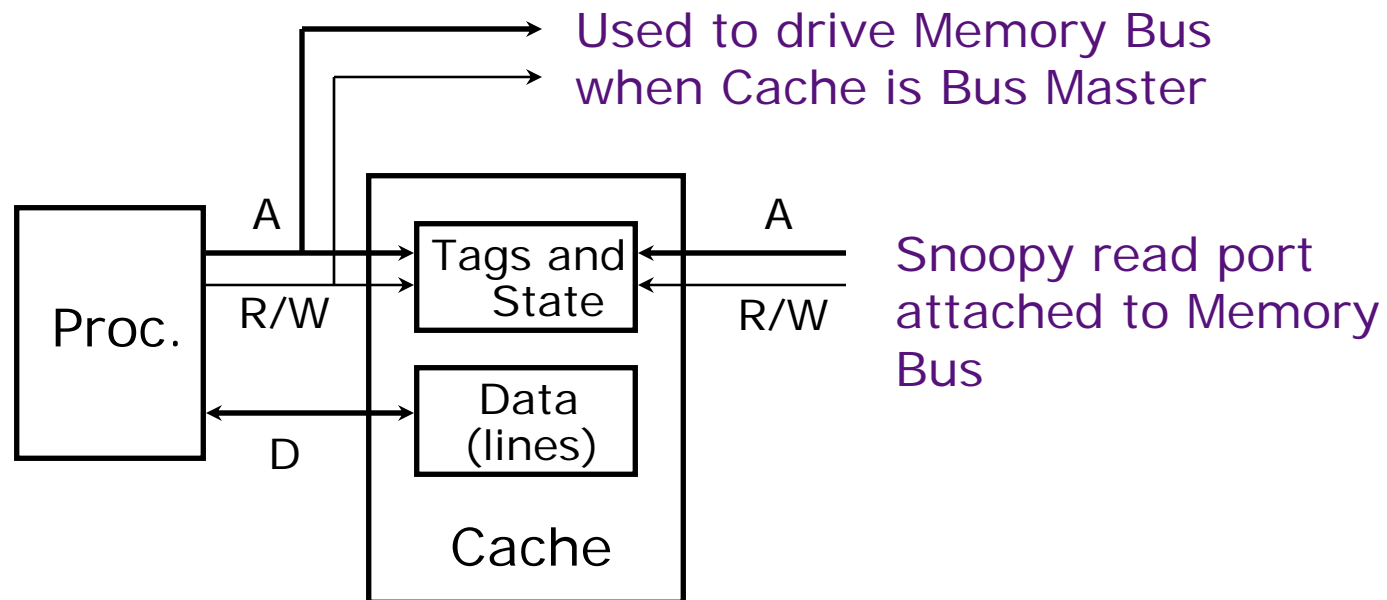
Physical
Memory

DMA

DMA transfers

DISK

Memory → Disk: Physical memory may be
stale if Cache copy is dirty

Disk → Memory: Cache may have data
corresponding to the memory

# Snoopy Cache *Goodman 1983*

- Idea: Have cache watch (or snoop upon) DMA transfers, and then "do the right thing"
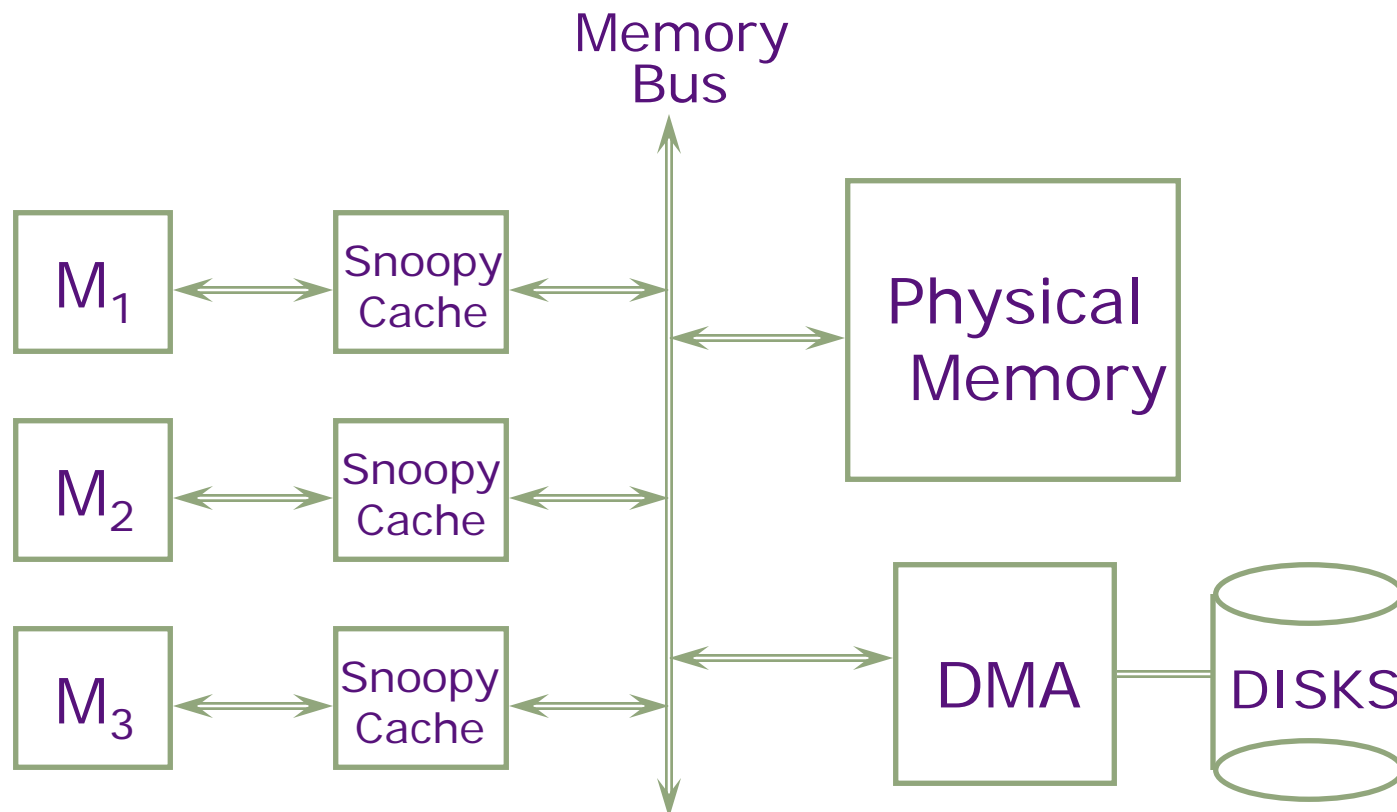- Snoopy cache tags are dual-ported

Used to drive Memory Bus
when Cache is Bus Master

Snoopy read port
attached to Memory
Bus

Proc.

A

R/W

D

Tags and State

Data (lines)

A

R/W

Cache

# Snoopy Cache Actions

| Observed Bus Cycle | Cache State | Cache Action |
|---|---|---|
| | Address not cached | No action |
| Read Cycle | Cached, unmodified | No action |
| Memory → Disk | Cached, modified | Cache intervenes |
| | Address not cached | No action |
| Write Cycle | Cached, unmodified | Cache purges its copy |
| Disk → Memory | Cached, modified | ??? |

November 9, 2005

# Shared Memory Multiprocessor

Memory
Bus

| | | |
|---|---|---|
| $M_1$ | Snoopy Cache | |
| | | Physical Memory |
| $M_2$ | Snoopy Cache | |
| | | |
| $M_3$ | Snoopy Cache | DMA — DISKS |

Use snoopy mechanism to keep all
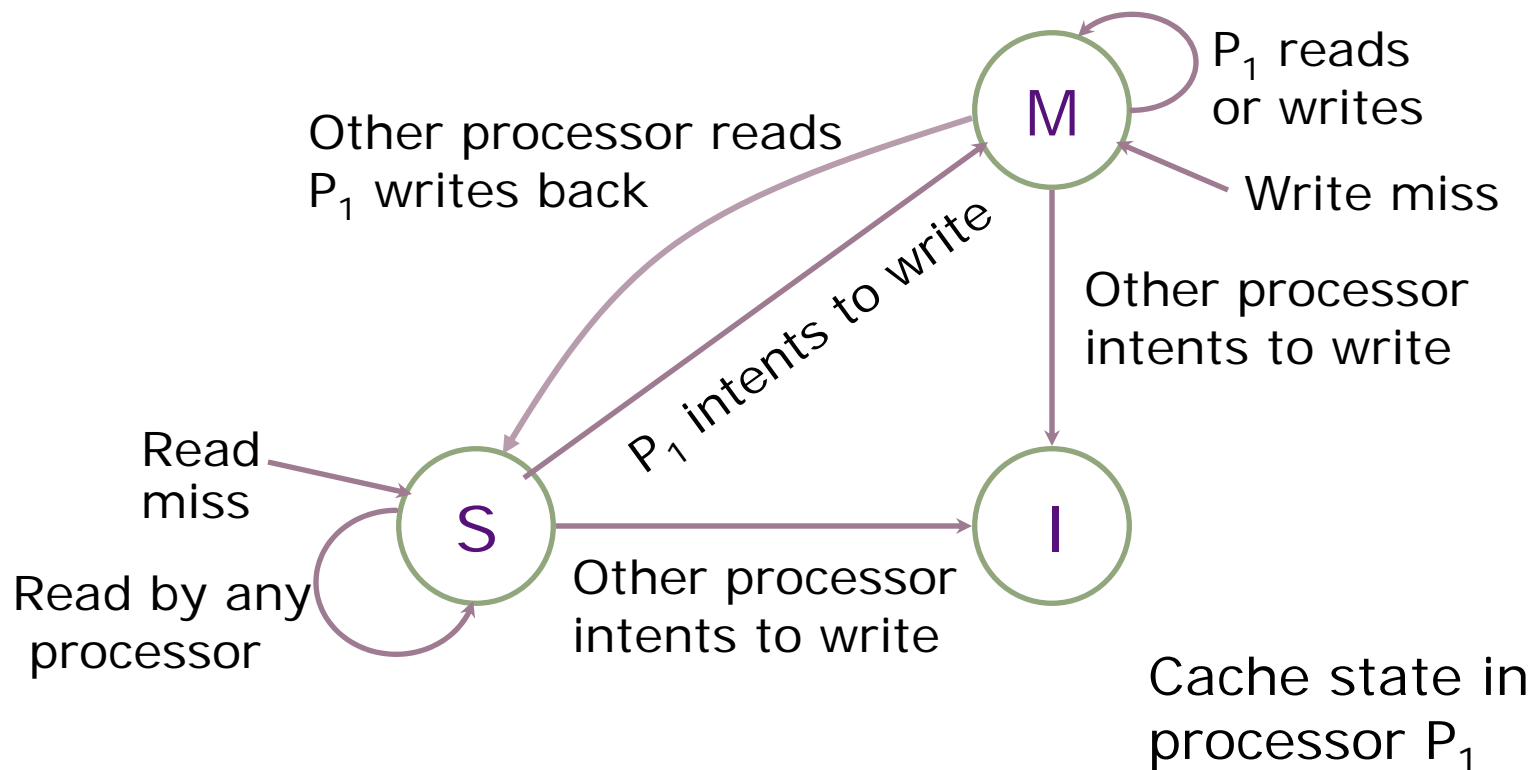processors' view of memory coherent

CSAIL

# Cache State Transition Diagram
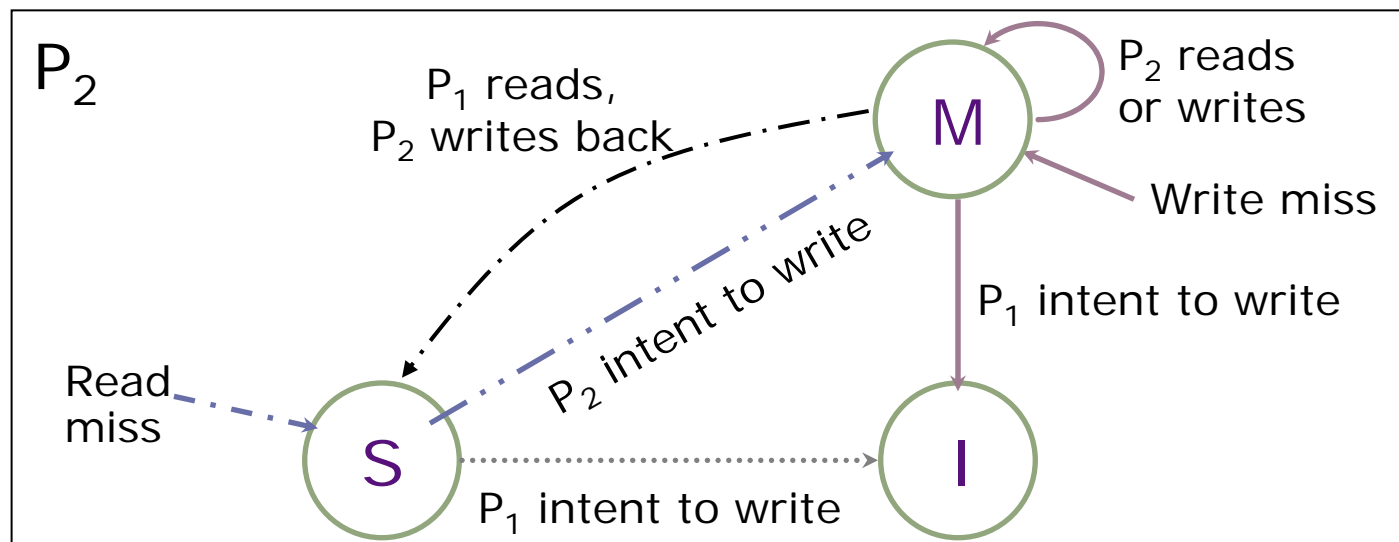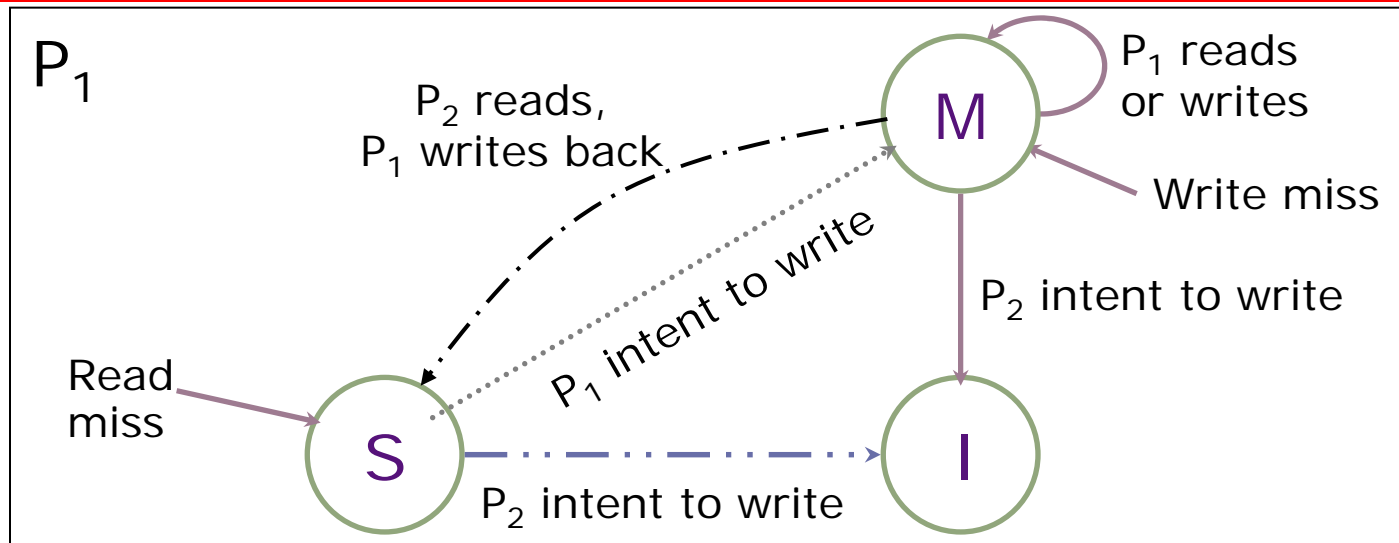## The MSI protocol

*Each* cache line has a tag

M: Modified
S: Shared
I: Invalid

| state bits | | Address tag |
|---|---|---|

Other processor reads
$P_1$ writes back

$P_1$ intents to write

M

$P_1$ reads or writes

Write miss

Other processor intents to write

Read miss

Read by any processor

S

Other processor intents to write

I

Cache state in processor $P_1$

CSAIL

# 2 Processor Example

P$_1$ reads
P$_1$ writes
P$_2$ reads
P$_2$ writes
P$_1$ reads
P$_1$ writes
P$_2$ writes
P$_1$ writes

P$_1$

P$_2$ reads,
P$_1$ writes back

M

P$_1$ reads
or writes

Write miss

P$_1$ intent to write

P$_2$ intent to write

Read miss

S

I

P$_2$ intent to write

P$_2$

P$_1$ reads,
P$_2$ writes back

M

P$_2$ reads
or writes

Write miss

P$_2$ intent to write

P$_1$ intent to write

Read miss

S

I

P$_1$ intent to write

CSAIL

# Observation

Other processor reads
$P_1$ writes back

M

$P_1$ reads
or writes

Write miss

$P_1$ intents to write

Other processor
intents to write

Read
miss

S
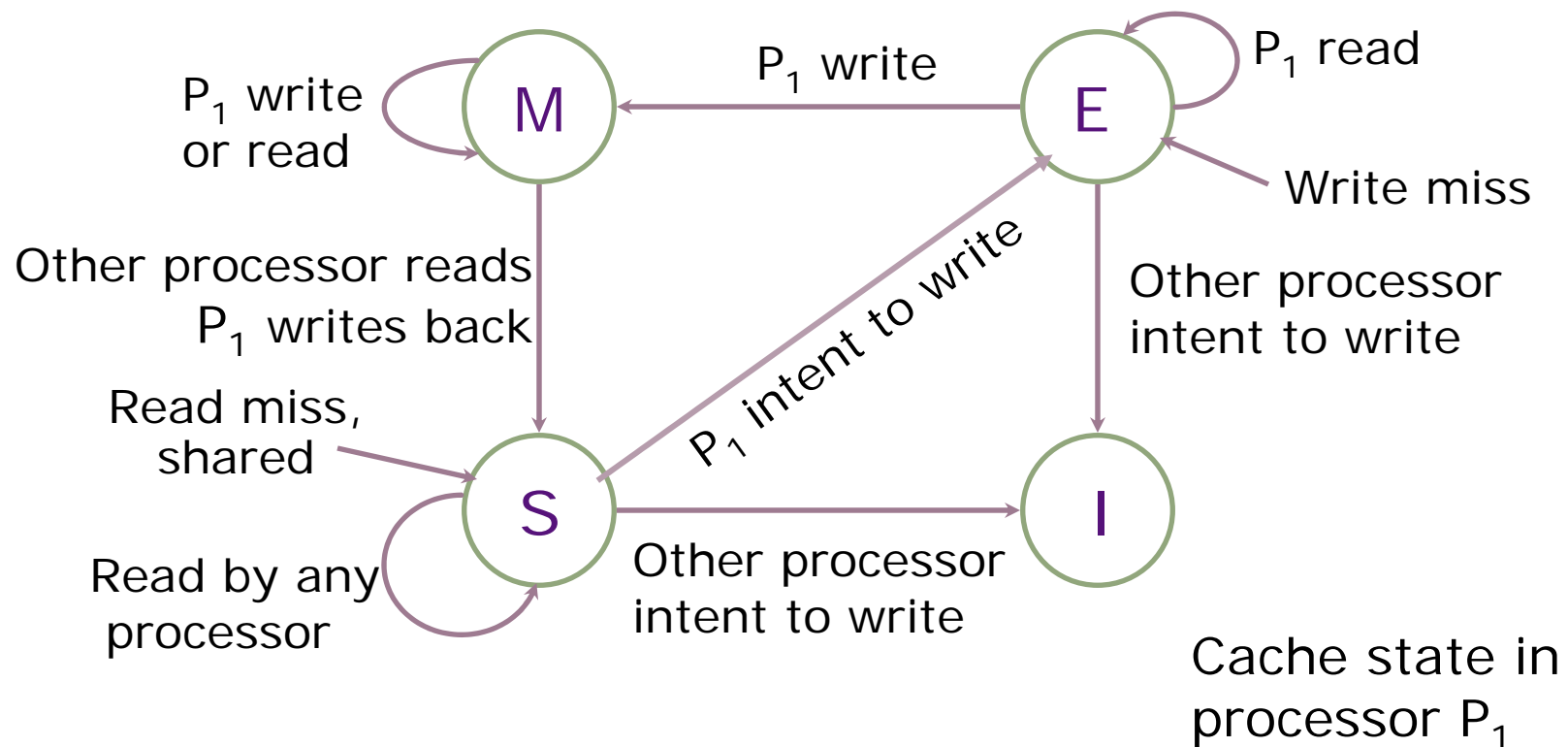
I

Read by any
processor

Other processor
intents to write

- If a line is in the M state then no other cache can have a copy of the line!
  - Memory stays coherent, multiple differing copies cannot exist

CSAIL

# MESI: An Enhanced MSI protocol

*Each* cache line has a tag

| state bits | | Address tag |
|---|---|---|

M: Modified Exclusive
E: Exclusive, unmodified
S: Shared
I: Invalid

$P_1$ write or read → M

M → E : $P_1$ write

E → E : $P_1$ read

Other processor reads $P_1$ writes back

M → S

Read miss, shared → S

Read by any processor (S self-loop)

S → E : $P_1$ intent to write

S → I : Other processor intent to write

E → I : Other processor intent to write

Write miss → E

Cache state in processor $P_1$

CSAIL

# Five-minute break to stretch your legs

# Cache Coherence State Encoding

block Address

| tag | | index$_m$ | offset |
|-----|--|-----------|--------|

| tag | V | M | data block |
|-----|---|---|------------|

=

*Valid and dirty bits can be used
to encode S, I, and (E, M) states*

V=0, D=x ⇒ Invalid
V=1, D=0 ⇒ Shared *(not dirty)*
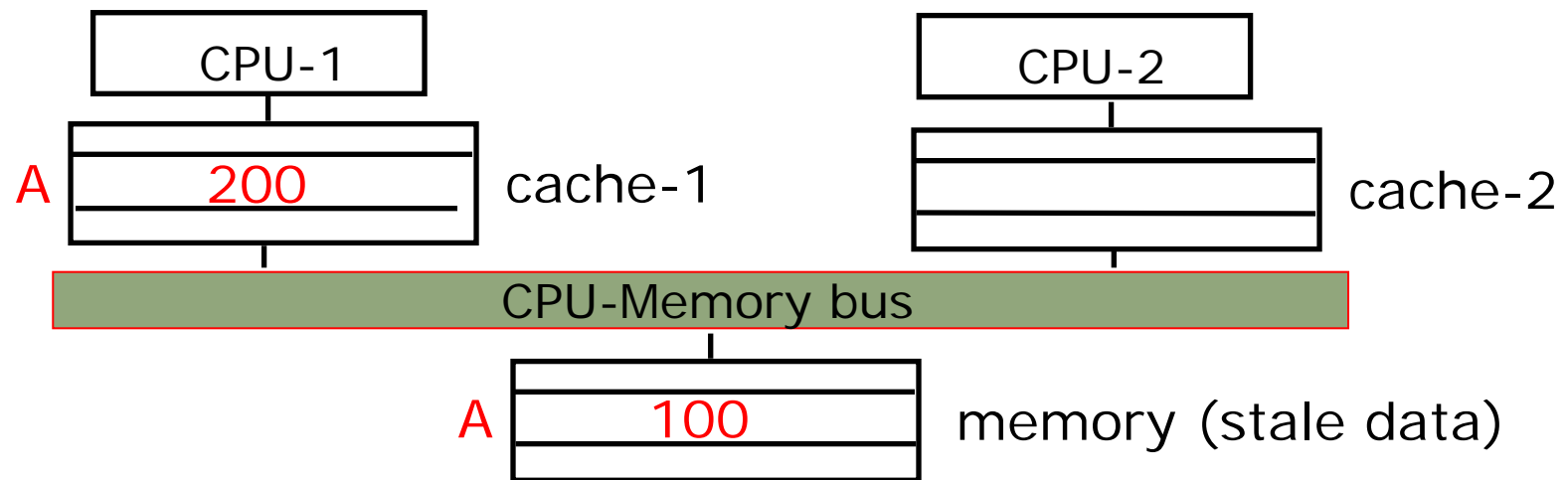V=1, D=1 ⇒ Exclusive *(dirty)*

Hit?          word

# 2-Level Caches



- Processors often have two-level caches
  - Small L1 on chip, large L2 off chip
- *Inclusion property:* entries in L1 must be in L2
    invalidation in L2 $\Rightarrow$ invalidation in L1
- Snooping on L2 does not affect CPU-L1 bandwidth

*What problem could occur?*

# Intervention



When a read-miss for A occurs in cache-2,
a read request for A is placed on the bus

- Cache-1 needs to supply & change its state to shared
- The memory may respond to the request also!

*Does memory know it has stale data?*

Cache-1 needs to intervene through memory
controller to supply correct data to cache-2

# False Sharing

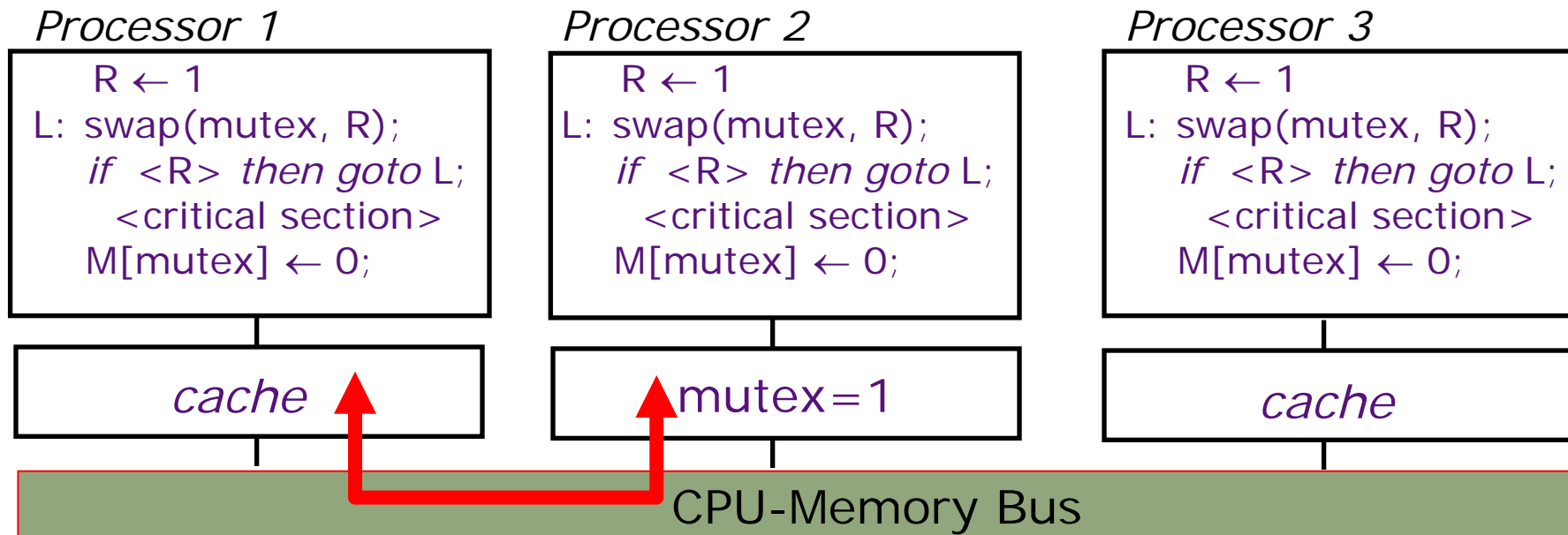| state | blk addr | data0 | data1 | ... | dataN |
|-------|----------|-------|-------|-----|-------|

A cache block contains more than one word

Cache-coherence is done at the block-level and not word-level

Suppose $M_1$ writes $word_i$ and $M_2$ writes $word_k$ and both words have the same block address.

*What can happen?*

CSAIL

# Synchronization and Caches:
## *Performance Issues*

| Processor 1 | Processor 2 | Processor 3 |
|---|---|---|

Processor 1:
```
      R ← 1
L: swap(mutex, R);
   if <R> then goto L;
      <critical section>
   M[mutex] ← 0;
```

Processor 2:
```
      R ← 1
L: swap(mutex, R);
   if <R> then goto L;
      <critical section>
   M[mutex] ← 0;
```

Processor 3:
```
      R ← 1
L: swap(mutex, R);
   if <R> then goto L;
      <critical section>
   M[mutex] ← 0;
```

*cache*   |   mutex=1   |   *cache*

**CPU-Memory Bus**

Cache-coherence protocols will cause mutex to *ping-pong* between P1's and P2's caches.

Ping-ponging can be reduced by first reading the mutex location *(non-atomically)* and executing a swap only if it is found to be zero.

CSAIL

# Performance Related to Bus occupancy

In general, a *read-modify-write* instruction requires two memory (bus) operations without intervening memory operations by other processors

In a multiprocessor setting, bus needs to be locked for the entire duration of the atomic read and write operation

$\Rightarrow$ expensive for simple buses
$\Rightarrow$ *very expensive* for split-transaction buses

modern processors use

*load-reserve*
*store-conditional*

# Load-reserve & Store-conditional

Special register(s) to hold reservation flag and
address, and the outcome of store-conditional

Load-reserve(R, a):
   <flag, adr> ← <1, a>;
   R ← M[a];

Store-conditional(a, R):
   *if* <flag, adr> == <1, a>
   *then* cancel other procs'
      reservation on a;
      M[a] ← <R>;
      status ← succeed;
   *else* status ← fail;

If the snooper sees a store transaction to the address
in the reserve register, the reserve bit is set to 0
   • Several processors may reserve 'a' simultaneously
   • These instructions are like ordinary loads and stores
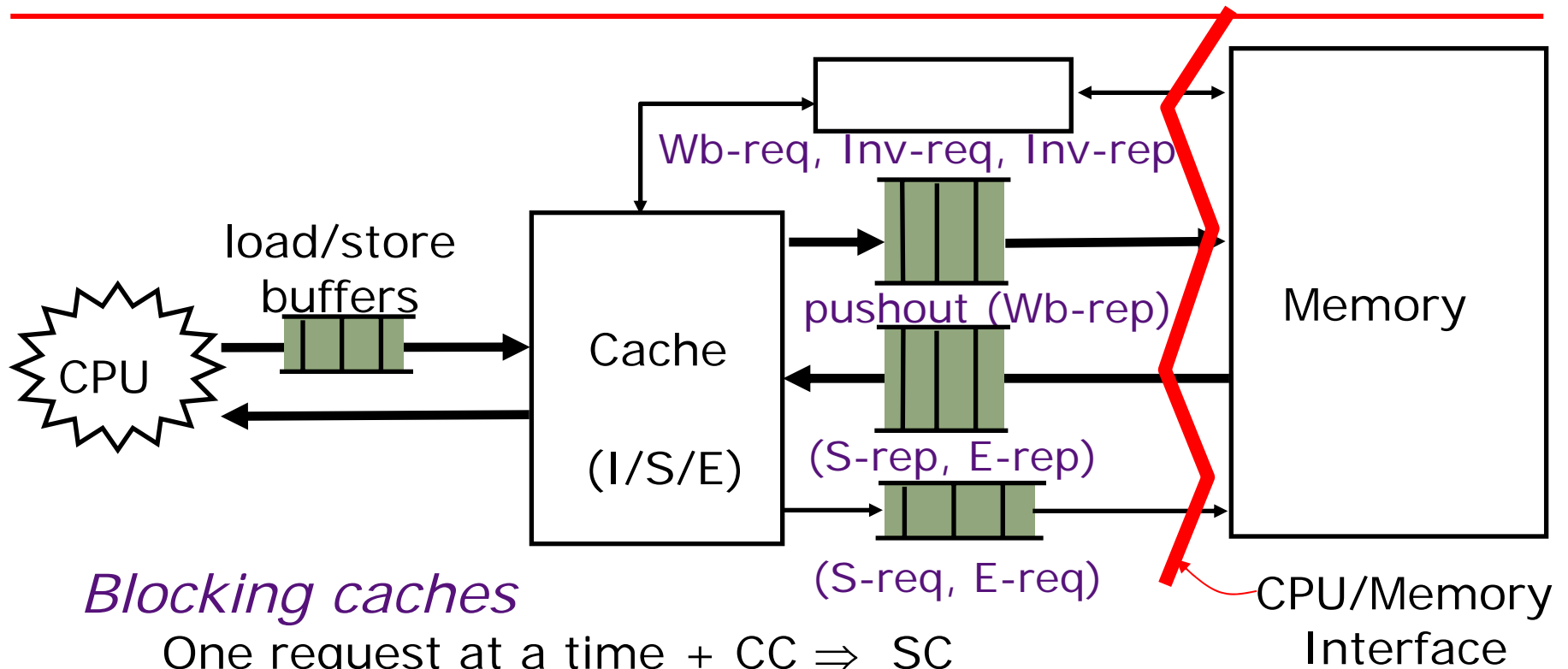      with respect to the bus traffic

# Performance:
## *Load-reserve & Store-conditional*

The total number of memory (bus) transactions is not necessarily reduced, but splitting an atomic instruction into load-reserve & store-conditional:

- *increases bus utilization* (and reduces processor stall time), especially in split-transaction  buses

- *reduces cache ping-pong effect* because processors trying to acquire a semaphore do not have to perform a store each time

# Out-of-Order Loads/Stores & CC

Wb-req, Inv-req, Inv-rep

load/store buffers

CPU

Cache

(I/S/E)

pushout (Wb-rep)

Memory

(S-rep, E-rep)

(S-req, E-req)

CPU/Memory Interface

*Blocking caches*
  One request at a time + CC $\Rightarrow$ SC

*Non-blocking caches*
  Multiple requests (different addresses) concurrently + CC
      $\Rightarrow$ Relaxed memory models

CC ensures that all processors observe the same order of loads and stores to an address

*next time*

# Designing a Cache Coherence Protocol

CSAIL

*Thank you !*

# 2 Processor Example