# CS 152 Computer Architecture and Engineering

# Lecture 21: Directory-Based Cache Protocols

Krste Asanovic

Electrical Engineering and Computer Sciences
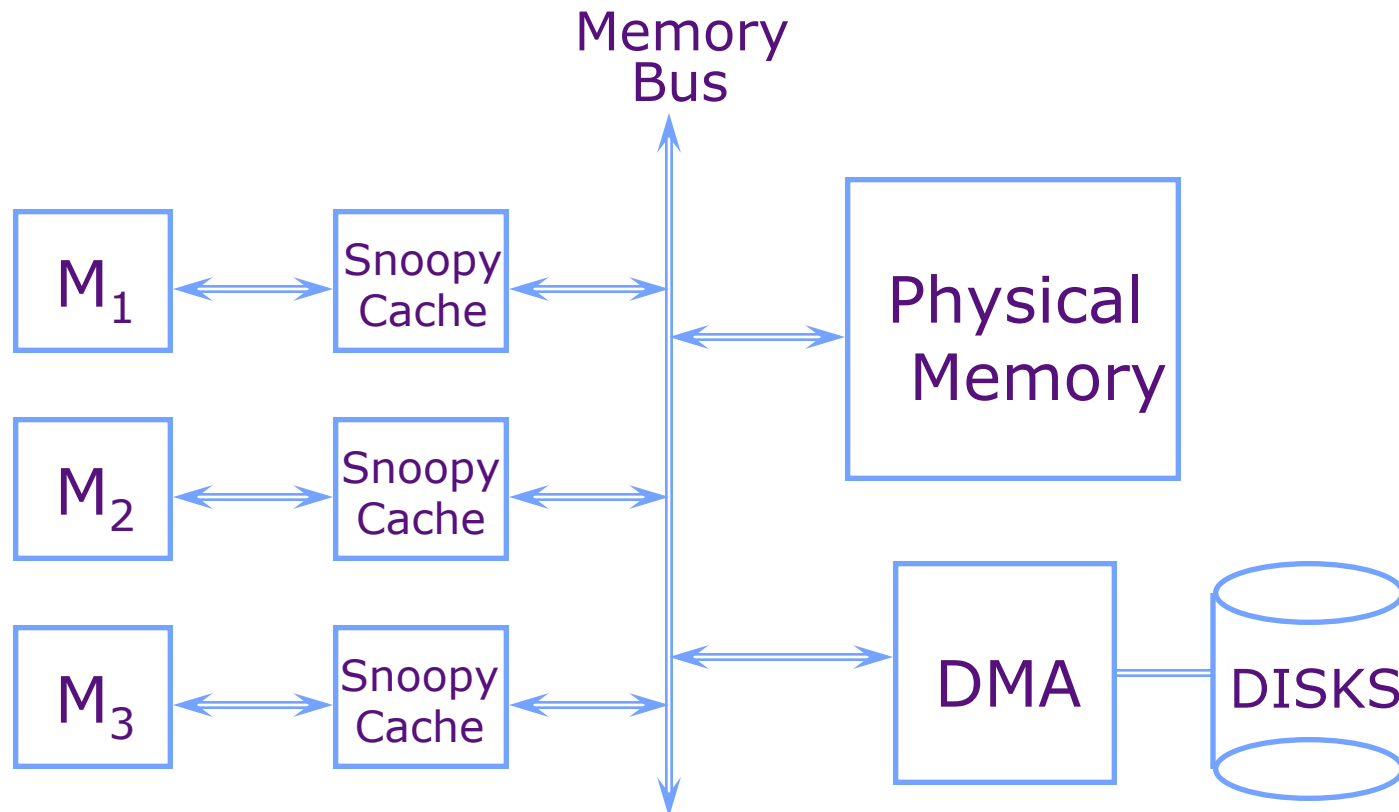University of California, Berkeley

**http://www.eecs.berkeley.edu/~krste**
**http://inst.cs.berkeley.edu/~cs152**

# Recap: Snoopy Cache Protocols

Memory
Bus

| $M_1$ | Snoopy Cache |
| --- | --- |

| $M_2$ | Snoopy Cache |
| --- | --- |

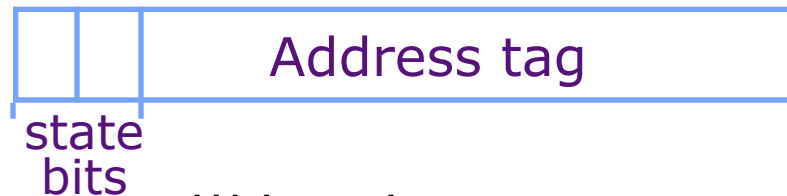| $M_3$ | Snoopy Cache |
| --- | --- |

Physical
Memory

DMA    DISKS

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

# MESI: An Enhanced MSI protocol
## increased performance for private data

*Each* cache line has a tag

| | | Address tag |
|---|---|---|

state bits

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



Write miss

$P_1$ write or read

M

$P_1$ write

$P_1$ read

E

Read miss, not shared

Other processor reads

$P_1$ intent to write

Other processor reads

Other processor intent to write

Other processor reads
$P_1$ writes back

Other processor intent to write, P1 writes back

Read miss, shared

S

I

Read by any processor

Other processor intent to write

Cache state in processor $P_1$

# Performance of Symmetric Shared-Memory Multiprocessors

Cache performance is combination of:

1. Uniprocessor cache miss traffic

2. Traffic caused by communication

   – Results in invalidations and subsequent cache misses

- Adds 4th C: *coherence miss*

   – Joins Compulsory, Capacity, Conflict

   – (Sometimes called a *Communication* miss)

# Coherency Misses

1. **True sharing misses** arise from the communication of data through the cache coherence mechanism
   - Invalidates due to 1st write to shared block
   - Reads by another CPU of modified block in different cache
   - Miss would still occur if block size were 1 word

2. **False sharing misses** when a block is invalidated because some word in the block, other than the one being read, is written into
   - Invalidation does not cause a new value to be communicated, but only causes an extra cache miss
   - Block is shared, but no word in block is actually shared
     $\Rightarrow$ miss would not occur if block size were 1 word
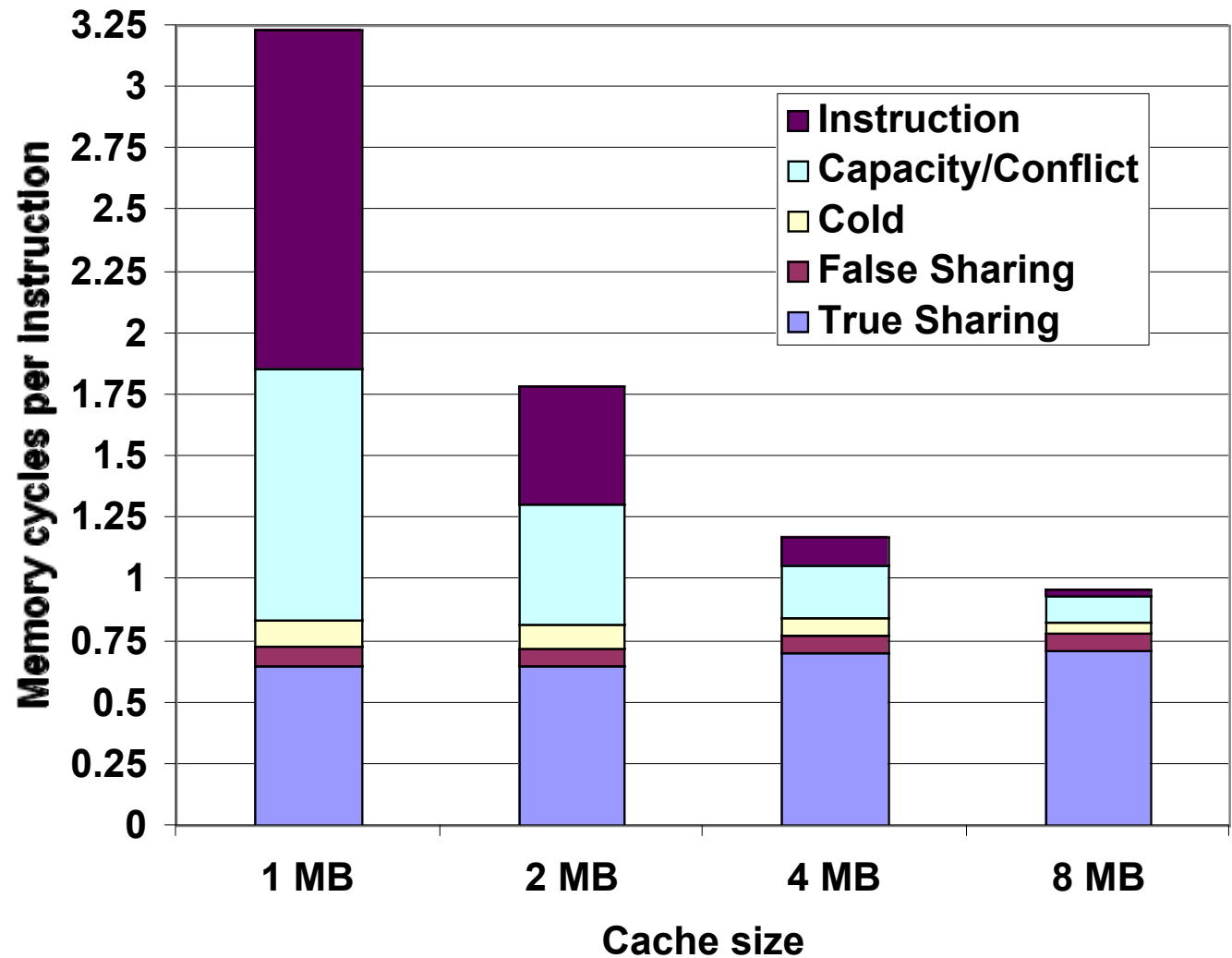
# Example: True v. False Sharing v. Hit?

- **Assume x1 and x2 in same cache block.**
  **P1 and P2 both read x1 and x2 before.**

| Time | P1 | P2 | True, False, Hit? Why? |
|------|----|----|------------------------|
| 1 | Write x1 | | True miss; invalidate x1 in P2 |
| 2 | | Read x2 | False miss; x1 irrelevant to P2 |
| 3 | Write x1 | | False miss; x1 irrelevant to P2 |
| 4 | | Write x2 | False miss; x1 irrelevant to P2 |
| 5 | Read x2 | | True miss; invalidate x2 in P1 |

# MP Performance 4 Processor Commercial Workload: OLTP, Decision Support (Database), Search Engine

• True sharing and false sharing unchanged going from 1 MB to 8 MB (L3 cache)

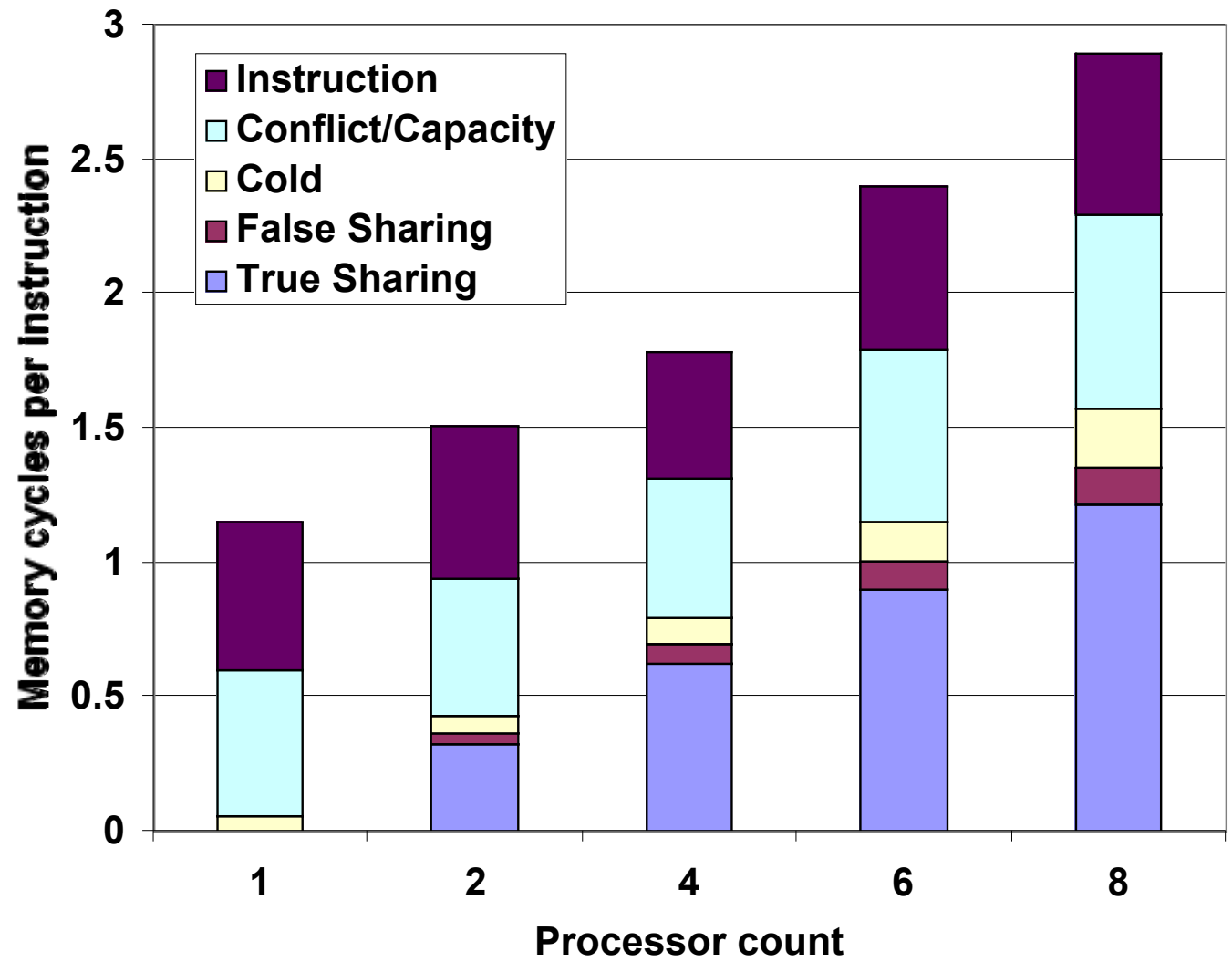• Uniprocessor cache misses improve with cache size increase (Instruction, Capacity/Conflict, Compulsory)

# MP Performance 2MB Cache Commercial Workload: OLTP, Decision Support (Database), Search Engine

- True sharing, false sharing increase going from 1 to 8 CPUs



Legend:
- Instruction
- Conflict/Capacity
- Cold
- False Sharing
- True Sharing

Y-axis: Memory cycles per instruction (0 to 3)
X-axis: Processor count (1, 2, 4, 6, 8)

# A Cache Coherent System Must:

- Provide set of states, state transition diagram, and actions

- Manage coherence protocol
  - (0) Determine when to invoke coherence protocol
  - (a) Find info about state of address in other caches to determine action
    - » whether need to communicate with other cached copies
  - (b) Locate the other copies
  - (c) Communicate with those copies  (invalidate/update)

- (0) is done the same way on all systems
  - state of the line is maintained in the cache
  - protocol is invoked if an "access fault" occurs on the line

- Different approaches distinguished by (a) to (c)

# 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 number of processors, 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
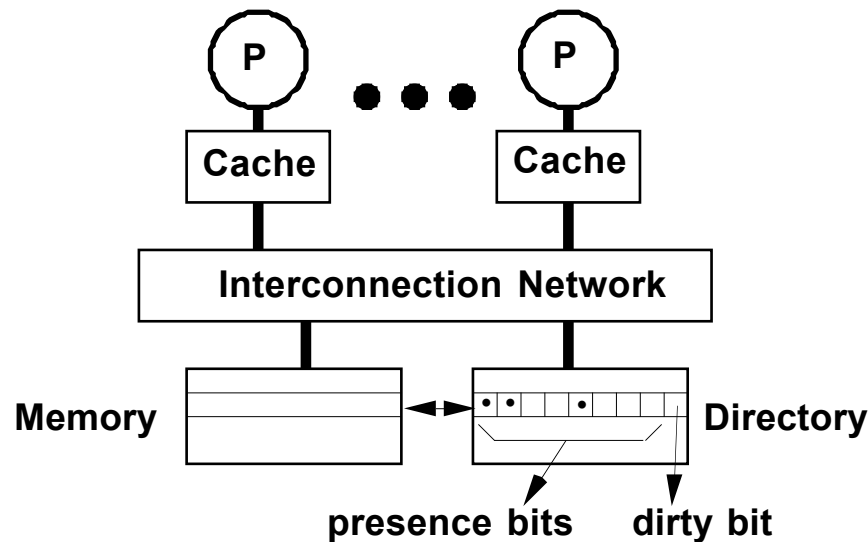    - different mechanisms to manage protocol

# 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 (downgrade 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:

  - If dirty-bit OFF then {send invalidations to all caches that have the block; turn dirty-bit ON; supply data to i; turn p[i] ON; ... }
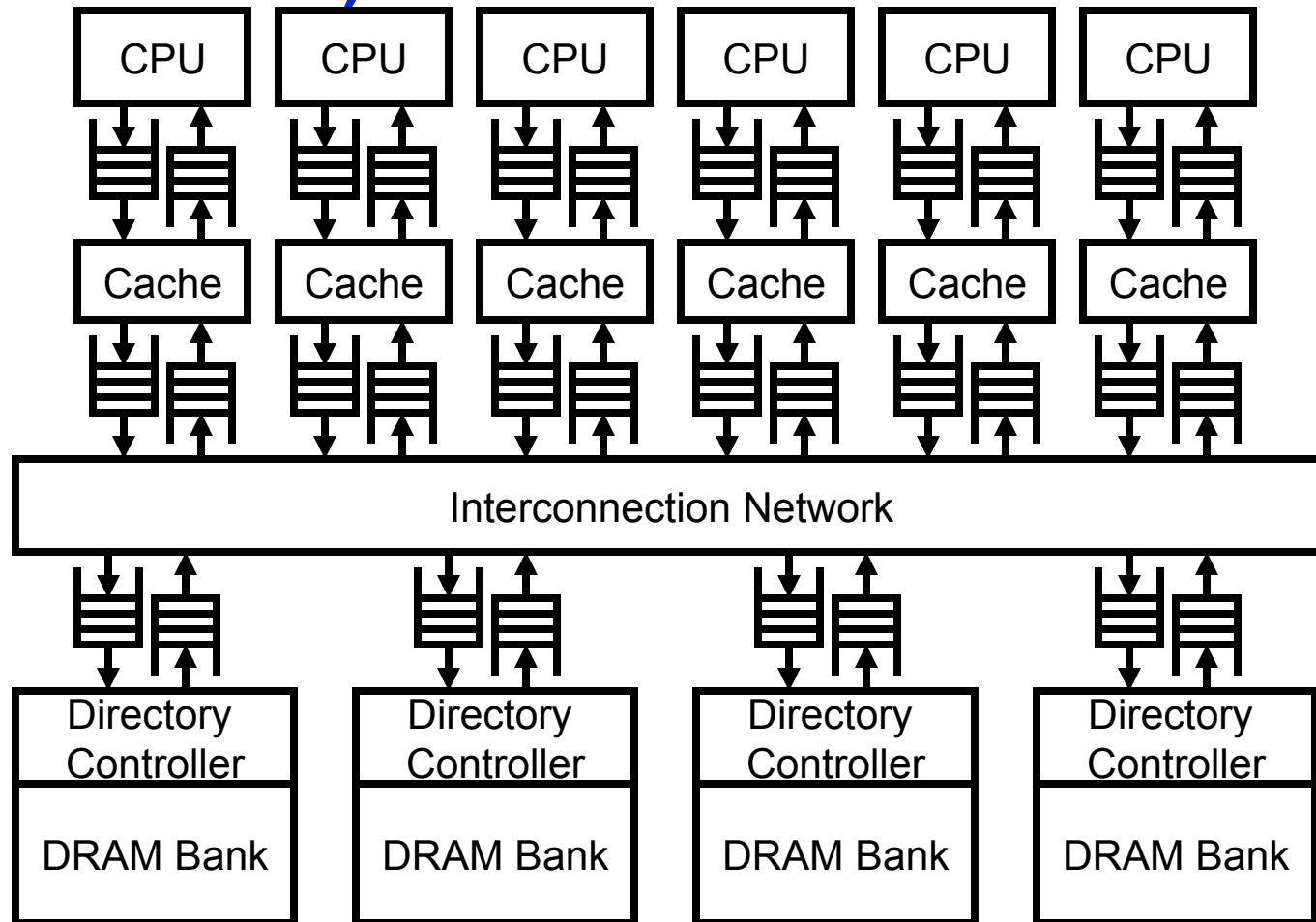
# CS152 Administrivia

- Final quiz, Thursday April 29
  - Multiprocessors, Memory models, Cache coherence
  - Lectures 19-21, PS 5, Lab 5

- Next lecture, "Virtual Machines", Thursday April 22

- Last lecture, "Putting it all Together", Tuesday April 27
  - Summary of the course
  - Case Study: Intel Nehalem
  - HKN Course Survey

# Directory Cache Protocol (Handout 6)



- Assumptions: Reliable network, FIFO message delivery between any given source-destination pair

# Cache States

For each cache line, there are 4 possible states:

- C-invalid (= Nothing): The accessed data is not resident in the cache.

- C-shared (= Sh): The accessed data is resident in the cache, and possibly also cached at other sites. The data in memory is valid.

- C-modified (= Ex): The accessed data is exclusively resident in this cache, and has been modified. Memory does not have the most up-to-date data.

- C-transient (= Pending): The accessed data is in a *transient* state (for example, the site has just issued a protocol request, but has not received the corresponding protocol reply).

# Home directory states

- ## For each memory block, there are 4 possible states:

  - R(dir): The memory block is shared by the sites specified in dir (dir is a set of sites). The data in memory is valid in this state.  If dir is empty (i.e., dir = ε), the memory block is not cached by any site.

  - W(id): The memory block is exclusively cached at site id, and has been modified at that site. Memory does not have the most up-to-date data.

  - TR(dir): The memory block is in a transient state waiting for the acknowledgements to the invalidation requests that the home site has issued.

  - TW(id): The memory block is in a transient state waiting for a block exclusively cached at site id (i.e., in C-modified state) to make the memory block at the home site up-to-date.

# Protocol Messages

There are 10 different protocol messages:

| Category | Messages |
|---|---|
| Cache to Memory Requests | ShReq, ExReq |
| Memory to Cache Requests | WbReq, InvReq, FlushReq |
| Cache to Memory Responses | WbRep(v), InvRep, FlushRep(v) |
| Memory to Cache Responses | ShRep(v), ExRep(v) |

# Cache State Transitions (from invalid state)

| No. | Current State | Handling Message | Next State | Dequeue Message? | Action |
|-----|---------------|------------------|------------|------------------|--------|
| 1 | C-nothing | Load | C-pending | No | ShReq(id,Home,a) |
| 2 | C-nothing | Store | C-pending | No | ExReq(id,Home,a) |
| 3 | C-nothing | WbReq(a) | C-nothing | Yes | None |
| 4 | C-nothing | FlushReq(a) | C-nothing | Yes | None |
| 5 | C-nothing | InvReq(a) | C-nothing | Yes | None |
| 6 | C-nothing | ShRep (a) | C-shared | Yes | updates cache with prefetch data |
| 7 | C-nothing | ExRep (a) | C-exclusive | Yes | updates cache with data |

# Cache State Transitions (from shared state)

| No. | Current State | Handling Message | Next State | Dequeue Message? | Action |
|-----|---------------|------------------|------------|------------------|--------|
| 8 | C-shared | Load | C-shared | Yes | Reads cache |
| 9 | C-shared | WbReq(a) | C-shared | Yes | None |
| 10 | C-shared | FlushReq(a) | C-nothing | Yes | InvRep(id, Home, a) |
| 11 | C-shared | InvReq(a) | C-nothing | Yes | InvRep(id, Home, a) |
| 12 | C-shared | ExRep(a) | C-exclusive | Yes | None |
| 13 | C-shared | (Voluntary Invalidate) | C-nothing | N/A | InvRep(id, Home, a) |

# Cache State Transitions (from exclusive state)

| No. | Current State | Handling Message | Next State | Dequeue Message? | Action |
|-----|---------------|------------------|------------|------------------|--------|
| 14 | C-exclusive | Load | C-exclusive | Yes | reads cache |
| 15 | C-exclusive | Store | C-exclusive | Yes | writes cache |
| 16 | C-exclusive | WbReq(a) | C-shared | Yes | WbRep(id, Home, data(a)) |
| 17 | C-exclusive | FlushReq(a) | C-nothing | Yes | FlushRep(id, Home, data(a)) |
| 18 | C-exclusive | (Voluntary Writeback) | C-shared | N/A | WbRep(id, Home, data(a)) |
| 19 | C-exclusive | (Voluntary Flush) | C-nothing | N/A | FlushRep(id, Home, data(a)) |

# Cache Transitions (from pending)

| No. | Current State | Handling Message | Next State | Dequeue Message? | Action |
|-----|---------------|------------------|------------|------------------|--------|
| 20 | C-pending | WbReq(a) | C-pending | Yes | None |
| 21 | C-pending | FlushReq(a) | C-pending | Yes | None |
| 22 | C-pending | InvReq(a) | C-pending | Yes | None |
| 23 | C-pending | ShRep(a) | C-shared | Yes | updates cache with data |
| 24 | C-pending | ExRep(a) | C-exclusive | Yes | update cache with data |

# Home Directory State Transitions

| No. | Current State | Message Received | Next State | Dequeue Message? | Action |
|---|---|---|---|---|---|
| 1 | R(dir) & (dir = ε) | ShReq(a) | R({id}) | Yes | ShRep(Home, id, data(a)) |
| 2 | R(dir) & (dir = ε) | ExReq(a) | W(id) | Yes | ExRep(Home, id, data(a)) |
| 3 | R(dir) & (dir = ε) | (Voluntary Prefetch) | R({id}) | N/A | ShRep(Home, id, data(a)) |
| 4 | R(dir) & (id ∉ dir) & (dir ≠ ε) | ShReq(a) | R(dir + {id}) | Yes | ShRep(Home, id, data(a)) |
| 5 | R(dir) & (id ∉ dir) & (dir ≠ ε) | ExReq(a) | Tr(dir) | No | InvReq(Home, dir, a) |
| 6 | R(dir) & (id ∉ dir) & (dir ≠ ε) | (Voluntary Prefetch) | R(dir + {id}) | N/A | ShRep(Home, id, data(a)) |

Messages sent from site *id*

# Home Directory State Transitions

| No. | Current State | Message Received | Next State | Dequeue Message? | Action |
|---|---|---|---|---|---|
| 7 | R(dir) & (dir = {id}) | ShReq(a) | R(dir) | Yes | None |
| 8 | R(dir) & (dir = {id}) | ExReq(a) | W(id) | Yes | ExRep(Home, id, data(a)) |
| 9 | R(dir) & (dir = {id}) | InvRep(a) | R(ε) | Yes | None |
| 10 | R(dir) & (id ∈ dir) & (dir ≠ {id}) | ShReq(a) | R(dir) | Yes | None |
| 11 | R(dir) & (id ∈ dir) & (dir ≠ {id}) | ExReq(a) | Tr(dir-{id}) | No | InvReq(Home, dir - {id}, a) |
| 12 | R(dir) & (id ∈ dir) & (dir ≠ {id}) | InvRep(a) | R(dir - {id}) | Yes | None |

Messages sent from site *id*

# Home Directory State Transitions

| No. | Current State | Message Received | Next State | Dequeue Message? | Action |
|-----|---------------|------------------|------------|------------------|--------|
| 13 | W(id') | ShReq(a) | Tw(id') | No | WbReq(Home, id', a) |
| 14 | W(id') | ExReq(a) | Tw(id') | No | FlushReq(Home, id', a) |
| 15 | W(id) | ExReq(a) | W(id) | Yes | None |
| 16 | W(id) | WbRep(a) | R({id}) | Yes | data -> memory |
| 17 | W(id) | FlushRep(a) | R(ε) | Yes | data -> memory |

Messages sent from site *id*

# Home Directory State Transitions

| No. | Current State | Message Received | Next State | Dequeue Message? | Action |
|-----|---------------|------------------|------------|------------------|--------|
| 18 | Tr(dir) & (id ∈ dir) | InvRep(a) | Tr(dir - {id}) | Yes | None |
| 19 | Tr(dir) & (id ∉ dir) | InvRep(a) | Tr(dir) | Yes | None |
| 20 | Tw(id) | WbRep(a) | R({id}) | Yes | data-> memory |
| 21 | Tw(id) | FlushRep(a) | R($\epsilon$) | Yes | data-> memory |

Messages sent from site *id*

# Acknowledgements

- These slides contain material developed and copyright by:
  - Arvind (MIT)
  - Krste Asanovic (MIT/UCB)
  - Joel Emer (Intel/MIT)
  - James Hoe (CMU)
  - John Kubiatowicz (UCB)
  - David Patterson (UCB)

- MIT material derived from course 6.823
- UCB material derived from course CS252