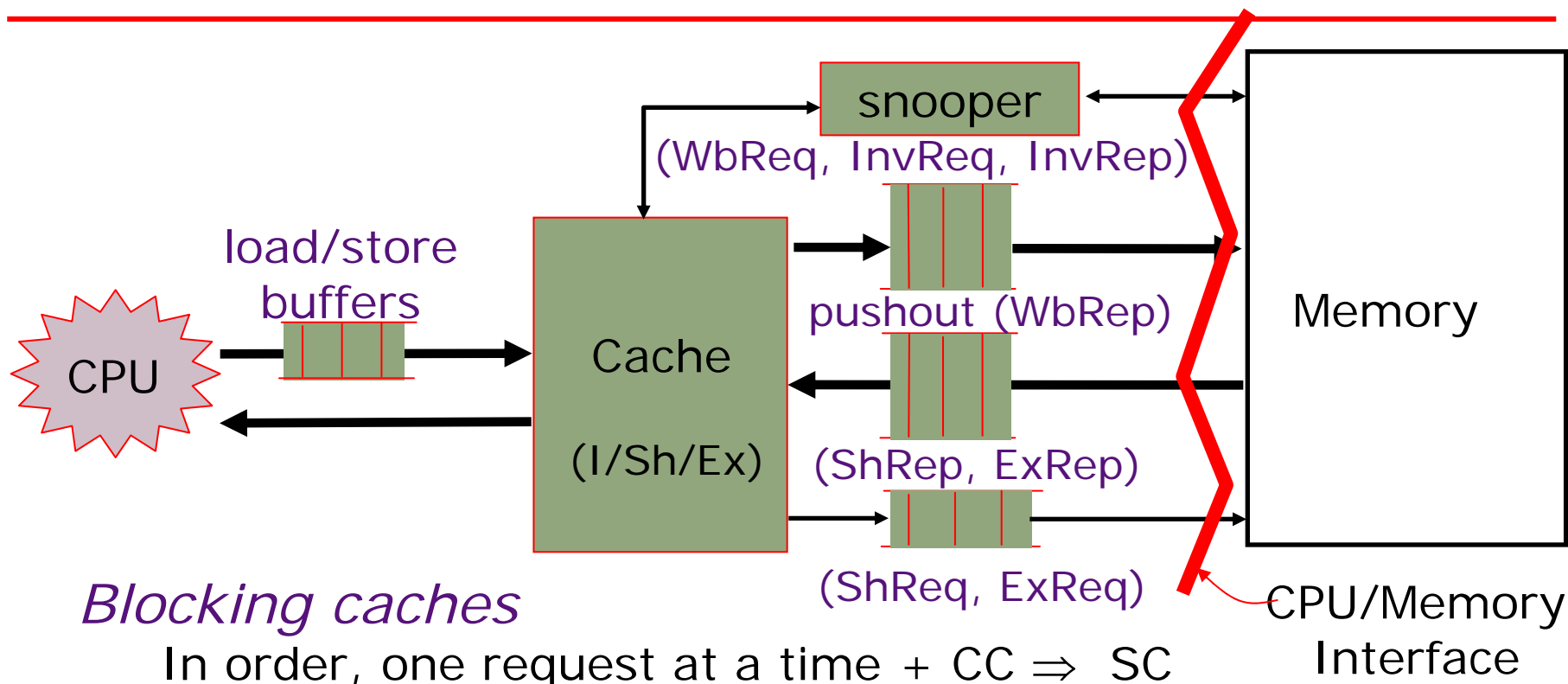# Cache Coherence Protocols
# for
# Sequential  Consistency

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

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

# Systems view



snooper

(WbReq, InvReq, InvRep)

load/store buffers

pushout (WbRep)

Memory

CPU

Cache

(I/Sh/Ex)

(ShRep, ExRep)

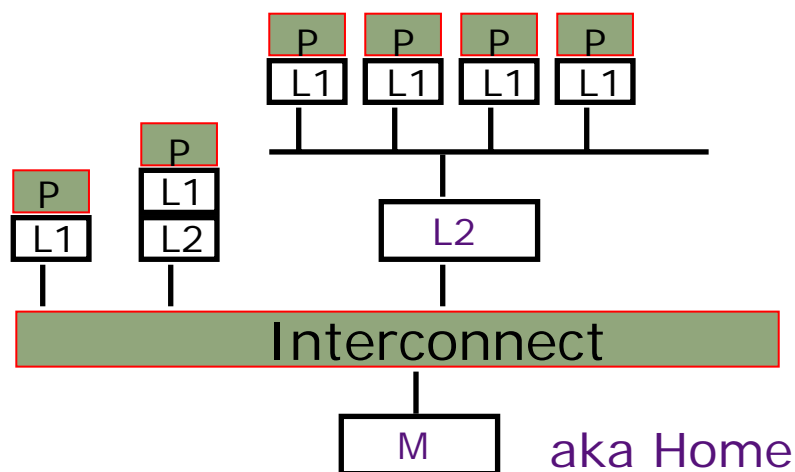(ShReq, ExReq)

CPU/Memory Interface

*Blocking caches*
   In order, 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

# A System with Multiple Caches



aka Home

Assumptions: Caches are organized in a hierarchical manner

- 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.,

$$a \in L_i \quad \Rightarrow \quad a \in L_{i+1}$$

# Maintaining Cache Coherence

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
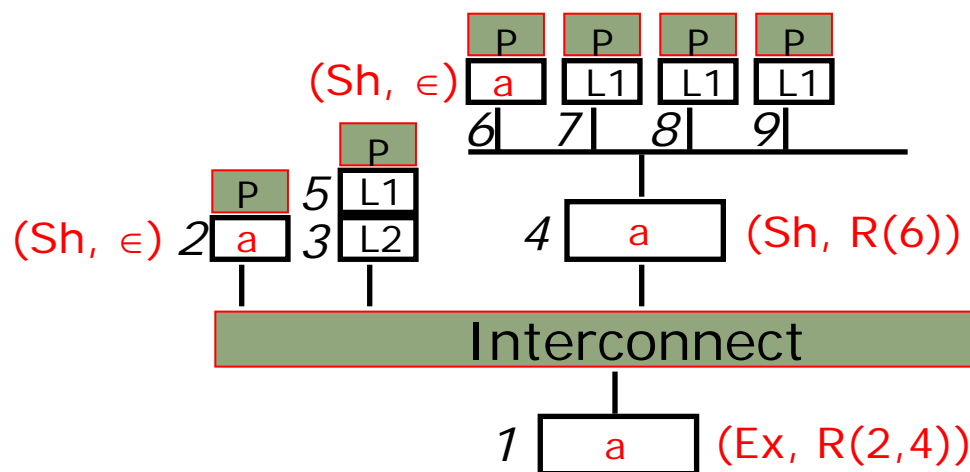
$\Rightarrow$

write request:
The address is *invalidated* in all other caches *before* 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

CSAIL

# State Encoding



Each address in a cache keeps two types of state info

- *sibling info:* do my siblings have a copy of address a
    - Ex (means no), Sh (means may be)
- *children info:* has this address been passed on to any of my children
    - W(id) means child id has a writable version
    - R(dir) means only children named in the *directory* dir have copies
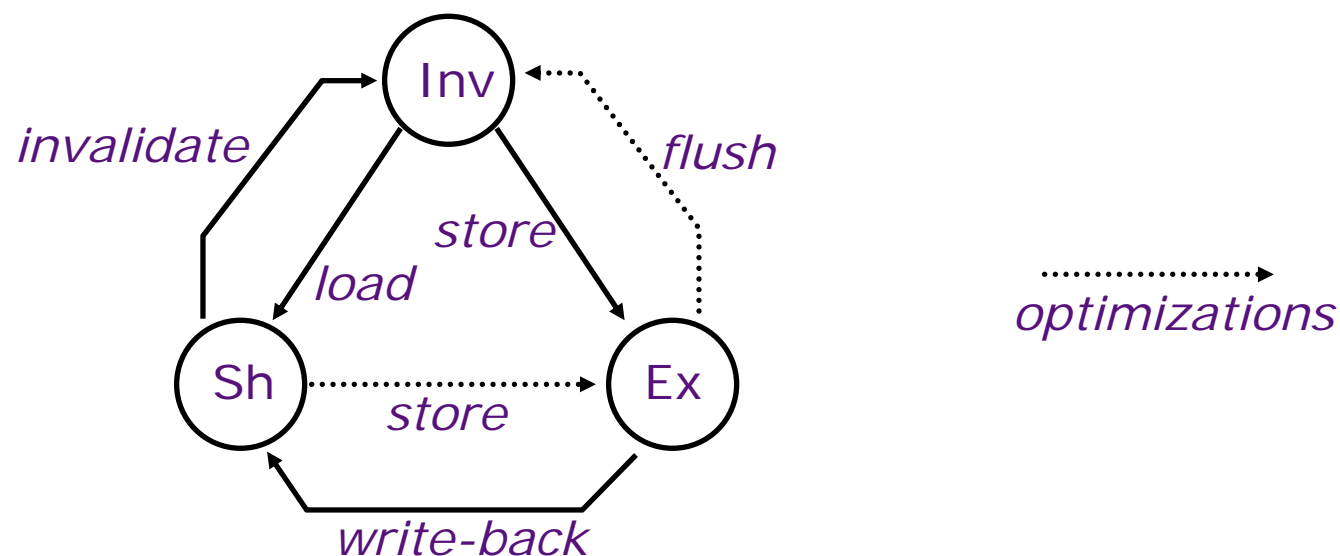
# Cache State Implications

Sh $\Rightarrow$ cache's siblings and decedents can only
have Sh copies

Ex $\Rightarrow$ each ancestor of the cache must be in Ex
$\Rightarrow$ *either* all children can have Sh copies
*or* one child can have an Ex copy

- Once a parent gives an Ex copy to a child, the parent's data is considered stale
- A processor cannot overwrite data in Sh state in L1
- By definition all addresses in the home are in the Ex state

CSAIL

# Cache State Transitions



This state diagram is helpful as long as one remembers that each transition involves cooperation of other caches and the main memory.

# High-level Invariants in Protocol Design

November 14, 2005

CSAIL

# Guarded Atomic Actions

- Rules specified using guarded atomic actions:

  <guard predicate>

  $\rightarrow$ {set of state updates that must occur
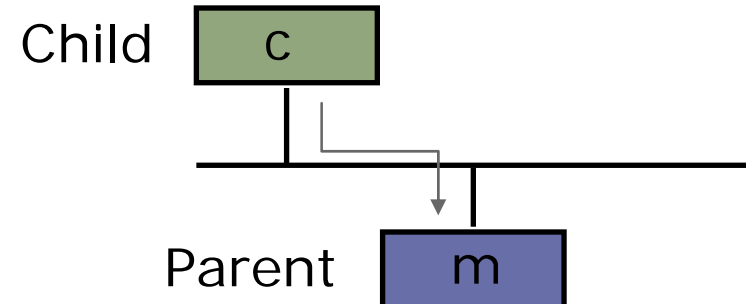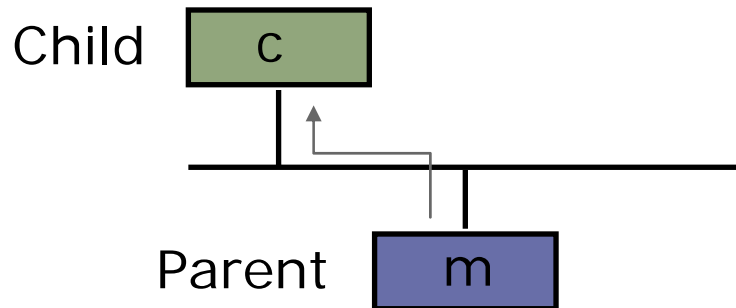  atomically with respect to other rules}

- E.g.:

  m.state(a) == R(dir) & $id_c \notin$ dir

  $\rightarrow$ m.setState(a, R(dir+ $id_c$)),
  c.setState(a, Sh); c.setData(a, m.data(a));

CSAIL

# Data Propagation Between Caches
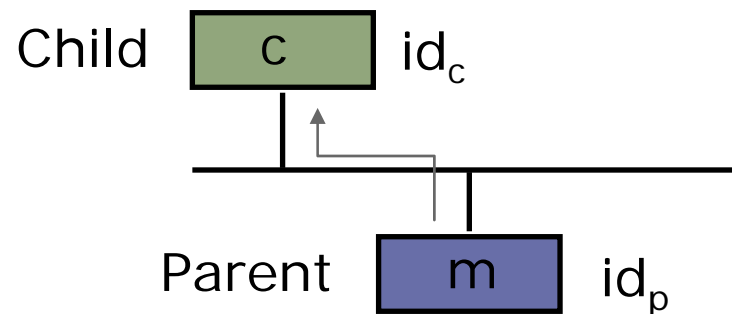
Child   c

Parent   m

Child   c

Parent   m

Caching rules
- *Read caching rule*
- *Write caching rule*

De-caching rules
- *Write-back rule*
- *Invalidate rule*

CSAIL

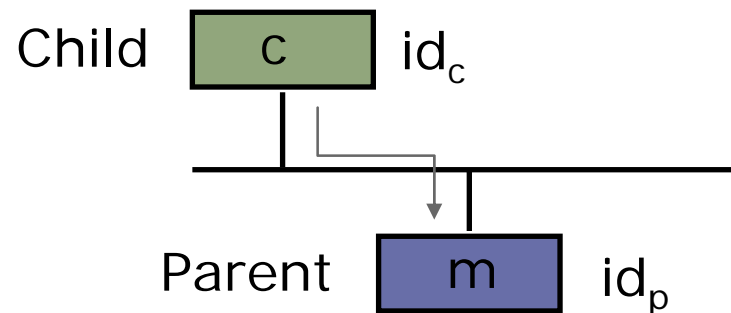# Caching Rules: *Parent to Child*

Child [ c ]  $id_c$

Parent [ m ]  $id_p$

- Read caching rule

$R(dir) == m.state(a)$  &  $id_c \notin dir$

$\rightarrow$  $m.setState(a, R(dir+ id_c))$

$c.setState(a, Sh);$    $c.setData(a, m.data(a));$

- Write caching rule

$\varepsilon == m.state(a)$

$\rightarrow$  $m.setState(a, W(id_c))$

$c.setState(a, Ex);$    $c.setData(a, m.data(a));$

November 14, 2005

# De-caching Rules: *Child to Parent*



- Writeback rule

    $W(id_c) == m.state(a)$    $\&$  $Ex == c.state(a)$

    $\rightarrow$ $m.setState(a, R(\{id_c\}))$

    $msetData(a, c.data(a));$

    $c.setState(a, Sh);$

- Invalidate rule

    $R(dir) == m.state(a) \&$ $id_c \in dir$ $\&$ $Sh == c.state(a)$

    $\rightarrow$ $m.setState(a, R(dir - id_c))$

    $c.invalidate(a);$

# Making the Rules *Local & Reactive*

Child $\quad$ c $\quad$ id$_c$

Parent $\quad$ m $\quad$ id$_p$

- Some rules require observing and changing the state of multiple caches simultaneously (atomically).
  - very difficult to implement, especially if caches are separated by a network
- Each rule must be triggered by some action
- Split rules are into multiple rules – "request for an action" followed by "an action and an ack".
  - ultimately all actions are triggered by some processor

November 14, 2005

# Protocol Design

*Note*
We will not be able to finish this part today.
(The rest of the material will be covered during the next lecture.)

# Protocol Processors *an abstract view*
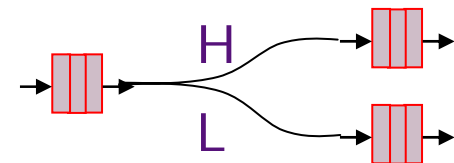


- Each cache has 2 pairs of queues
  - one pair (c2m, m2c) to communicate with the memory
  - one pair (p2m, m2p) to communicate with the processor

- Messages format:
  msg(idsrc,iddest,cmd,priority,a,v)

- FIFO messages passing between each (src,dest) pair except a Low priority (L) msg cannot block a high priority (H) msg

# H and L Priority Messages

- At the memory unprocessed requests cannot block the result messages. Hence all messages are classified as H or L priority.
  – all messages carrying results are classified as high priority

- Accomplished by having separate paths for H and L priority
  – In Theory: separate networks
  – In Practice:
    - Separate Queues
    - Shared buses for both networks

CSAIL

# A Protocol for a system with two memory levels (L1 + M)

Cache states:          Sh, Ex, Pending, Nothing

Memory states:        $R(dir)$,   $W(id)$, $T_R(dir)$, $T_W(id)$

*If dir is empty then $R(dir)$ and $T_R(dir)$ represent the same state*

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)

Operations on cache:

cache.state(a) – returns state s

cache.data(a) - returns data v

cache.setState(a,s), cache.setData(a,v), cache.invalidate(a)

inst = first(p2m);   msg= first(m2c);   mmsg = first(in)

CSAIL

# Voluntary rules: Cache must be able to evict values to create space

Invalidate rule

   cache.state(a) is Sh

→  cache.invalidate(a)

   c2m.enq (Msg(id, Home, InvRep, a))

> *It would be good to have "silent drops" but difficult in a directory-based protocol*

Flush rule

   cache.state(a) is Ex

→  cache.invalidate(a)

   c2m.enq (Msg(id, Home, FlushRep, a, cache.data(v))

Writeback rule

   cache.state(a) is Ex

→  cache.setState(a, Sh)

   c2m.enq (Msg(id, Home, WbRep, a, cache.data(v)))

  *This rule may be applied if the cache/processor knows it is the "last store" operation to the location.*

  *Such voluntary rules can be used to construct adaptive protocols.*

# Voluntary rules: Memory should be able to send more values than requested

Cache Rule

$\quad$ m.state(a) is R(dir) &  id $\notin$ dir

$\rightarrow\quad$ m.setState(a, R(id+dir))

$\qquad$ out.enq(Msg(Home,id,ShRep, a,m.data(a)))

*It is a rule like this that allows us to fetch locations a+1, a+2, … when a processor requests address a.*

November 14, 2005

CSAIL

# Five-minute break to stretch your legs

# Load Rules (at cache)

- Load-hit rule

  Load(a)==inst

  &   cache.state(a) is Sh or Ex

  →      p2m.deq

  m2p.enq(cache.data(a))

- Load-miss rule

  Load(a)==inst

  &   cache.state(a) is Nothing

  →      c2m.enq(Msg(id, Home, ShReq, a)

  cache.setState(a,Pending)

  This is blocking cache because the Load miss
  rule does not remove the request from the
  input queue (p2m)    ... *more later*

# Store Rules (at cache)

- ## Store-hit rule

  Store(a,v)==inst
  & cache.state(a) is Ex
  → p2m.deq;
  m2p.enq(Ack)
  cache.setData(a, v)

- ## Store-miss rules

  Store(a,v)==inst
  & cache.state(a) is Nothing
  → c2m.enq(Msg(id, Home, ExReq, a);
  cache.setState(a,Pending)

  Store(a,v)==inst
  & cache.state(a) is Sh
  → c2m.enq(Msg(id, Home, InvRep, a);
  cache.setState(a,Nothing)

**Already covered by the Invalidate voluntary rule**

CSAIL

# Processing ShReq Messages (at Home)

**Uncached or Outstanding Shared Copies**

$\quad$ Msg(id,Home,ShReq,a) ==mmsg

& $\quad$ m.state(a) is R(dir) & id $\notin$ dir

$\rightarrow \quad$ in.deq;

$\quad$ m.setState(a, R(dir+{id}));

$\quad$ out.enq(Msg(Home,id,ShRep, a,m.data(a)))


**Outstanding Exclusive Copy**

$\quad$ Msg(id,Home,ShReq,a) ==mmsg

& $\quad$ m.state(a) is W(id$'$) & (id$'$ is not id)

$\rightarrow \quad$ m.setState(a, $T_W$(id$'$));

$\quad$ out.enq(Msg(Home,id$'$,WbReq, a))

# Processing ExReq Messages (at home)

Uncached or cached only at the requester cache

Msg(id,Home,ExReq,a) ==mmsg
& m.state(a) is R(dir) & (dir is empty or has only id)
→ in.deq
m.setState(a, W(id))
out.enq(Msg(Home,id,ExRep, a, m.data(a))

Outstanding Shared Copies

Msg(id,Home,ExReq,a) ==mmsg
& m.state(a) is R(dir) & !(dir is empty or has only id)
→ m.setState(a, $T_R$(dir-{id}))
out.enq(*multicast*(Home,dir-{id},InvReq, a)

Outstanding Exclusive Copy

Msg(id,Home,ExReq,a) ==mmsg
& m.state(a) is W(id′) & (id′ is not id)
→ m.setState(a, $T_W$(id′))
out.enq(Msg(Home,id′,FlushReq, a)

CSAIL

# Processing Reply Messages (at cache)

ShRep

$\quad$ Msg(Home, id, ShRep, a, v) == msg

-- cache.state(a) must be Pending or Nothing

$\rightarrow \quad$ m2c.deq

$\quad$ cache.setState(a, Sh)

$\quad$ cache.setData(a, v)

ExRep

$\quad$ Msg(Home, id, ExRep, a, v) == msg

-- cache.state(a) must be Pending or Nothing

$\rightarrow \quad$ m2c.deq

$\quad$ cache.setState(a, Ex)

$\quad$ cache.setData(a, v)

$\qquad$ -- In general only a part of v will be
$\qquad$ overwritten by the Store instruction.

CSAIL

# Processing InvReq Message (at cache)

InvReq

        Msg(Home,id,InvReq,a) == msg

  &   cache.state(a) is Sh

$\rightarrow$     m2c.deq

        cache.invalidate(a)

        c2m.enq (Msg(id, Home, InvRep, a))

        Msg(Home,id,InvReq,a) == msg

  &   cache.state(a) is Nothing or Pending

$\rightarrow$     m2c.deq

CSAIL

# Processing WbReq Message (at cache)

WbReq

Msg(Home,id,WbReq,a) == msg
& cache.state(a) is Ex
→ m2c.deq
cache.setState(a, Sh)
c2m.enq (Msg(id, Home, WbRep, a, cache.data(v)))

Msg(Home,id,WbReq,a) == msg
& cache.state(a) is Sh or Nothing or Pending
→ m2c.deq

CSAIL

# Processing FlushReq Message (at cache)

FlushReq

Msg(Home,id,FlushReq,a) == msg
&   cache.state(a) is Ex
→      m2c.deq
cache.invalidate(a)
c2m.enq (Msg(id, Home, FlushRep, a, cache.data(v)))

Msg(Home,id,FlushReq,a) == msg
&   cache.state(a) is Sh
→      m2c.deq
cache.invalidate(a)
c2m.enq (Msg(id, Home, InvRep, a))

Msg(Home,id,FlushReq,a) == msg
&   cache.state(a) is Nothing or Pending
→      m2c.deq

CSAIL

# Processing Reply InvRep Messages
## (at home)

InvRep

$\quad\quad$ Msg(id,Home,InvRep,a) == mmsg

$\quad$ & $\;$ m.state(a) is $T_R$(dir)

$\rightarrow\quad\quad$ deq mmsg;

$\quad\quad$ m.setState(a, $T_R$(dir-{id}))


$\quad\quad$ Msg(id,Home,InvRep,a) == mmsg

$\quad$ & $\;$ m.state(a) is R(dir)

$. \rightarrow\quad\quad$ deq mmsg;

$\quad\quad$ m.setState(a, R(dir-{id}))

CSAIL

# Processing Reply WbRep Messages
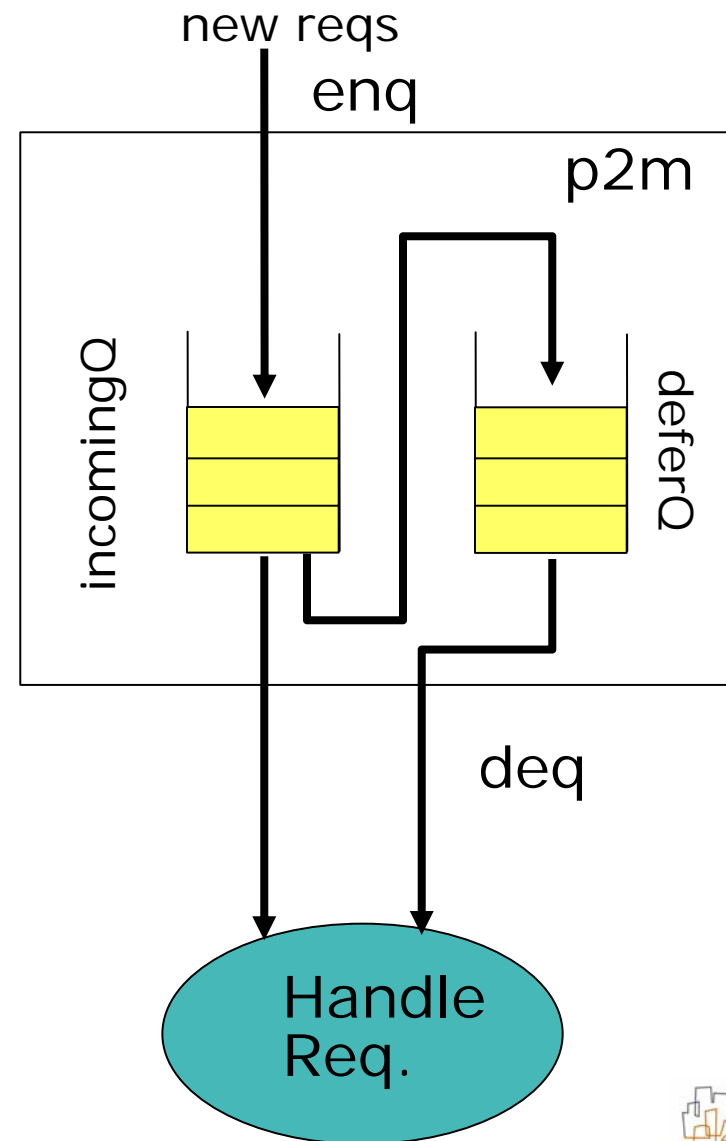## (at home)

### WbRep

Msg(id,Home,WbRep,a,v) == mmsg

-- m.state(a) must be $T_W$(id) or W(id)

$\rightarrow$    deq mmsg;

m.setState(a, R(id))

m.setData(a,v)


### FlushRep

Msg(id,Home,FlushRep,a,v) == mmsg

-- m.state(a) must be $T_W$(id) or W(id)

$\rightarrow$    deq mmsg;

m.setState(a, R(Empty))

m.setData(a,v)

# Non-Blocking Caches

- Non-blocking caches are needed to tolerate large memory latencies

- To get non-blocking property we implement p2m with 2 FIFOs (deferQ, incomingQ)

- Requests moved to deferQ when:
  – address not there
  – needed for consistency

new reqs

enq

p2m

incomingQ

deferQ

deq

Handle Req.

# Conclusion

- This protocol with its voluntary rules captures many other protocols that are used in practice.
    - we will discuss a bus-based version of this protocol in the next lecture
- We need policies and mechanisms to invoke voluntary rules to build truly adaptive protocols.
    - search for such policies and mechanisms in an active area of research
- Quantitative evaluation of protocols or protocol features is extremely difficult.
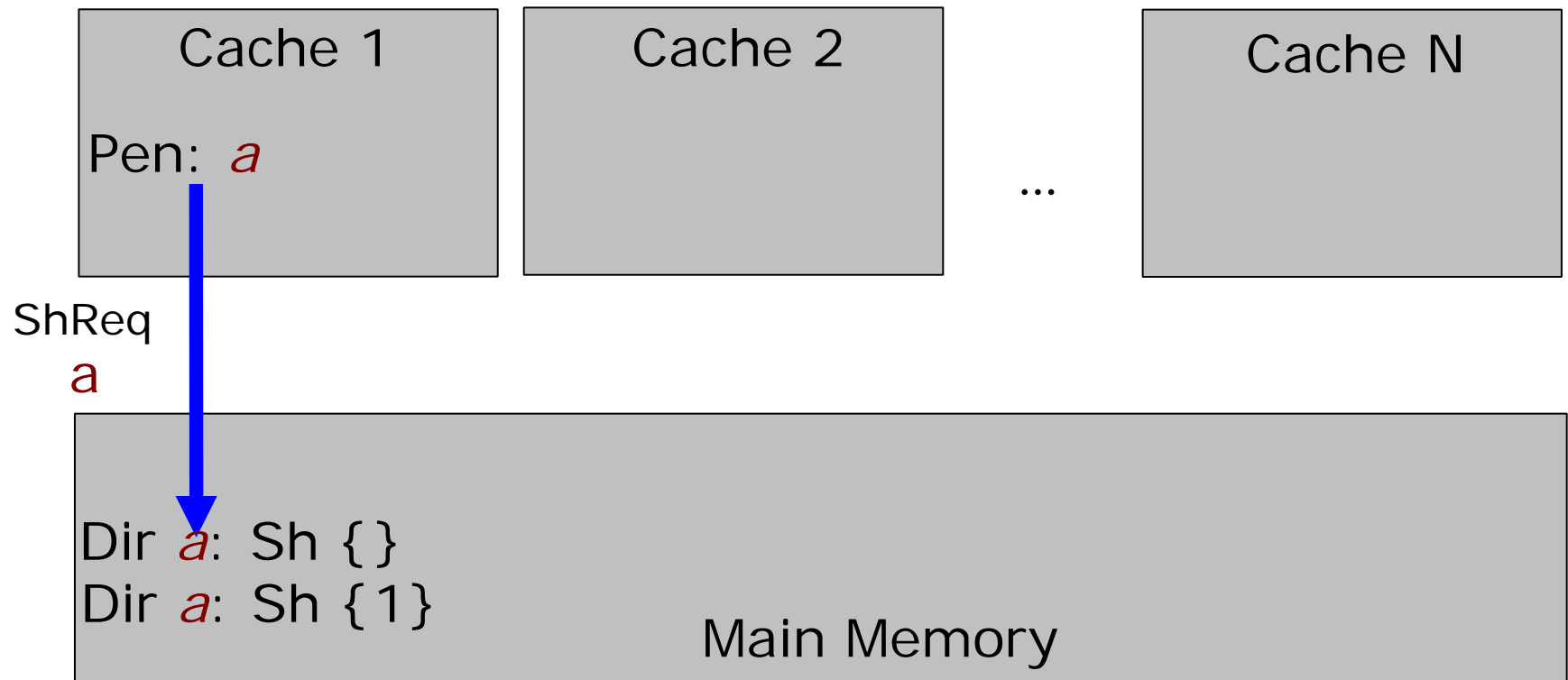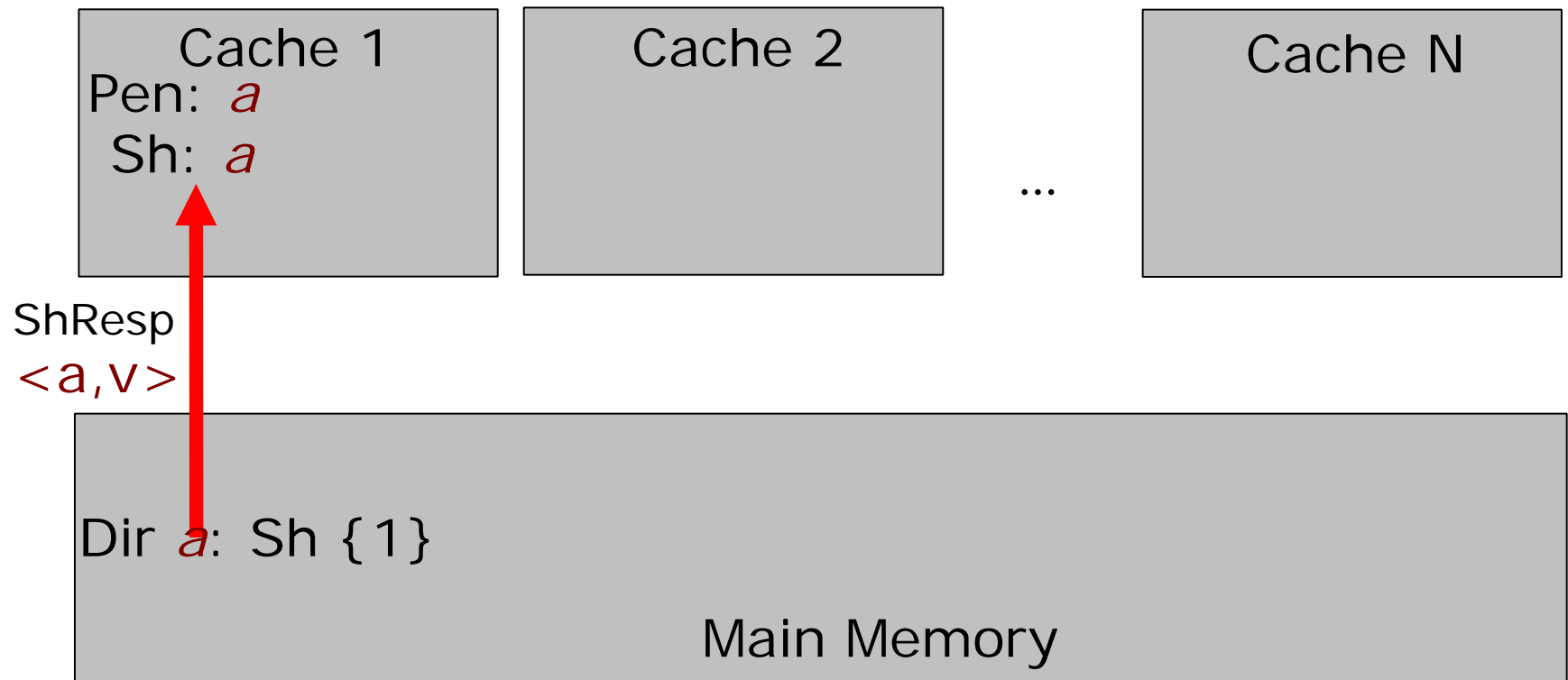
*Thank you !*

# Protocol Diagram

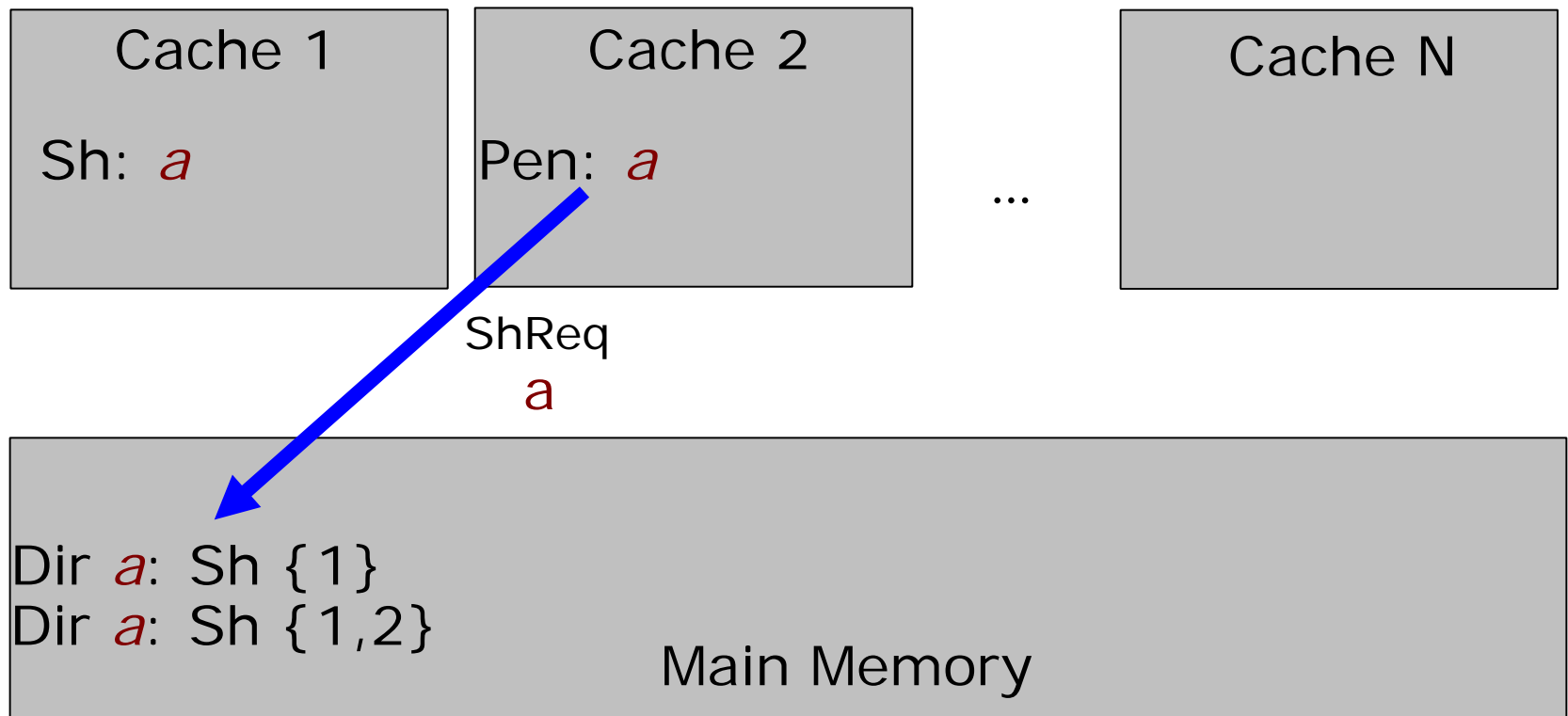| Cache 1 | Cache 2 | ... | Cache N |
|---------|---------|-----|---------|

Pen: *a*

ShReq
a

Dir *a*: Sh { }
Dir *a*: Sh {1}

Main Memory

# Protocol Diagram

| Cache 1 | Cache 2 | | Cache N |
|---|---|---|---|
| Pen: *a* | | ... | |
| Sh: *a* | | | |

ShResp
<a,v>

Dir *a*: Sh {1}

Main Memory

# Protocol Diagram

| Cache 1 | Cache 2 | ... | Cache N |
|---------|---------|-----|---------|
| Sh: $a$ | Pen: $a$ | | |

ShReq
a

Dir $a$: Sh {1}
Dir $a$: Sh {1,2}

Main Memory

CSAIL

# Protocol Diagram

| Cache 1 | Cache 2 | | Cache N |
|---------|---------|---|---------|
| Sh: *a* | Pen: *a*<br>Sh: *a* | ... | |

ShResp
<a,v>

Dir *a*: Sh {1,2}

Main Memory

CSAIL

# Protocol Diagram



Cache 1

Sh: *a*

Cache 2

Sh: *a*

...

Cache N

Pen: *a*

InvReq
a

InvReq
a

ExReq
a

Dir *a*: Sh {1,2}

Main Memory

CSAIL

# Protocol Diagram

# Protocol Diagram

| Cache 1 | Cache 2 | | Cache N |
|---------|---------|---|---------|
| Pen: *a* | | ... | Ex: *a* |

ShReq
a

WBReq
a

Dir *a*: Ex {N}

Main Memory

CSAIL

# Protocol Diagram