

# Lecture 4: Directory-Based Coherence

---

- Details of memory-based (SGI Origin) and cache-based (Sequent NUMA-Q) directory protocols

# Handling Reads

---

- When the home receives a read request, it looks up memory (speculative read) and directory in parallel
- Actions taken for each directory state:
  - shared or unowned: memory copy is clean, data is returned to requestor, state is changed to excl if there are no other sharers
  - busy: a NACK is sent to the requestor
  - exclusive: home is not the owner, request is fwded to owner, owner sends data to requestor and home

# Inner Details of Handling the Read

---

- The block is in exclusive state – memory may or may not have a clean copy – it is speculatively read anyway
- The directory state is set to busy-exclusive and the presence vector is updated
- In addition to fwding the request to the owner, the memory copy is speculatively forwarded to the requestor
  - Case 1: excl-dirty: owner sends block to requestor and home, the speculatively sent data is over-written
  - Case 2: excl-clean: owner sends an ack (without data) to requestor and home, requestor waits for this ack before it moves on with speculatively sent data

# Inner Details II

---

- Why did we send the block speculatively to the requestor if it does not save traffic or latency?
  - the R10K cache controller is programmed to not respond with data if it has a block in excl-clean state
  - when an excl-clean block is replaced from the cache, the directory need not be updated – hence, directory cannot rely on the owner to provide data and speculatively provides data on its own

# Handling Write Requests

---

- The home node must invalidate all sharers and all invalidations must be acked (to the requestor), the requestor is informed of the number of invalidates to expect
- Actions taken for each state:
  - shared: invalidates are sent, state is changed to excl, data and num-sharers is sent to requestor, the requestor cannot continue until it receives all acks (Note: the directory does not maintain busy state, subsequent requests will be fwded to new owner and they must be buffered until the previous write has completed)

# Handling Writes II

---

- Actions taken for each state:
  - unowned: if the request was an upgrade and not a read-exclusive, is there a problem?
  - exclusive: is there a problem if the request was an upgrade? In case of a read-exclusive: directory is set to busy, speculative reply is sent to requestor, invalidate is sent to owner, owner sends data to requestor (if dirty), and a “transfer of ownership” message (no data) to home to change out of busy
  - busy: the request is NACKed and the requestor must try again

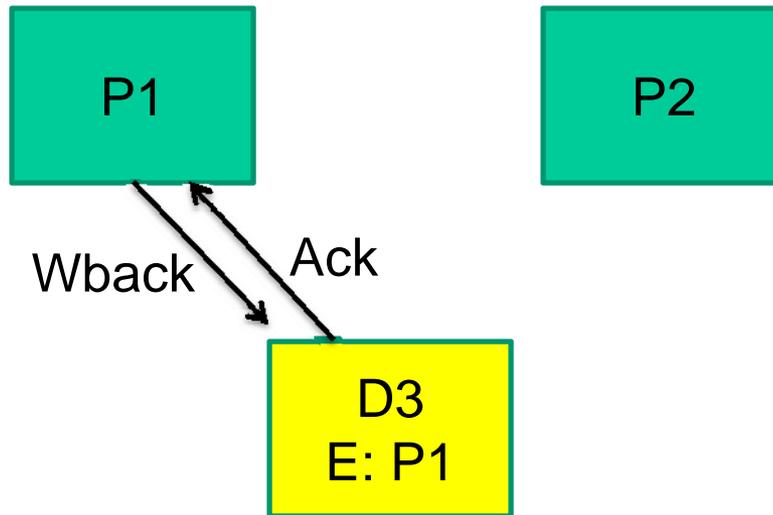
# Handling Write-Back

---

- When a dirty block is replaced, a writeback is generated and the home sends back an ack
- Can the directory state be shared when a writeback is received by the directory?
- Actions taken for each directory state:
  - exclusive: change directory state to unowned and send an ack
  - busy: a request and the writeback have crossed paths: the writeback changes directory state to shared or excl (depending on the busy state), memory is updated, and home sends data to requestor, the intervention request is dropped

# Writeback Cases

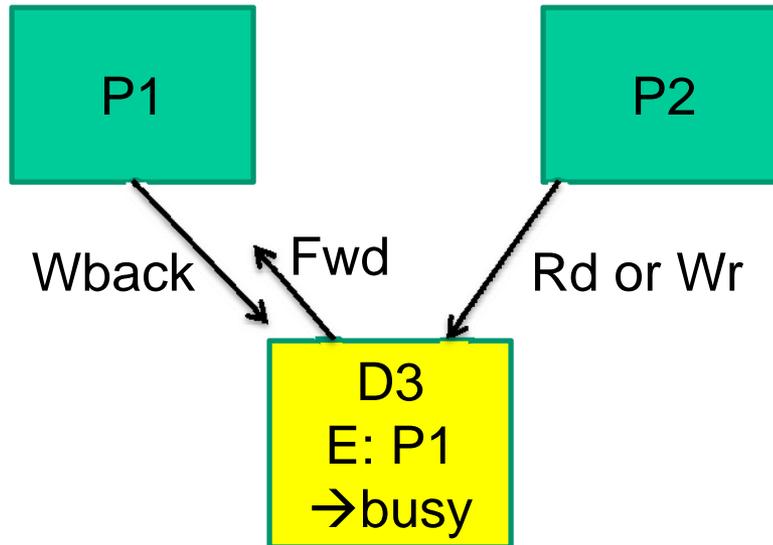
---



This is the “normal” case  
D3 sends back an Ack

# Writeback Cases

---



If someone else has the block in exclusive, D3 moves to busy

If Wback is received, D3 serves the requester

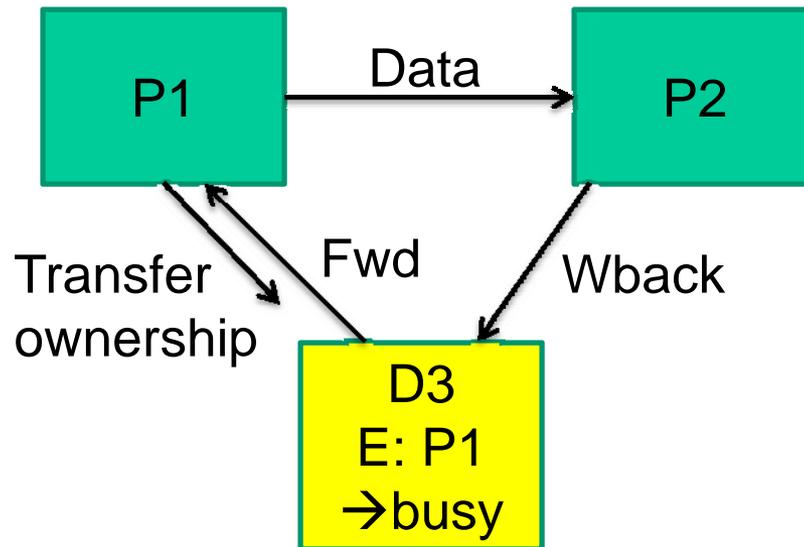
If we didn't use busy state when transitioning from E:P1 to E:P2,  
D3 may not have known who to service

(since ownership may have been passed on to P3 and P4...)

(although, this problem can be solved by NACKing the Wback  
and having P1 buffer its “strange” intervention requests)

# Writeback Cases

---



If Wback is from new requester, D3 sends back a NACK

Floating unresolved messages are a problem

Alternatively, can accept the Wback and put D3 in some new busy state

Conclusion: could have got rid of busy state between E:P1 → E:P2, but with Wback ACK/NACK and other buffering could have kept the busy state between E:P1 → E:P2, could have got rid of ACK/NACK, but need one new busy state<sup>10</sup>

# Serialization

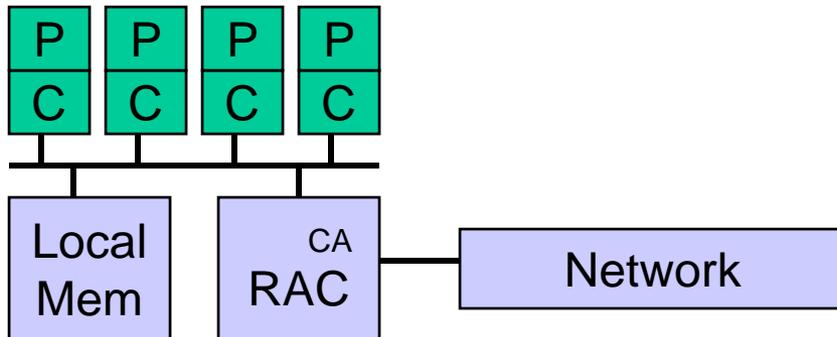
---

- Note that the directory serializes writes to a location, but does not know when a write/read has completed at any processor
- For example, a read reply may be floating on the network and may reach the requestor much later – in the meantime, the directory has already issued a number of invalidates, the invalidate is overwritten when the read reply finally shows up – hence, each node must buffer its requests until outstanding requests have completed

# Sequent NUMA-Q

---

- Employs a flat cache-based directory protocol between nodes – IEEE standard SCI (Scalable Coherent Interface) protocol
- Each node is a 4-way SMP with a bus-based snooping protocol
- The communication assist includes a large “remote access cache” – the directory protocol tries to keep the remote caches coherent, while the snooping protocol ensures that each processor cache is kept coherent with the remote access cache and local-mem



# Directory Structure

---

- The physical address identifies the home node – the home node directory stores a pointer to the head of a linked list – each cache stores pointers to the next and previous sharer
- A main memory block can be in three directory states:
  - Home: (similar to unowned) the block does not exist in any remote access cache (may be in the home node's processor caches, though)
  - Fresh: (similar to shared) read-only copies exist in remote access caches and memory copy is up-to-date
  - Gone: (similar to exclusive) writeable copy exists in some remote cache

# Cache Structure

---

- 29 stable states and many more pending/busy states!
- The stable states have two descriptors:
  - position in linked list: ONLY, HEAD, TAIL, MID
  - state within cache: dirty, clean, fresh, valid, etc.
- SCI defines and implements primitive operations to facilitate linked list manipulations:
  - List construction: add a new node to the list head
  - Rollout: remove a node from a list
  - Purging: invoked by the head to invalidate all other nodes

# Handling Read Requests

---

- On a read miss, the remote cache sets up a block in busy state and other requests to the block are not entertained
- The requestor sends a “list construction request” to the home and the steps depend on the directory state:
  - Home: state updated to fresh, head updated to requestor, data sent to requestor, state at requestor is set to ONLY\_FRESH
  - Fresh: head updated to requestor, home responds with data and pointer to old head, requestor moves to a different busy state, sends list construction request to old head, old head moves from HEAD\_FRESH to MID\_VALID, sends ack, requestor → HEAD\_FRESH

# Handling Read Requests II

---

- Gone: home does not reply with data, it remains in Gone state, sends old head pointer to requestor, requestor moves to a different busy state, asks old head for data and “list construction”, old head moves from HEAD\_DIRTY to MID\_VALID, returns data, requestor moves to HEAD\_DIRTY (note that HEAD\_DIRTY does not mean exclusive access; the head can write without talking to the home, but sharers must be invalidated)
- Home keeps forwarding requests to head even if head is busy – this results in a pending linked list that is handled as transactions complete

# Handling Write Requests

---

- At all times, the head of a list is assumed to have the latest copy and only the head is allowed to write
- The writer starts by moving itself to the head of the list; actions depend on the state in the cache:
  - HEAD\_DIRTY: the home is already in GONE state, so home is not informed, sharing list is purged (each list element invalidates itself and informs the requestor of the next element – simple, but slow – works well for small invalidation sizes)

# Handling Write Requests II

---

- HEAD\_FRESH: home directory is updated from FRESH to GONE, sharing list is purged; if the home directory is not in FRESH state, some other node's request is in flight – the requestor will have to move to the head again and retry
- ONLY\_DIRTY: the write happens without generating any interconnect traffic

# Writeback & Replacement

---

- Replacements are no longer “quiet” as the linked lists have to be updated – the “rollout” operation is used
- To rollout, a node must set itself to pending, inform the neighbors, and set itself to invalid – to prevent deadlock in the case of two neighbors attempting rollout, the node closer to the tail is given priority
- If the node is the head, it makes the next element the head and informs home

# Writeback & Replacement II

---

- If the head is attempting a rollout, it sends a message home, but the home is pointing to a different head: the old head will eventually receive a request from the new head – at this point, the writeback is complete, and the new head is instead linked with the next node
- To reduce buffering needs, the writeback happens before the new block is fetched

# Serialization

---

- The home serves as the point of serialization – note that requests are almost never NACKed – requests are usually re-directed to the current head – helps avoid race conditions
- Since requests get queued in a pending list and buffers are rarely used, the protocol is less prone to starvation, unfairness, deadlock, and livelock problems

# Title

---

- Bullet