

## Protocol Design Space of Snooping Cache Coherent Multiprocessors

CS 258, Spring 99  
David E. Culler  
Computer Science Division  
U.C. Berkeley

## Recap

- **Snooping cache coherence**
  - solve difficult problem by applying extra interpretation to naturally occurring events
    - » state transitions, bus transactions
  - write-thru cache
    - » 2-state: invalid, valid
    - » no new transaction, no new wires
    - » coherence mechanism provides consistency, since all writes in bus order
    - » poor performance
- **Coherent memory system**
- **Sequential Consistency**

2/12/99

CS258 S99

2

## Sequential Consistency

- Memory operations from a proc become visible (to itself and others) in program order
- There exist a total order, consistent with this partial order - i.e., an interleaving
  - the position at which a write occurs in the hypothetical total order should be the same with respect to all processors
- **Sufficient Conditions**
  - every process issues mem operations in program order
  - after a write operation is issued, the issuing process waits for the write to complete before issuing next memory operation
  - after a read is issued, the issuing process waits for the read to complete and for the write whose value is being returned to complete (globally) before issuing its next operation
- **How can compilers violate SC? Architectural enhancements?**

2/12/99

CS258 S99

3

## Outline for Today

- **Design Space of Snoopy-Cache Coherence Protocols**
  - write-back, update
  - protocol design
  - lower-level design choices
- **Introduction to Workload-driven evaluation**
- **Evaluation of protocol alternatives**

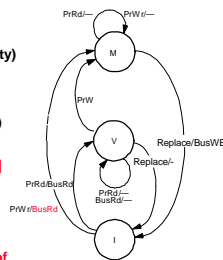
2/12/99

CS258 S99

4

## Write-back Caches

- **2 processor operations**
    - PrRd, PrWr
  - **3 states**
    - invalid, valid (clean), modified (dirty)
    - ownership: who supplies block
  - **2 bus transactions:**
    - read (BusRd), write-back (BusWB)
    - only cache-block transfers
- => treat Valid as "shared" and Modified as "exclusive"**
- => introduce one new bus transaction**
- read-exclusive: read for purpose of modifying (read-to-own)



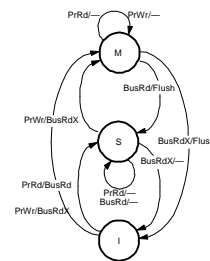
2/12/99

CS258 S99

5

## MSI Invalidate Protocol

- **Read obtains block in "shared"**
  - even if only cache copy
- **Obtain exclusive ownership before writing**
  - BusRdx causes others to invalidate (demote)
  - If M in another cache, will flush
  - BusRdx even if hit in S
    - » promote to M (upgrade)
- **What about replacement?**
  - S->I, M->I as before

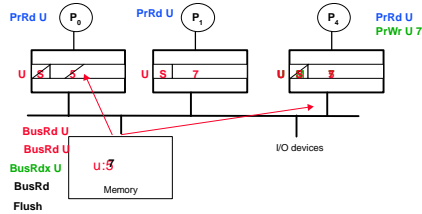


2/12/99

CS258 S99

6

## Example: Write-Back Protocol



2/12/99

CS258 S99

7

## Correctness

- When is write miss performed?
  - How does writer “observe” write?
  - How is it “made visible” to others?
  - How do they “observe” the write?
- When is write hit made visible?

2/12/99

CS258 S99

8

## Write Serialization for Coherence

- Writes that appear on the bus (BusRdX) are ordered by bus
  - performed in writer’s cache before other transactions, so ordered same w.r.t. all processors (incl. writer)
  - Read misses also ordered wrt these
- Write that don’t appear on the bus:
  - P issues BusRdX B.
  - further mem operations on B until next transaction are from P
    - » read and write hits
    - » these are in program order
  - for read or write from another processor
    - » separated by intervening bus transaction
- Reads hits?

2/12/99

CS258 S99

9

## Sequential Consistency

- Bus imposes total order on bus xactions for all locations
- Between xactions, procs perform reads/writes (locally) in program order
- So any execution defines a natural partial order
  - $M_i$  subsequent to  $M_j$  if
    - » (i) follows in program order on same processor,
    - » (ii)  $M_j$  generates bus xaction that follows the memory operation for  $M_i$
- In segment between two bus transactions, any interleaving of local program orders leads to consistent total order
- w/i segment writes observed by proc P serialized as:
  - Writes from other processors by the previous bus xaction P issued
  - Writes from P by program order

2/12/99

CS258 S99

10

## Sufficient conditions

- Sufficient Conditions
  - issued in program order
  - after write issues, the issuing process waits for the write to complete before issuing next memory operation
  - after read is issues, the issuing process waits for the read to complete and for the write whose value is being returned to complete (globally) before issuing its next operation
- Write completion
  - can detect when write appears on bus
- Write atomicity:
  - if a read returns the value of a write, that write has already become visible to all others already

2/12/99

CS258 S99

11

## Lower-level Protocol Choices

- BusRd observed in M state: what transition to make?
  - M  $\rightarrow$  I
  - M  $\rightarrow$  S
  - Depends on expectations of access patterns
- How does memory know whether or not to supply data on BusRd?
- Problem: Read/Write is 2 bus xactions, even if no sharing
  - » BusRd (I $\rightarrow$ S) followed by BusRdX or BusUpgr (S $\rightarrow$ M)
  - » What happens on sequential programs?

2/12/99

CS258 S99

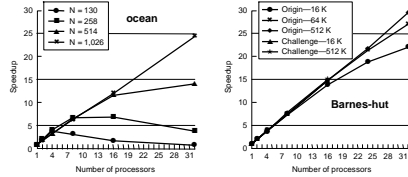
12





## A Lot Depends on Sizes

- Application parameters and no. of procs affect inherent properties
  - Load balance, communication, extra work, temporal and spatial locality
- Interactions with organization parameters of extended memory hierarchy affect artifactual communication and performance
- Effects often dramatic, sometimes small: application-dependent



Understanding size interactions and scaling relationships is key  
2/12/99 CS258 S99 25

## Scaling: Why Worry?

- Fixed problem size is limited
- Too small a problem:
  - May be appropriate for small machine
  - Parallelism overheads begin to dominate benefits for larger machines
    - Load imbalance
    - Communication to computation ratio
  - May even achieve slowdowns
  - Doesn't reflect real usage, and inappropriate for large machines
    - Can exaggerate benefits of architectural improvements, especially when measured as percentage improvement in performance
- Too large a problem
  - Difficult to measure improvement (next)

2/12/99 CS258 S99 26

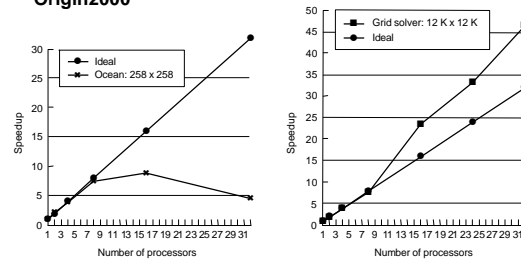
## Too Large a Problem

- Suppose problem realistically large for big machine
- May not "fit" in small machine
  - Can't run
  - Thrashing to disk
  - Working set doesn't fit in cache
- Fits at some  $p$ , leading to *superlinear speedup*
- Real effect, but doesn't help evaluate effectiveness
- Finally, users want to scale problems as machines grow
  - Can help avoid these problems

2/12/99 CS258 S99 27

## Demonstrating Scaling Problems

- Small Ocean and big equation solver problems on SGI Origin2000



2/12/99 CS258 S99 28

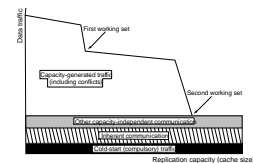
## Communication and Replication

- View parallel machine as extended memory hierarchy
  - Local cache, local memory, remote memory
  - Classify "misses" in "cache" at any level as for uniprocessors
    - compulsory or cold misses (no size effect)
    - capacity misses (yes)
    - conflict or collision misses (yes)
    - communication or coherence misses (no)
- Communication induced by finite capacity is most fundamental artifact
  - Like cache size and miss rate or memory traffic in uniprocessors

2/12/99 CS258 S99 29

## Working Set Perspective

- At a given level of the hierarchy (to the next further one)



- Hierarchy of working sets
- At first level cache (fully assoc, one-word block), inherent to algorithm
  - working set curve for program
- Traffic from any type of miss can be local or nonlocal (communication)

2/12/99 CS258 S99 30