

Lecture 19: Scalable Protocols & Synch

- Topics: coherence protocols for distributed shared-memory multiprocessors and synchronization (Sections 4.4-4.5)

Coherence Protocols

- Two conditions for cache coherence:
 - write propagation
 - write serialization
- Cache coherence protocols:
 - snooping
 - directory-based

 - write-update
 - write-invalidate

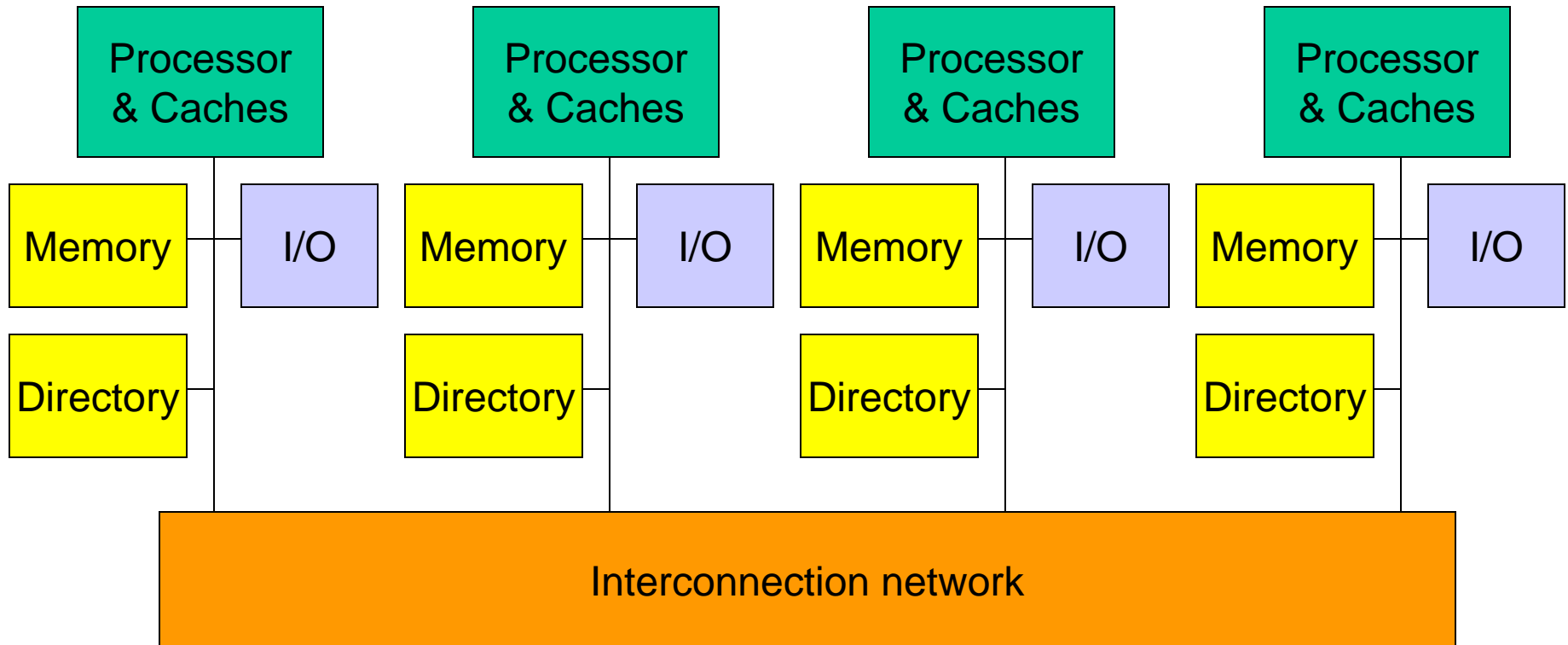
Coherence in Distributed Memory Multiprocs

- Distributed memory systems are typically larger → bus-based snooping may not work well
- Option 1: software-based mechanisms – message-passing systems or software-controlled cache coherence
- Option 2: hardware-based mechanisms – directory-based cache coherence

Directory-Based Cache Coherence

- The physical memory is distributed among all processors
- The directory is also distributed along with the corresponding memory
- The physical address is enough to determine the location of memory
- The (many) processing nodes are connected with a scalable interconnect (not a bus) – hence, messages are no longer broadcast, but routed from sender to receiver – since the processing nodes can no longer snoop, the directory keeps track of sharing state

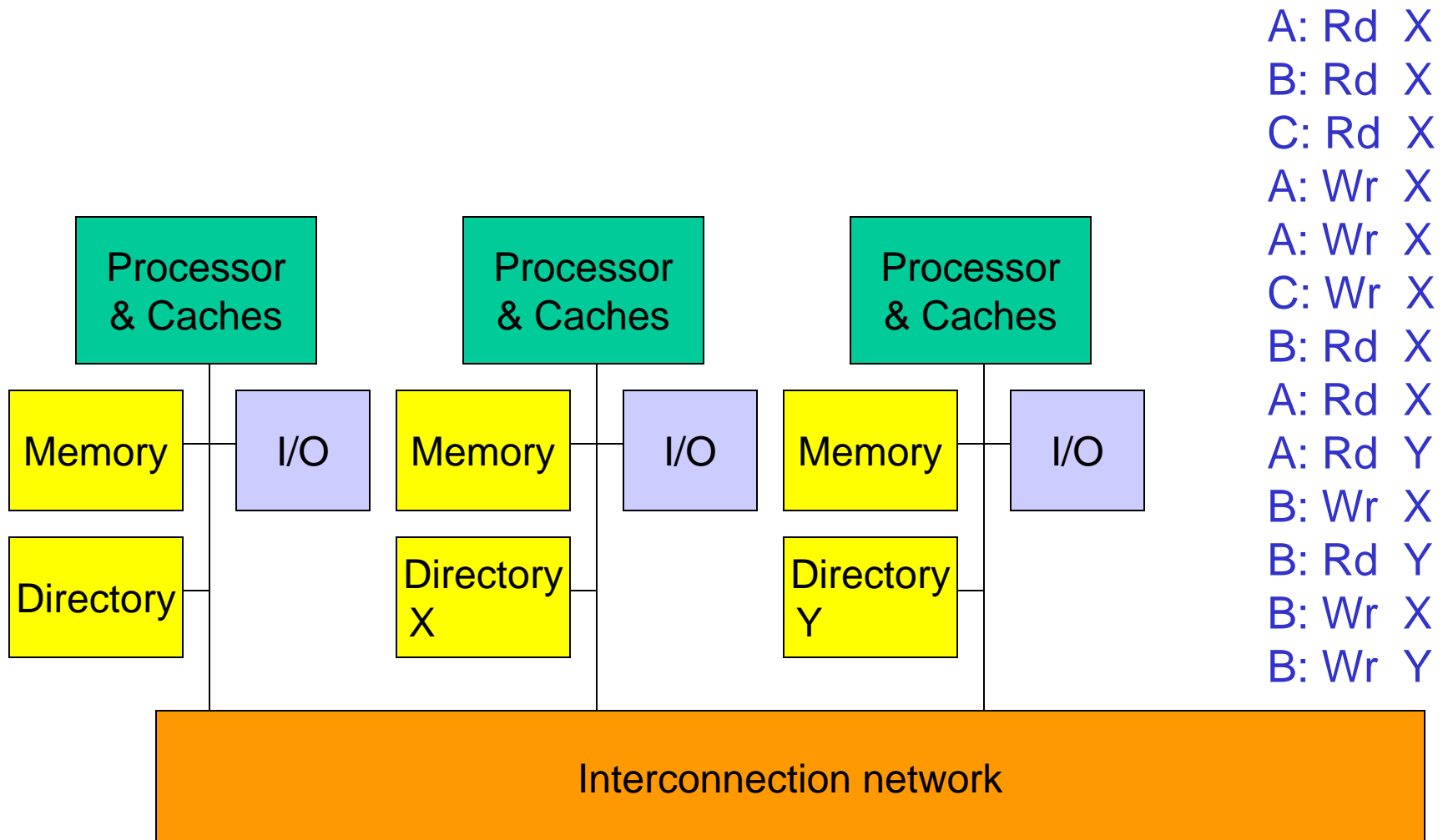
Distributed Memory Multiprocessors



Cache Block States

- What are the different states a block of memory can have within the directory?
- Note that we need information for each cache so that invalidate messages can be sent
- The block state is also stored in the cache for efficiency
- The directory now serves as the arbitrator: if multiple write attempts happen simultaneously, the directory determines the ordering

Directory-Based Example



Directory Actions

- If block is in uncached state:
 - Read miss: send data, make block shared
 - Write miss: send data, make block exclusive
- If block is in shared state:
 - Read miss: send data, add node to sharers list
 - Write miss: send data, invalidate sharers, make excl
- If block is in exclusive state:
 - Read miss: ask owner for data, write to memory, send data, make shared, add node to sharers list
 - Data write back: write to memory, make uncached
 - Write miss: ask owner for data, write to memory, send data, update identity of new owner, remain exclusive

Constructing Locks

- Applications have phases (consisting of many instructions) that must be executed atomically, without other parallel processes modifying the data
- A lock surrounding the data/code ensures that only one program can be in a critical section at a time
- The hardware must provide some basic primitives that allows us to construct locks with different properties
- Lock algorithms assume an underlying cache coherence mechanism – when a process updates a lock, other processes will eventually see the update

Synchronization

- The simplest hardware primitive that greatly facilitates synchronization implementations (locks, barriers, etc.) is an atomic read-modify-write
- Atomic exchange: swap contents of register and memory
- Special case of atomic exchange: test & set: transfer memory location into register and write 1 into memory
- lock: t&s register, location
 bnz register, lock
 CS
 st location, #0

Caching Locks

- Spin lock: to acquire a lock, a process may enter an infinite loop that keeps attempting a read-modify till it succeeds
- If the lock is in memory, there is heavy bus traffic → other processes make little forward progress
- Locks can be cached:
 - cache coherence ensures that a lock update is seen by other processors
 - the process that acquires the lock in exclusive state gets to update the lock first
 - spin on a local copy – the external bus sees little traffic

Coherence Traffic for a Lock

- If every process spins on an exchange, every exchange instruction will attempt a write → many invalidates and the locked value keeps changing ownership
- Hence, each process keeps reading the lock value – a read does not generate coherence traffic and every process spins on its locally cached copy
- When the lock owner releases the lock by writing a 0, other copies are invalidated, each spinning process generates a read miss, acquires a new copy, sees the 0, attempts an exchange (requires acquiring the block in exclusive state so the write can happen), first process to acquire the block in exclusive state acquires the lock, others keep spinning

Test-and-Test-and-Set

- lock: test register, location
bnz register, lock
t&s register, location
bnz register, lock
CS
st location, #0

Load-Linked and Store Conditional

- LL-SC is an implementation of atomic read-modify-write with very high flexibility
- LL: read a value and update a table indicating you have read this address, then perform any amount of computation
- SC: attempt to store a result into the same memory location, the store will succeed only if the table indicates that no other process attempted a store since the local LL (success only if the operation was “effectively” atomic)
- SC implementations may not generate bus traffic if the SC fails – hence, more efficient than test&test&set

Spin Lock with Low Coherence Traffic

```
lockit:  LL      R2, 0(R1)  ; load linked, generates no coherence traffic
         BNEZ   R2, lockit  ; not available, keep spinning
         DADDUI R2, R0, #1  ; put value 1 in R2
         SC     R2, 0(R1)  ; store-conditional succeeds if no one
                               ; updated the lock since the last LL
         BEQZ   R2, lockit  ; confirm that SC succeeded, else keep trying
```

- If there are i processes waiting for the lock, how many bus transactions happen?

Spin Lock with Low Coherence Traffic

```
lockit: LL      R2, 0(R1) ; load linked, generates no coherence traffic
        BNEZ   R2, lockit ; not available, keep spinning
        DADDUI R2, R0, #1 ; put value 1 in R2
        SC     R2, 0(R1) ; store-conditional succeeds if no one
                          ; updated the lock since the last LL
        BEQZ   R2, lockit ; confirm that SC succeeded, else keep trying
```

- If there are i processes waiting for the lock, how many bus transactions happen?
 - 1 write by the releaser + i read-miss requests + i responses + 1 write by acquirer + 0 ($i-1$ failed SCs) + $i-1$ read-miss requests

Further Reducing Bandwidth Needs

- Ticket lock: every arriving process atomically picks up a ticket and increments the ticket counter (with an LL-SC), the process then keeps checking the now-serving variable to see if its turn has arrived, after finishing its turn it increments the now-serving variable
- Array-Based lock: instead of using a “now-serving” variable, use a “now-serving” array and each process waits on a different variable – fair, low latency, low bandwidth, high scalability, but higher storage
- Queueing locks: the directory controller keeps track of the order in which requests arrived – when the lock is available, it is passed to the next in line (only one process sees the invalidate and update)

Title

- Bullet