

<p>Performance Metrics for Communication Mechanisms</p> <p>Communication bandwidth - limited by time for which resources within a node are tied up (occupancy), in addition to processor, memory, and interconnection bandwidths</p> <p>Communication latency - sender overhead + time of flight + transport latency + receiver overhead (overhead and occupancy closely related)</p> <p>Communication latency hiding - how well can the mechanism hide latency by overlapping communication with computation or other communication</p>	<p>Multiprocessor Systems Design: Coherence, Synchronization, and Consistency Issues</p> <p>Sandhya Dwarkadas Department of Computer Science University of Rochester</p>
<p>Advantages of Shared Memory</p> <p>Ease of programming when communication patterns are complex or vary dynamically during execution</p> <p>Better use of bandwidth when communicating small items (lower overhead) - protection implemented in hardware, communication implicit</p> <p>Use of hardware-controlled automatic caching with reduced latency and contention for accessing shared data</p>	<p>Flynn's Taxonomy of Parallel Architectures</p> <p>Single instruction stream, single data stream (SISD) - uniprocessor</p> <p>Single instruction stream, multiple data streams (SIMD)</p> <p>Multiple instruction streams, single data stream (MISD) - no existing commercial implementation</p> <p>Multiple instruction streams, multiple data streams (MIMD)</p>
<p>Advantages of Message Passing</p> <p>Simpler hardware - better scalability</p> <p>Explicit communication, focusing attention on its cost</p>	<p>MIMD Issues</p> <p>Data sharing - single versus multiple address spaces</p> <p>Process coordination - synchronization using lock or messages</p> <p>Distributed versus centralized memory</p> <p>Connectivity - single shared bus versus network with many different topologies</p> <p>Fault tolerance</p>
<p>Problems</p> <p>Amdahl's law - insufficient parallelism</p> $Speedup = \frac{1}{\frac{Fraction_{enhanced}}{Speedup_{enhanced}} + (1 - Fraction_{enhanced})}$ <p>High-latency communication</p> <p>Portability - knowledge of architecture required</p>	<p>MIMD Machines</p> <p>Classification based on memory organization</p> <ul style="list-style-type: none"> • Centralized shared memory architecture - uniform memory access (UMA) • distributed memory architectures <ul style="list-style-type: none"> - distributed (or scalable) shared memory architectures (NUMA) - Multi-computers <p>Models for communication</p> <ul style="list-style-type: none"> • shared memory • message passing

<p style="text-align: center;">Cache Coherence Protocols</p> <p>Snooping - sharing status broadcast to all copies</p> <p>Directory-based - sharing status of each block of data kept in one location</p>	<p style="text-align: center;">Message Passing Machines</p> <p>Multiple private address spaces - MIMD connected by a network (no lockstep operation)</p> <p>Interconnection used only for inter-processor communication</p> <p>Processes communicate via explicit messages (sends and recvs)</p> <p>Synchronization implicit in the message</p>
<p style="text-align: center;">Snoopy-Based Protocols</p> <p>Write Invalidate Protocol</p> <p>Write Update Protocol</p>	<p style="text-align: center;">Distributed Shared Memory Without Coherence</p> <p>Example - Cray T3D</p>
<p style="text-align: center;">A Simple Invalidate-Based Snoopy Protocol</p> <p>Finite State Machine with three states -</p> <ul style="list-style-type: none"> • Invalid • Shared/read-only (clean) • Modified/read-write (not shared) 	<p style="text-align: center;">Single Bus Shared Memory Architecture - MIMD</p>
<p style="text-align: center;">Mechanisms for Efficiency</p> <p>Dual cache tags</p> <p>Multi-level cache hierarchy with inclusion</p>	<p style="text-align: center;">Coherence Versus Consistency</p> <p>Cache coherence - ensuring that modifications made by a processor propagate to all (cached) copies of the data - mechanism used</p> <ul style="list-style-type: none"> • program order preserved • Writes to the same location by different processors serialized • defines the behavior of reads and writes to the same memory location <p>Consistency - defines when and in what order modifications are propagated to other processors</p> <ul style="list-style-type: none"> • defines the behavior of reads and writes with respect to accesses to other memory locations

Implementing Locks Using Coherence

On a SPARC: *ldstub* moves an unsigned byte into the destination register, and rewrites the same byte in memory to all 1s.

```
_Lock:
    ldstub [%o0], %o1
    addcc %g0, %o1, %g0
    bne _Lock
    nop
fin:   jmpl %r15+8, %g0
    nop
_UnLock:
    st %g0, [%o0]
    jmpl %r15+8, %g0
    nop
```

Directory-Based Cache Coherence Protocols

A directory keeps track of the state of every block that may be cached - sharing status in single location

Directory entries distributed along with the memory

States

- Shared
- Uncached
- Exclusive

How would you improve scalability?

Artificially delay processor - exponential backoff

Use a software queue

Example Directory Protocol

Possible messages

- Read miss - request data from home node
- Write miss - request data from home node
- Invalidate - invalidate a shared copy of data
- Fetch - fetch a block from remote cache to home node
- Fetch/invalidate - fetch and invalidate block at remote cache
- Data value reply - return data value from home node
- Data write back - write back data to home node

Hardware Primitives

Queueing locks instead of letting all processors try again

Assumption: directory-based coherence

Synchronization controller returns lock if free on a miss

Else, create a record of the node's request - let's the processor spin on locked value

When freed, controller updates lock variable in selected processor's cache

Basic Hardware Mechanisms for Synchronization

Must provide the ability to atomically read and modify a memory location

- Test-and-set - atomic exchange
- Fetch-and-increment - returns value and atomically increments it
- Load-locked/store conditional - pair of instructions - deduce atomicity if second instruction returns correct value

Memory Consistency Models

When must a processor see a value that has been updated by another processor?

```
P1:  A = 0;          P2:  B = 0;
    ...             ...
    A = 1;          B = 1;
L1:  if (B == 0) ... L2:  if (A == 0) ...
```

Using ll/sc for Atomic Exchange

Swap the content of R4 with the memory location specified by R1

```
try:  mov  R3, R4      ; mov exchange value
      ll   R2, 0(R1)  ; load linked
      sc  R3, 0(R1)  ; store conditional
      beqz R3, try    ; branch if store fails
      mov  R4, R2     ; put load value in R4
```

<p style="text-align: center;">Review</p> <p>A broad overview of multiprocessor systems design issues</p> <ul style="list-style-type: none">• Flynn's classification• shared memory vs. message passing• Snoopy-based versus directory-based coherence• Synchronization primitives• Consistency models	<p style="text-align: center;">Consistency Model Classification</p> <p>Models vary along the following dimensions</p> <ul style="list-style-type: none">• Local order - order of a processor's accesses as seen locally• Global order - order of a single processor's accesses as seen by each of the other processors• Interleaved order - order of interleaving of different processor's accesses on other processors
	<p style="text-align: center;">Sequential Consistency</p> <p>"A system is sequentially consistent if the result of any execution is the same as if the operations of all the processors were executed in some sequential order, and the operations of each individual processor appear in this sequence in the order specified by its program." [Lamport 79]</p>
	<p style="text-align: center;">Processor Consistency</p> <p>Also called Total Store Order (TSO - IBM 370, DEC VAX, SPARC)</p> <p>"A multiprocessor is said to be processor consistent if the result of any execution is the same as if the operations of each individual processor appear in the sequential order specified by its program." [Goodman 91]</p>
	<p style="text-align: center;">Release Consistency</p> <p>Implemented in DEC Alpha</p> <p>Assume data-race-free (synchronized) programs [Ghara. 90]</p> <p>Distinguish ordinary accesses and synchronization accesses</p> <p>Separate synchronization accesses into acquires and releases</p> <p>Guarantee consistency of ordinary data only after release ~→ architectural optimizations such as buffering to hid latency</p>