

ECE437: Introduction to Digital Computer Design

Chapter 5b (multicores)

Fall 2016

Outline

- Multiprocessors
 - Why?
 - Power changes everything
 - What about performance
 - Two processors of speed X instead of one processor of speed 2X
 - Challenges in harnessing parallelism
 - Is there a single vendor with uniprocessor products? (non-embedded)
- Key concepts
 - Programming models
 - Cache coherence
 - Synchronization
 - Consistency

ECE437, Fall 2016

(2)

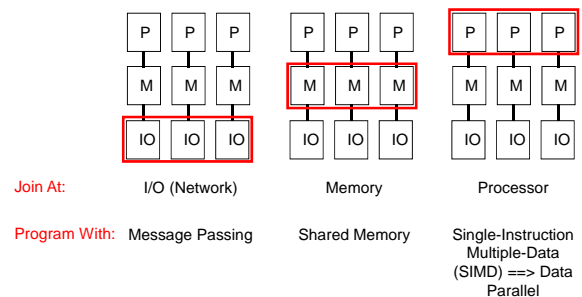
Why Multicores

- In the past, we got faster clock with every microprocessor generation
- But faster clock led to higher power which cannot be sustained
- Power = CV^2f where C is circuit capacitance, V is voltage, f is clock frequency
- V is linear in f , so power is cubic in f
- 2x clock \Rightarrow 8x power but 2x cores at same clock \Rightarrow 2x power
- We have been lowering voltage but that causes other problems (leakage, reliability)

ECE437, Fall 2016

(3)

Various Programming Views



- Our focus: Shared Memory
 - Most representative design point for multicores

ECE437, Fall 2016

(4)

Various Programming Views

In most interesting parallel applications

- parallel processes (threads) need to communicate
 - NO assumptions about threads' relative speeds
- communication method leads to another division
- message passing
 - shared memory
 - In message passing and shared memory - programmer writes explicitly programs
 - Shared memory is most general
 - In SIMD either parallelized by compiler or written by programmer (eg CUDA)
 - Limited applicability due to SIMD requirement but works very well if applicable

ECE437, Fall 2016

(5)

Shared memory

shared memory - all processors see one global memory

- programs use loads/stores to access data
- + conceptually compatible with uniprocessors
- + ease of programming even if communication complex/dynamic
- + lower latency communicating small data items
 - Real programs tend to communicate small fine-grain data
- + hardware controlled sharing allows automatic data motion
- The pluses/minuses on this and next slide are from programmer's point of view

ECE437, Fall 2016

(6)

Message passing

message passing - each processor sees only its private memory

- programs use sends/receives to access data
- + very simple hardware (almost nothing extra)
- + communication pattern explicit
 - but they MUST be explicit
- + least common denominator
- + shared memory MP can emulate message passing easily
- - (BIG MINUS) biggest programming burden: managing communication artifacts
 - Not bad for fixed simple communication patterns
 - Very hard for arbitrary, dynamic communication patterns found in real programs - databases, middleware,

ECE437, Fall 2016

(7)

Message Passing

- distribute data carefully to threads
 - no automatic data motion
- partition data if possible, replicate if not
 - replicate in s/w not automatic (so extra work to ensure ok)
- coalesce small mesgs into large, gather separate data into one mesg to send and scatter received mesg into separate data

ECE437, Fall 2016

(8)

Shared memory

- | | |
|--------------------------|-------------------------|
| • Thread1 | Thread2 |
| • | |
| • compute (data) | synchronize |
| • store(A,B, C, D, ...) | load (A, B, C, D,) |
| • synchronize | compute |
| • | |
- A B C D SAME in both threads- SINGLE shared memory

ECE437, Fall 2016

(9)

Message Passing

- | | |
|------------------------------|------------------------------|
| • Thread1 | Thread2 |
| • | |
| • compute (data) | receive (mesg) |
| • store (A,B, C, D ..) | scatter (mesg to A B C D ..) |
| • gather (A B C D into mesg) | load (A, B, C, D,) |
| • send (mesg) | compute |
| • | |
| • | |
- A B C D are DIFFERENT in each thread -- PRIVATE memory

ECE437, Fall 2016

(10)

Example

- Ocean - Compute temperature of a part of an ocean
- Very simple example so you can see
- Real parallel programs have much more complex communication patterns - databases, online transaction processing (reservations), decision support (inventory control), web servers, Java middleware
 - in other words the programs that make the world go round!

ECE437, Fall 2016

(11)

Sequential Ocean

```

main()
begin
  read(n);
  A = malloc(n * n);
  initialize(A);
  Solve(A);
end main

Solve(float **A)
begin
  while (!done)
    diff = 0;
    for i=1 to n do
      for j = 1 to n do
        temp = A(i,j);
        A(i,j)=0.2*(A(i,j)+A(i,j-1)+
                    A(i,j+1)+A(i+1,j)+
                    A(i-1,j));
        diff += abs(A(i,j) - temp);
      end for
    end for
    if (diff / (n*n) < TOL) then done = 1;
  end while
end Solve
    
```

ECE437, Fall 2016

(12)

Shared Memory Ocean

```

main()
begin
  p = NUM_PROCS();
  pid = MY_PROC();
  read(n);
  A = G_MALLOC(n * n);
  initialize(A);
  CREATE(p);
  Solve(A);
  WAIT_FOR_END(p-1);
end main

Solve(float **A)
begin
  start_row = 1 + (pid * n/p);
  end_row = start_row + n/p - 1;
  while (!done)
    mydiff = diff = 0;
    for i=start_row to end_row do
      for j = 1 to n do
        temp = A(i,j);
        A(i,j)=0.2*(A(i,j)+A(i,j-1)+ A(i+1,j)+
                    A(i-1,j));
        mydiff += abs(A(i,j) - temp);
      end for
    end for
    LOCK(dlock); diff += mydiff; UNLOCK(dlock);
    BARRIER();
    if (diff / (n*n) < TOL) then done = 1;
    BARRIER();
  end while
end Solve

```

ECE437, Fall 2016

(13)

Message Passing Ocean

```

main()
begin
  p = NUM_PROCS();
  read(n);
  CREATE(p-1, solve);
  solve(A);
  WAIT_FOR_END(p-1);
end main

Solve(float **A)
begin
  myA = malloc(n/p + 2 by n array);
  while (!done)
    mydiff = 0;
    send (myA 1 row); send (myA n/p row);
    receive (myA 0 row); receive (myA n/p+1 row);
    for i= 1 to mymax row do
      for j = 1 to n do
        temp = myA(i,j);
        myA(i,j)=0.2*(myA(i,j)+myA(i,j-1)+myA(i-1,j)
                    +myA(i,j+1) + myA(i+1, j));
        mydiff += abs(myA(i,j) - temp);
      end for
    end for
  end for
end Solve

```

ECE437, Fall 2016

(14)

Message Passing Ocean

```

If (pid != 0)
  send (mydiff)
  receive (done);
else
  for (p-1 times)
    receive (tempdiff);
    mydiff += tempdiff;
  end for
  if (mydiff/n*n < TOL) done = 1;
  for (p-1 procs)
    send (done);
  end for
end if
end while

```

ECE437, Fall 2016

(15)

Ocean

- Think of how each core's cache would look in the shared memory case
- Think of how each core's messages would look in message passing case

ECE437, Fall 2016

(16)

SIMD

- Single instruction multiple data
 - Shared memory/msg passing are MIMD -
- One instruction operates on multiple data - less instruction overhead
- Deep, fast pipelines
- No hazards
 - Parallelism guaranteed by compiler/programmer
- No need to check for hazards

ECE437, Fall 2016

(17)

SIMD (Vector)

- do i = 1,64 - $a * x + y$ (daxpy)
 - $y[i] = a * x[i] + y[i]$
- ld f0, a
- lv, v1, rx 64-element vector
- multsv v2, f0, v1
- lv v3, ry
- addv v4, v2, v3
- sv ry, v4
- Above is vector example
- Intel MMX similar - except narrower vector

ECE437, Fall 2016

(18)

Common Parallelization Strategies

- Think of parallelism at the algorithm level
- Possibly different
 - Best parallel algorithm
 - Best parallelization of best sequential algorithm
- Think of partitioning
 - Tasks? (do webpage serving and data-base lookup for different requests in parallel.)
 - Data? (Do similar operations on different data in parallel.)

ECE437, Fall 2016

(19)

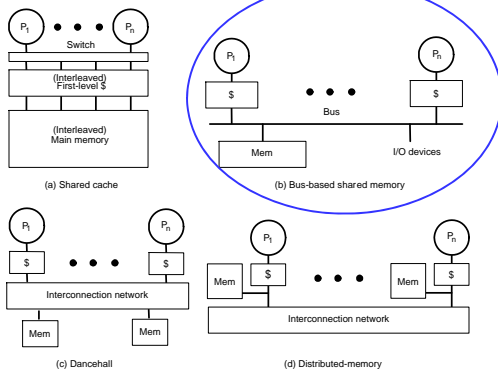
Programming models

- Pthreads - Shared memory processors **
- MPI (mesg passing interface) - clusters **
- OpenMP - Shared memory
 - Marking loops as parallel loops
- Higher level primitives
 - Streaming, MapReduce, Parallel Recursion (Cilk C)
 - Customized to architecture
 - CUDA for nVidia GPUs

ECE437, Fall 2016

(20)

Shared address space MPs



ECE437, Fall 2016

(21)

AMD and Intel quad-core designs

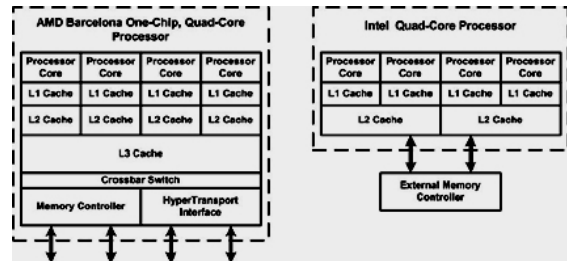


Figure © Embedded.com

ECE437, Fall 2016

(22)

Terminology: Process vs. Thread

- "Heavyweight" **processes**
 - Different Thread of control
 - PC, registers, stack
 - Different Address space (page table, heap)
- "Lightweight" processes, a.k.a. "**threads**"
 - Different PC, register values, and stack allocation
 - "Context"
 - Same address space (thus same page table, same heap)
- Shared regions across heavyweight processes possible.

ECE437, Fall 2016

(23)

How do threads communicate?

- Our focus:
 - Shared memory, bus-based systems
 - Shared address space
 - Memory does not have to physically shared
 - Communication via loads and stores

ECE437, Fall 2016

(24)

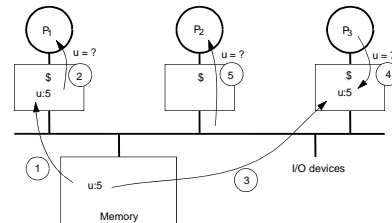
Several Interesting Problems

- Does everything work intuitively***
 - No.
 - Reason about coherence/consistency
- Sequencing
 - Synchronization
- Today: Cache Coherence

ECE437, Fall 2016

(25)

Example Cache Coherence Problem



- Is there a problem?
 - Replication + writes = coherence problem

ECE437, Fall 2016

(26)

Cache coherence

- Previous slide makes intuitive sense but EXACTLY what should coherence guarantee?
- Is the previous slide the only way to do business? Are there many options in what coherence can and cannot do?
- What are the rules of the game?

ECE437, Fall 2016

(27)

Coherence

- $A = 0$ (initial)
- Thread 1 thread 2 thread 3
-
- $A = 10$
- .. $A = 5$...
- read A
- What should thread3 read?
 - No relative speeds assumed
- Does this mean 'anything goes' in shared memory?

ECE437, Fall 2016

(28)

Coherence

- A = 0 (initial)
- Thread 1 thread 2 thread 3
-
- A = 10 ..
- .. Read A
- .. read A
- What should thread2 read?
- What should thread3 read if thread2 got 10?
- Does this mean 'anything goes' in shared memory?

ECE437, Fall 2016

(29)

Coherence

- A = 0 (initial)
- Thread 1 thread 2 thread 3
-
- A = 10 ..
- .. Read A wait thread2
- .. Signal thread 3 read A
- What should thread2 read?
- What should thread3 read if thread2 got 10?
 - If thread2 reads 0 then thread3 can read 0 or 10
- Do you still think 'anything goes' in shared memory?

ECE437, Fall 2016

(30)

Coherence

- Hard to distinguish between previous two slides
- Hard to track arbitrary synchronization among arbitrary threads in a large distributed system
 - Previous slide threads 2 and 3 synchronization affects thread1's write - hard for thread1 to track threads 2&3
- → the sanest option is that whenever a write is done it is done for ALL threads
- Whenever thread1's write to A is done it should be done for ALL threads
 - Ie the previous value is not visible to anyone
 - Called **write atomicity** and is provided by cache coherence
 - Does this mean thread3's read will ALWAYS return 5 in the FIRST eg? Thread2's read will ALWAYS return 10 in SECOND eg? Thread3 will ALWAYS get 10 in SECOND eg?
 - what does this mean?

ECE437, Fall 2016

(31)

Coherence

- Fuzzy idea
 - Necessary for correctness
- Precise definition
- For any memory location X
 - [within one thread - easy] Operations issued by any one thread **occur** in the order in which they issued
 - [across threads - hard] Value returned by read is the value written by **last** write (but "last" is relative!)
- Derivative properties
 - Write propagation : writes become visible to everyone (**when?**)
 - Write serialization/atomicity: writes to a single location seen in same order by **every** processor OR writes occur for everyone or no one (ie indivisibly)
- If you understand the above you will understand why coherence does what it does

ECE437, Fall 2016

(32)

Coherence

- cache coherence suggests absolute time scale
 - not necessary
- what is required is appearance of coherence
 - not instantaneous coherence
- Bart Simpson's famous words -
 - "nobody saw me do it so I didn't do it!"
- e.g. temporary incoherence between
 - writeback cache and memory ok
- Remember that no assumptions about threads' relative speed

ECE437, Fall 2016

(33)

Coherence

- The ONLY thing coherence provides is write atomicity
 - ie when a write happens, it happens for all so nobody can see the previous value
 - even this suggests instantaneous and global update of all copies but cannot be implemented that way
 - writes take a non-zero time window to complete (ie not instantaneous) and after this window nobody can see the old value and during this window any access to the new value is blocked so nobody can see the new value WHILE the old value is not all gone
 - But during the window the old value may be visible because the write reaches different cores at different times

ECE437, Fall 2016

(34)

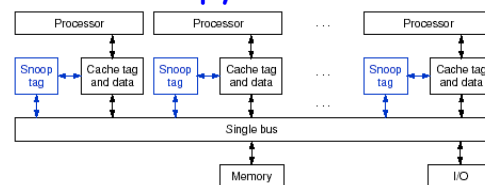
How to enforce coherence

- Correct actions:
 - On write, **invalidate** other copies
 - On write, **update** other copies
- Snooping:
 - Hits are guaranteed "not stale" → misses have to get latest data
 - Every node can "see" the bus
 - If all transactions can be observed
 - Each can maintain coherence

ECE437, Fall 2016

(35)

Snoopy coherence



- Hits are guaranteed "not stale" → misses have to get latest data
- ALL cache misses go on the bus which snoops into ALL the other caches (ie look at tags)
- Cache controller updates blocks' state for processor/snoop events and generates bus actions (**Ignore "snoop tag" for now**)

ECE437, Fall 2016

(36)

Snoopy Design Choices

- Snoopy protocol
 - set of states
 - state-transition diagram
 - actions
- Basic Choices
 - write-through vs. **write-back**
 - **invalidate** vs. update
- Multiple reader single writer protocol

ECE437, Fall 2016

(37)

A 3-State Write-Back Invalidation Protocol

- 3-State Protocol (MSI) - KEY invariants
 - **M**odified
 - one cache has valid/latest copy
 - memory is stale
 - **S**hared
 - one or more caches have valid copy
 - Memory is up-to-date
 - **I**nvalid
- Must invalidate ALL other copies before entering modified state (multi reader single writer)
- Requires bus transaction to invalidate
- Show previous 'u' block eg

ECE437, Fall 2016

(38)

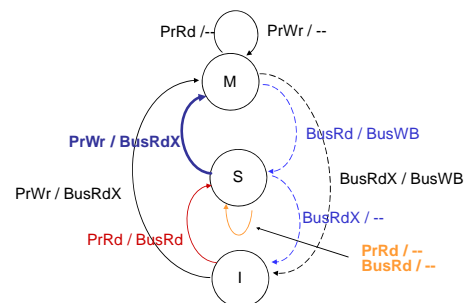
MSI Processor and Bus Actions

- Processor:
 - PrRd
 - PrWr
 - Writeback on replacement of modified block
- Bus
 - PrRD miss → Bus Read (**BusRd**) Read **without** intent to modify, data could come from memory or another cache
 - PrWr miss → Bus Read-Exclusive (**BusRdX**) Read **with** intent to modify, must invalidate all other caches copies
 - Writeback (**BusWB**) cache controller puts contents on bus and memory is updated (no atomicity issues)
 - Definition: **cache-to-cache transfer** occurs when another cache satisfies BusRd or BusRdX request
- Let's draw it!

ECE437, Fall 2016

(39)

MSI State Diagram



- A/B means action A causes transaction B on the bus
- BusWB does TWO things TOGETHER
 - requestor gets block (cache-to-cache) and memory is updated

ECE437, Fall 2016

(40)

An example

Proc Action	P1 State	P2 state	P3 state	Bus Act	Data from
1. P1 read u	S	--	--	BusRd	Memory
2. P3 read u	S	--	S	BusRd	Memory
3. P3 write u	I	--	M	BusRdX	Memory (?)
4. P1 read u	S	--	S	BusRd	P3's cache
5. P2 read u	S	S	S	BusRd	Memory

- **Single writer, multiple reader protocol**
- **Why Modified to Shared?**
 - Why not to Invalid?
 - Give up as little as possible. Retain option to read.
 - Why throw away write permission?
 - Ownership
- **What if not in any cache?**
 - Read, followed by Write produces 2 bus transactions!

ECE437, Fall 2016

(41)

MSI is a distributed FSM

- The slide shows one state machine but really there is a state machine per block per cache
 - Of course the state machines in all the caches are identical (but the state may not be identical)
 - States of different blocks obviously different
 - But state of same block in different caches may or may not be identical
- Though the slide shows one state machine, on a bus snoop, **MULTIPLE** caches with the same block transit from one state to another
 - Eg upon busrdx, ALL copies other than the requestor transit to I and requestor goes to M

ECE437, Fall 2016

(42)

What happens where (relevant for lab)

EACH block's tag entry in EACH cache has the MSI state bits

- In each cache, only one block accessed at a time by either the core or bus, so only one set of combinational logic to implement MSI state machine

Cache controller (for EACH cache) and bus controller (one) and memory controller (one) - ALL SEPARATE (merging these will cause immense suffering with one exception)

On the bus, BusRd or BusRX have an accompanying reply (the requested block), BUSWB has no reply

Sequence of operations

1. Core request ld/st ->
2. cache hit/miss (cache controller) ->
3. bus request (to bus controller) ->
4. grant & bus transaction (BusRd or BusRdX) (on bus) ->
5. snoop into other caches (each cache controller) ->
6. look up MSI state in EACH cache in parallel (each cache controller) ->

ECE437, Fall 2016

(43)

What happens where (relevant for lab)

7. perform needed MSI state transition and update state (each cache controller) ->
8. access memory if needed (memory controller) ->
9. either memory or a cache place block (BusWB) on bus ->
10. requestor cache places block in cache, updates its MSI state, and complete ld/st (cache controller) ->
11. data returned to core, if any

Till the reply comes, the requestor cache (and core) simply waits

In step 9 if a cache supplies the block (Modified) via a BusWB then in step 10 the requestor cache gets block (cache-to-cache) AND memory is updated (occurs TOGETHER in ONE BusWB)

ECE437, Fall 2016

(44)

What happens where (relevant for lab)

Bus operation broken into multiple cycles to maintain fast bus clock

- DO NOT put all of previous slide into one bus cycle - kills bus clock speed

5 bus cycles-

1. Bus req (bus controller) - step 3 on previous slide
2. arbitrate and bus grant (bus controller) - step 4 prev. slide
3. snoop and state transit (ALL cache controllers) - steps 5-7
4. reply block (ONE cache or memory) - steps 8-9 prev. slide
 - If memory is accessed (via memory controller), many wait cycles on bus for slow memory
5. requestor complete and relinquish bus (ONE cache controller) - step 10 prev. slide

Pipeline-cache interaction

- MEM stage goes to cache controller (as before), upon a miss (ie need a bus transaction), cache controller requests bus controller (instead of memory controller, as in uncore) and waits for reply

ECE437, Fall 2016

(45)

What happens where (relevant for lab)

Though multiple cycles, NO OTHER bus request allowed until previous is done - ie atomic bus

- ie ONLY ONE bus transaction at a time - no pipelining
 - This is an exception - in atomic bus, it MAY be ok to merge the bus and memory controllers but real buses are not atomic and don't do this merge
- Remember: Snoops/going to the bus ONLY on cache misses (writes to S blocks are misses though data is present but not have permission to write) - cache hits proceed at full speed without going to the bus

ECE437, Fall 2016

(46)

BusWB

- BusWB is also for writing back replaced M blocks
 - BusWB can also occur as a reply to BusRD or BusRDX
- There is no snooping
- WBs are not writes - there is no atomicity requirement - why?
- You just write back the block to memory

ECE437, Fall 2016

(47)

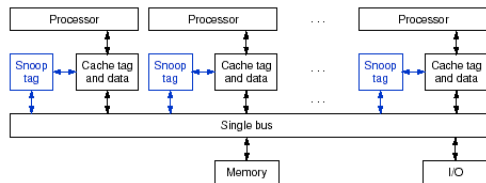
Writes and invalidations

- So much talk about invalidations, but invalidation means many cache misses
 - Then why bother with coherence or even caches if you are going to miss all the time
 - Whoa!
- Think about shared memory Ocean - would the above occur in Ocean?

ECE437, Fall 2016

(48)

Duplicate tag



- If only one tag array then both processor and bus snoop contend for it - structural hazard
- → use a duplicate tag array so one tag array serves the processor and other serves snoops
 - Only tag array copies but not data array - why not?

ECE437, Fall 2016

(49)

Duplicate tag

- But the tag copies have to be updated together
 - any state transitions due to processor AND bus have to update both tag copies
 - Eg processor reads ($I \rightarrow s$) or writes ($I \rightarrow M$ or $S \rightarrow M$)
 - Bus invalidations
- Wait a minute - if both copies have to be updated then won't we have structural hazards all over again?
 - Tag updates are infrequent - why?

ECE437, Fall 2016

(50)

Duplicate tag

- While both tag copies are being updated, a bus snoop can't occur (snoop tag busy)
- Then extend the bus transaction's snoop (bus clk #3) from one cycle to multiple cycles until all caches' snoop tags are free (may have to block processor request to avoid snoop waiting for ever)
 - The bus transaction goes from 5 cycles to more
 - Like adding "bus wait cycles" in 362
- If memory supplies a block, bus transaction's "reply" has to be extended from one cycle to many cycles, anyways (memory is slow)

ECE437, Fall 2016

(51)

Connect to previous egs

- Write atomicity
- Thread 1 thread2 thread3
 - Wr A rd A rd A
- Coherence does not mean instantaneous - wr A will take time to reach thread2 &3 and may not reach both threads at the same time even for snoopy bus coherence

ECE437, Fall 2016

(52)

Connect to previous eggs

- The write will reach thread2 and 3 at different times (may be thread3 tag is busy) so thread3 could read old A for some time even if thread2 is invalidated (of course thread2 cannot read new A until write is done when all invalidations completes - the "non-zero time window" in long-ago slide)
- So write atomicity is not a straight line through the threads in physical time
 - It bends and curves through the threads
 - It is a straight line in virtual time
 - See previous slide
 - Profound!

ECE437, Fall 2016

(53)

Atomic bus

- Coherence makes writes atomic
 - Not instantaneous but atomic
 - Hits guaranteed to have correct data
 - Misses go thru snoopy bus
 - But if two misses overlap write atomicity can be violated
 - Eg 2 write misses to same block
 - Make the bus atomic - ie only one miss at a time
 - Each miss has to be complete before next miss starts
 - → if any memory access then bus is blocked for whole access (request to response)
 - Bus becomes bottleneck
 - Bus arbiter ensures one at a time - first come first served
 - What order are writes seen by everyone (everyone sees the same order)
 - In the order in which they get to the bus

ECE437, Fall 2016

(54)

Atomic bus

- Real systems don't use atomic buses
 - Either use a "split transaction bus" which splits a request and its response → other requests/responses can come in between
 - write atomicity complicated
 - Many small-scale multiprocessors - Sun Wildfire
 - Recent multicores - Intel Sandy Bridge, Ivy Bridge AMD Bulldozer
 - OR use no buses at all - called distributed shared memory where write atomicity is quite challenging
 - Large-scale shared memory multiprocessors
 - SGI Origin
 - ASCI Red, White, Blue machines from Intel, IBM, Cray

ECE437, Fall 2016

(55)

Invalidate vs. Update

- Pattern 1:


```
for i = 1 to k
    P1(write, x);           // one write before many reads
    P2--PN-1(read, x);
end for i
```
- Pattern 2:


```
for i = 1 to k
    for j = 1 to m
        P1(write, x); // many writes before one read
    end for j
    P2(read, x);
end for i
```

ECE437, Fall 2016

(56)

Invalidate vs. Update, cont.

- Pattern 1 (one write before reads)
 - $N = 16, M = 10, K = 10$
 - Update
 - Iteration 1: N regular cache misses (70 bytes)
 - Remaining iterations: update per iteration (14 bytes; 6 ctrl, 8 data)
 - Total Update Traffic = $16 \times 70 + 9 \times 14 = 1246$ bytes
 - Invalidate
 - Iteration 1: N regular cache misses (70 bytes)
 - Remaining: P1 generates upgrade (6), 15 others Read miss (70)
 - Total Invalidate Traffic = $16 \times 70 + 9 \times 6 + 15 \times 9 \times 70 = 10,624$ bytes
- Pattern 2 (many writes before reads)
 - Update = 1400 bytes
 - Invalidate = 824 bytes
- No clear winner PLUS in updates write atomicity is hard in distrib. Sh. Mem - hence not used in practice

ECE437, Fall 2016

(57)

Coherence: Performance Issues

- Access Pattern specific optimizations
- Update vs. Invalidate
- Cache Block size
 - Remember 3C classification
 - 4th C: "Coherence miss"
 - True vs. False Sharing
- Read Snarfing -
 - if tag match but invalid, "snarf" from bus
- Non-Atomic Bus

ECE437, Fall 2016

(58)

Why Synchronize?

lw \$r1, ac-balance	lw \$r1, ac-balance
add \$r1, \$r1, 100	add \$r1, \$r1, -50
sw \$r1, ac-balance	sw \$r1, ac-balance

- Two concurrent operations on the same bank account
 - Deposit \$100 check
 - Withdraw \$50 cash
- Correct result
 - $ac_balance(new) = ac_balance(old) + 50$
- What can go wrong?

ECE437, Fall 2016

(59)

Synchronization

- Solution: Locks/Mutex
 - "critical section"
 - Only one person can have a lock
- Other synchronization primitives
- Barriers
 - Wait for everyone before proceeding
 - Can we just use a counter?
- Event notification
 - Producer consumer
 - Just use flags?
 - Consistency issues

LOCK ac-lock

lw \$r1, ac-balance
add \$r1, \$r1, 100
sw \$r1, ac-balance
UNLOCK ac-lock

LOCK ac-lock

lw \$r1, ac-balance
add \$r1, \$r1, -50
sw \$r1, ac-balance
UNLOCK ac-lock

ECE437, Fall 2016

(60)

How Not To Implement Locks

- LOCK
while(lock_variable == 1); /* locked so wait */
lock_variable = 1;
- UNLOCK
lock_variable = 0;
- Implementation requires atomicity!
 - read and write of lock_variable must be atomic
 - Else two processes can acquire lock at same time
 - Another thread can come between read & write
 - Unlock is easy

ECE437, Fall 2016

(61)

Atomic Read-Modify-Write Operations

- Test&Set(r,x) -- r gets previous lock value
r = m[x] • r is register
m[x] = 1 • m[x] is memory location x
 - if r is 1 then already locked else accessor gets the lock
- Swap(r,x) --- r gets previous lock value
r = m[x], m[x] = r
- Compare&Swap(r1,r2,x) -- r2 gets previous lock value
if (r1 == m[x]) then
r2 = m[x], m[x] = r2
- Fetch&Op(r,x,op) -- r gets previous lock value
r = m[x], m[x] = op(m[x])
- In ALL cases, the operations have to be ATOMIC in hardware

ECE437, Fall 2016

(62)

R-M-W primitives

- Test&set (and others) does a strange thing: it ALWAYS writes a 1 whether a thread gets the lock or not
- Getting the lock is determined by the test&set return value which is the PREVIOUS value of the lock
 - If already locked t&s writes a 1 on top of 1 (previous value 1 which is returned)
 - If not t&s writes a 1 on top of 0 (previous value 0 which is returned)

ECE437, Fall 2016

(63)

Correct lock implementation

LOCK
while (test&set(x) == 1);

UNLOCK
x = 0;

ECE437, Fall 2016

(64)

LL-SC based locks

- All primitives on previous slides need atomic Read-Modify-Write support in h/w
- Cannot easily do atomic Read-modify-write in h/w
 - R-M-W is reading followed by writing
 - one way would be to read (load) and then "negative acknowledge (nack)" invalidations from others so their write cannot complete while you try to write to complete your R-M-W
 - But what happens if two threads load at the same time?
 - They will both read and then permanently nack each other
 - Called a DEADLOCK!
 - Chances may be low but bad if it happens
 - Rare case has to be correct and it is not so rare if lock is heavily contended
 - Harder to do in a pipelined bus/distributed system

ECE437, Fall 2016

(65)

LL-SC based locks

- Different approach
 - Instead of forcing a core's R-M-W to be atomic, detect non-atomicity between this core's R and W (of lock)
 - if detected, retry R-M-W else see if got the lock
 - If did not get the lock, retry getting lock
 - How to detect non-atomicity
 - i.e., if some other core got the lock BETWEEN this core's R and W of lock
 - lock = 0 if unlocked and 1 if locked → to get the lock, a core MUST write to lock
 - Writes invalidate in coherence
 - So if another core invalidates lock between this core's R and W of lock → non-atomicity

ECE437, Fall 2016

(66)

LL-SC based locks

- Invalidations go to cache (s/w not involved) but we want s/w to know if lock invalidated between R and W
- So need special instrs to catch such invalidations
- Instead of normal ld for R, use load-linked (ll) and instead of normal st for W, use store-conditional (sc)
- ll reads the lock as usual but also puts its address in a "link register" (next to cache - NOT CPU register)
- And sc first checks if the link register has not been invalidated (ie zeroed) and equals sc's address
- if so means no invalidation between R and W and sc simply stores; else atomicity violation between this core's R and W, so sc does NOT store (ie fails) and instead sets its source (CPU) register to 0 to signal failure (ie non-atomicity)

ECE437, Fall 2016

(67)

LL-SC based locks

- Link register is "implicit" in ll and sc, so not specified
- ```

Test&set: move $s2, 1 // 1 means locked
 ll $s1, X // $s1 ← M[X] and sets link reg ← X
 sc X,$s2 // $s2 contains 1 to signify locked
 // if link reg == X then
 // M[X] ← $s2 and $s2 ← 1
 // else sc fails and sets $s2 ← 0
 beq $s2, $0 T&S // retry if atomicity failure
 jr $31 // return $s1 which is previous lock
 // value

Unlock: sw $0, X // easy!! HUMBLE sw -- not fancy sc
 jr $31

```
- The returned previous lock value may already be locked in which case retry test&set in lock's "while loop" - this retry is not due to atomicity failure but due to not getting the lock

ECE437, Fall 2016

(68)

## Performance of Test & Set

### LOCK

```
while (test&set(x) == 1);
```

### UNLOCK

```
x = 0;
```

- High **contention** (many threads want lock while one thread holds it)
- All non-holders keep retrying T&S but each retry writes a 1 on top of 1 (sc within T&S is a write)
- Remember the **CACHE!**
- Each test&set writes and invalidates others - LOTS of cache misses while non-holders wait
- Several Optimizations

ECE437, Fall 2016

(69)

## Better Lock Implementations

- Two choices:
  - Don't execute test&set so much
  - Spin without generating bus traffic
- Test&Set with Backoff
  - Insert delay between test&set operations (not too long)
  - Exponential seems good ( $k^i c^i$ )
  - Not fair (a later contender may get lock before earlier)
- Test-and-Test&Set
  - **Intuition: No point in setting the lock with a 1 on top of 1 until we know that it's unlocked**
  - First keep testing (ie repeated NORMAL loads) until test succeeds (ie unlocked) then try test&set (for atomicity)
  - Testing (loading) gets cache hits without invalidations/misses until lock is unlocked (a store ) which will invalidate all contenders and then one test&set will succeed
  - So no cache misses while non-holders wait
  - Still contention at unlock
  - Still not fair

ECE437, Fall 2016

(70)

## Synchronization: Performance Issues

- Granularity
  - Lock single element vs. Lock whole array
  - Locking overhead vs. loss of concurrency
- Contention
  - TAS good at low contention
- Fairness
  - TAS offers no guarantees of fairness

ECE437, Fall 2016

(71)

## Synchronization

- Deadlocks: Huge Correctness Issue
  - Previous deadlock was hardware deadlock here it is software deadlocks
- Races
- May not be visible to the programmer
  - Hard to debug (irreproducible)
- "Locks don't compose with themselves"
  - Other synchronization mechanisms have other problems
- **We live with this problem**
  - Find a solution, become famous, rich.

ECE437, Fall 2016

(72)

## Coherence vs. Consistency

- Intuition says loads should return latest value
  - what is latest?
- **Coherence concerns only one memory location**
- **Consistency concerns apparent ordering for all locations**
- A Memory System is Coherent if
  - can serialize all operations to **that location** such that,
  - operations performed by any processor appear in program order
    - program order = order defined by program text or assembly code
  - value returned a read is value written by last store to that location

ECE437, Fall 2016

(73)

## Why Coherence != Consistency

/\* initial A = flag = 0 \*/

|           |                             |
|-----------|-----------------------------|
| <u>P1</u> | <u>P2</u>                   |
| A = 1;    | while (flag == 0); /*wait*/ |
| flag = 1; | print A;                    |

Intuition says print A = 1

- Intuition is based on an implicit assumption

Coherence doesn't say anything, why?

ECE437, Fall 2016

(74)

## Why Coherence != Consistency

/\* initial A = flag = 0 \*/

|           |                             |
|-----------|-----------------------------|
| <u>P1</u> | <u>P2</u>                   |
| A = 1;    | while (flag == 0); /*wait*/ |
| flag = 1; | print A;                    |

Consider coalescing write buffer

- If multiple writes to same location, the buffer coalesces so only one write goes to memory
- A and flag written atomically - coherence correctly invalidates other copies
- BUT flag = 1 coalesced with previous write to flag so flag and A reordered in buffer so flag updates memory before A → print A = 0!

ECE437, Fall 2016

(75)

## Why Coherence != Consistency

- Coherence concerns ONE location from ALL threads
- Consistency concerns ordering of ALL locations from ONE thread
- Coherence says WHAT action has to happen
- Consistency says WHEN action has to be complete
- Related through implementation but not the same

ECE437, Fall 2016

(76)

## Consistency affected by s/w too

- H/w is not the only culprit - s/w too
- Register Allocation can break SC

```

a = f(x) while (a == 0);
 b = g(a);

lw $arg1, x lw $1, a
jal f loop: beqz $1, $0, loop
sw a, $retval lw $arg1, $1
 jal g

```

Diagram illustrating a scenario where register allocation can break Sequential Consistency (SC). The code shows two parallel execution paths. The left path (P1) consists of `lw $arg1, x`, `jal f`, and `sw a, $retval`. The right path (P2) consists of `while (a == 0);`, `b = g(a);`, and a loop containing `beqz $1, $0, loop`, `lw $arg1, $1`, and `jal g`. A dashed line separates the two paths. Callouts 1, 2, and 3 highlight specific instructions: 1 points to `b = g(a);`, 2 points to `sw a, $retval`, and 3 points to `lw $arg1, $1` inside the loop.

- Use "volatile" declaration to prevent register allocation
  - `volatile int a;`

ECE437, Fall 2016

(77)

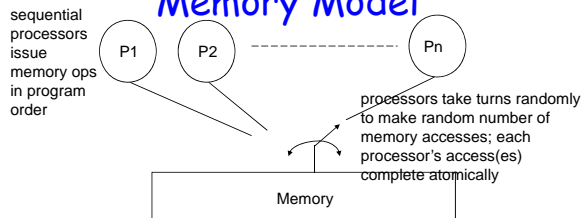
## Consistency Model

- Is a contract between the system and the programmer
  - The system includes hardware, compiler, runtime, OS, etc
- Consistency models says what behavior is to be expected
- Programmer programs based on that behavior
- Sequential consistency is the simplest/most intuitive such behavior/model

ECE437, Fall 2016

(78)

## Sequential Consistency (SC) Memory Model



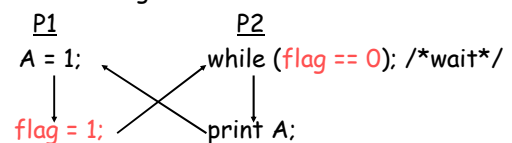
- Accesses from ONE thread in program order and done atomically one at a time
  - No relative speeds assumed
  - Each access is complete before next starts
  - Intuition assumes this - hence `print A = 1`

ECE437, Fall 2016

(79)

## Violation of SC

`/* initial A = flag = 0 */`



- Print 0 means a cycle which is impossible in SC
- ANY acyclic interleaving is ok but cyclic is not
  - SC does not say which acyclic interleaving

ECE437, Fall 2016

(80)

## Sequential Consistency

- Programmers think in SC
  - i.e. preserve memory order
- Performance maximized when memory order is flexible
  - i.e., do not preserve memory order
- Architect's solution: Don't preserve order
  - Our problem is now the programmer's problem
  - Done today but causes programmer grief
- Can we get intuitive programming, correctness and fast performance?
- Several interesting approaches proposed
  - Parallel programming remains harder than sequential programming.
- Come to grad school (666)

ECE437, Fall 2016

(81)

## Announcement

- DON'T GIVE UP!!
- Ask for help if the lab or lectures don't make sense
- TAs and I are here to help
- Come to lectures even if much lab work remains - your work will reduce by 10x
  - Listen/understand/internalize IN the lectures and not "go back to slides later" - your work cut by 10x
- We have simplified all lab workload
  - All the interfaces are pre-defined for both cache and multicore labs - previously not so
  - Final lab due last week - previously one week earlier
- 437 is a required course, so if you give up now it will bite you later
- DON'T GIVE UP!!
  - MANY CE grads get jobs based on 437

ECE437, Fall 2016

(82)

- Ch 5b done!

ECE437, Fall 2016

(83)