Shared Memory

- SMP: shared memory multiprocessor
 - Hardware provides single physical address space for all processors
 - Synchronize shared variables using locks
 - Memory access time
 - UMA (uniform) vs. NUMA (nonuniform)







Figure 5.1 Basic structure of a centralized shared-memory multiprocessor based on a multicore chip.

Multiple processor-cache subsystems share the same physical memory, typically with one level of shared cache on the multicore, and one or more levels of private per-core cache. The key architectural property is the uniform access time to all of the memory from all of the processors. In a multichip design, an interconnection network links the processors and the memory, which may be one or more banks. In a single-chip multicore, the interconnection network is simply the memory bus.



Figure 5.8 A single-chip multicore with a distributed cache. In current designs, the distributed shared cache is usually L3, and levels L1 and L2 are private. There are typically multiple memory channels (2–8 in today's designs). This design is NUCA, since the access time to L3 portions varies with faster access time for the directly attached core. Because it is NUCA, it is also NUMA.

Cache Coherence

Συνοχή Coherence

- All reads by any processor must return the most recently written value
- entralized Shared-Memory Architectures Writes to the same location by any two processors are seen in the same order by all processors
- Συνέπεια Consistency
 - When a written value will be returned by a read
 - If a processor writes location A followed by location B, any processor that sees the new value of B must also see the new value of A



Cache Coherence Problem

- Suppose two CPU cores share a physical address space
 - Write-through caches

Time step	Event	CPU A's cache	CPU B's cache	Memory
0				0
1	CPU A reads X	0		0
2	CPU B reads X	0	0	0
3	CPU A writes 1 to X	1	0	1



Coherence Defined

- Informally: Reads return most recently written value
- Formally:
 - P writes X; P reads X (no intervening writes)
 ⇒ read returns written value
 - P₁ writes X; P₂ reads X (sufficiently later)
 ⇒ read returns written value
 - c.f. CPU B reading X after step 3 in example
 - P₁ writes X, P₂ writes X
 - \Rightarrow all processors see writes in the same order
 - End up with the same final value for X



Cache Coherence Protocols

- Operations performed by caches in multiprocessors to ensure coherence
 - Migration of data to local caches
 - Reduces bandwidth for shared memory
 - Replication of read-shared data
 - Reduces contention for access
- Snoping protocols
 All activities that potentially affect other caches are broadcast onto the shared bus; all caches monitor ("snoop") that shared bus.
 - Each cache monitors bus reads/writes^{OK for few (4, 8, 16?)} sharers, but too much
 - traffic beyond ~8.
- Directory-based protocols
 - Caches and memory record sharing status of blocks in a directory ^{A central Directory} (may consist of interleaved banks) records which caches have copies of which blocks => only "bother" those caches that are affected



Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 95



Upon Writes: to Invalidate or to Update the Others? · Invalidate-based Protocols: When I write (modify) something that I have, and there is a danger others may copies of it, tell them to Invalidate their copies." (like telling them: "If you ever need it again, come back to me and ask me for its latert value") ho con • Advantage: from that point on, J Know that J have the only Copy, therefore J can freely "play with:t" as long as notody ask me for a copy again. When I write (modify) something that I have, and there is a danger others may have copies of it, broadcast the new value so as to (arely) update the values of all copies! · Advantage: if others will need the value again Soon, they will have it in their cacles (iterer) · Disadvantage: multiple apries will keep existing, hence used (statistics showed that disadvantage the need to continuously keep updating them is important in the general case (unless there is hint about some data by the algorithm (software)

Invalidating Snooping Protocols

- Cache gets exclusive access to a block when it is to be written
 - Broadcasts an invalidate message on the bus
 - Subsequent read in another cache misses

Owning cache supplies updated value

CPU activity	Bus activity	CPU A's cache	CPU B's cache	Memory
				0
CPU A reads X	Cache miss for X	0		0
CPU B reads X	Cache miss for X	0	0	0
CPU A writes 1 to X	Invalidate for X	1		0
CPU B read X	Cache miss for X	1	, 1	1



Write-Ba

Snoopy Coherence Protocols

- Locating an item when a read miss occurs
 - In write-back cache, the updated value must be sent to the requesting processor
- Cache lines marked as shared or exclusive/modified
 - Only writes to shared lines need an invalidate broadcast
 - After this, the line is marked as exclusive



Cache Line States: what do I know about it? "MDESI" I have modified M "Modified" A FI my copy, and I A for writing it Owned" back to memory + 3) It is guaranteed that the memory (and others) S "Shared E "Exclusive" have the same Value as I do. (potentially) Shared: Exclusive : I "Invalid" other caches It is guaranteed that may have copies my copy is the ONLY copy of this line currently existing in (not known for sure) any cache. nothing in this Line.

Simplification: MSI Protocol/No "E" info when Clean; M (Modified): Dirty and Exclusive (If Dirty unst be Exclusive) - with invalidate based protocol, when I write into my copy, (no "O" state) I have to invalidate all others, hence I am guaranteed to have the ONLY copy => Exclusive S (Shared): (potentially) Shared and guaranteed Clean -when first reading from memory -> S -for simplicity, I do not keep track whether or not others too may have april - if I had the line as MI, and another cache misses it, then: . I have to write it back to memory (no "O" state, for simplicity) • the other cache gets a copy automatically on the bus - I may have had the line as is, and all other caches many have existed their copies, so I may be the only one having a copy, but I do NOT know that for sure... - If I have the line a S, and I want to write into it, I must broad coust on Invalidate command, since I have no "E" info. I (Invalid)

from MSI to MESI protocol: - Add and maintain E (oxclusive and clean) info: . When I read miss, if no other cache responds (forster than memory) providing me the plata, and I have to wait for the memory to bringme the data => then I know I am "E". · Advantage versus MSI: If the line is "E" and I want to write into it, I do NOT need to broad cast any Invalidate: save bus traffic. MDESI protocol: - Add "O" (Owned) info: Line is showed, Memory is "Old", and the responsibility to write back is mine ! . When the line is "M" (Modified: dicty and exclusive), and another Conche read-misses on it, I supply the data to the other Cache (faster than memory), but I do not spend the time to write-back to memory (now): save time (for now). . The other caches that have copies, have them in S state, hence they are allowed to exict their copies without writing back: I have the (only) "O" wpy, hence the responsibility to write back is mine!



FIGURE e5.12.11 Cache coherence state diagram with the state transitions induced by the local processor shown in black and by the bus activities shown in gray. As in Figure e5.12.10, the activities on a transition are shown in bold.

Advanced Processors

Instruction-Level Parallelism (ILP)

Multiple Issue

Fetch multiple (e.g. 2, 4) instructions in parallel, and then consider how many and which

- Static multiple issue
 - Of them to execute in parallel
 Compiler groups instructions to be issued together
 - Packages them into "issue slots"
 - Compiler detects and avoids hazards available, fills-in noop's
- Dynamic multiple issue
 - CPU examines instruction stream and chooses instructions to issue each cycle
 - Compiler can help by reordering instructions
 - CPU resolves hazards using advanced techniques at runtime



Static Multiple Issue

- Compiler groups instructions into "issue packets"
 - Group of instructions that can be issued on a single cycle
 - Determined by pipeline resources required
- Think of an issue packet as a very long instruction
 - Specifies multiple concurrent operations
 - \Rightarrow Very Long Instruction Word (VLIW)



Scheduling Static Multiple Issue

- Compiler must remove some/all hazards
 - Reorder instructions into issue packets
 - No dependencies withⁱⁿ a packet
 - Possibly some dependencies between packets
 - Varies between ISAs; compiler must know!
 - Pad with nop if necessary





2 extra clock cycles lost

What if the program is?: RAW dependence?



- Does the compiler know for sure if i!=j
 - (OK to reorder sd–ld) or i==j (fwd in reg.)?
- If unknown to compiler, static sch. impossible
 - => dynamic scheduling at runtime (ooo pipe)

Dynamic Multiple Issue

- "Superscalar" processors
 CPU decides whether to issue 0, 1, 2, ...
 each cycle
 - Avoiding structural and data hazards
- Avoids the need for compiler scheduling
 - Though it may still help
 - Code semantics ensured by the CPU
 Allows executables to run on newer processors, with same ISA but different pipeline, without needing to be recompiled



Dynamic Pipeline Scheduling

- Allow the CPU to execute instructions out of order to avoid stalls
 - But commit result to registers in order
- Example
 - ld x31,20(x21) add x1,x31,x2 sub x23,x23,x3 andi x5,x23,20
- Out-of-Order (ooo) Execution
- In-Order Commit
 - (so as to flush results of misspeculated instructions, and also allow precise exceptions)
- Can start sub while add is waiting for Id



Why Do Dynamic Scheduling?

- Why not just let the compiler schedule code?
- Not all stalls are predicable
 - e.g., cache misses
- Can't always schedule around branches
 - Branch outcome is dynamically determined
- Different implementations of an ISA have different latencies and hazards



Does Multiple Issue Work?

The BIG Picture

- Yes, but not as much as we'd like
- Programs have real dependencies that limit ILP
- Some dependencies are hard to eliminate
 - e.g., pointer aliasing
- Some parallelism is hard to expose
 - Limited window size during instruction issue
- Memory delays and limited bandwidth
 - Hard to keep pipelines full
- Speculation can help if done well



Parallelism



Strong vs Weak Scaling

- Strong scaling: problem size fixed
 - As in example
- Weak scaling: problem size proportional to number of processors
 - 10 processors, 10 × 10 matrix
 - Time = $20 \times t_{add}$
 - 100 processors, 32 × 32 matrix
 - Time = 10 × t_{add} + 1000/100 × t_{add} = 20 × t_{add}
 - Constant performance in this example



Instruction and Data Streams

An alternate classification Streaming SIMD Extension (SSE)

		Data St	reams
		Single	Multiple
Instruction Streams	Single	SISD: Intel Pentium 4	SIMD : SSE Vector instructions of x86
	Multiple	MISD: No examples today	MIMD: Intel Xeon e5345

- SPMD: Single Program Multiple Data
 - A parallel program on a MIMD computer
 - Conditional code for different processors



often

SIMD

Operate elementwise on vectors of data

- E.g., MMX and SSE instructions in x86
 - Multiple data elements in 128-bit wide registers
- All processors execute the same instruction at the same time
 - Each with different data address, etc.
- Simplifies synchronization
- Reduced instruction control hardware
- Works best for highly data-parallel applications



Vector Processors

Data-Level Parallelism Identical & Independent operations on all elements of a vector (array) - one vector instr. replaces a loop

- Highly pipelined function units
- Stream data from/to vector registers to units
 - Data collected from memory into registers
 - Results stored from registers to memory
- Example: Vector extension to RISC-V
 - v0 to v31: 32 × 64-element registers, (64-bit elements)

Vector Length Register assists in counting the number of remaining elements to process
 Vector instructions
 Fld.v, fsd.v: load/store vector
 Sequential memory addresses, or strided (e.g. for 2D/3D arrays), or scatter-gather (via array of pointers
 fadd.d.v: add vectors of double
 fadd.d.vs: add scalar to each element of vector of double
 Significantly reduces instruction-fetch bandwidth



Vector vs. Scalar

- Vector architectures and compilers
 - Simplify data-parallel programming
 - Explicit statement of absence of loop-carried dependences
 - Reduced checking in hardware
 - Regular access patterns benefit from interleaved and burst memory
 - Avoid control hazards by avoiding loops
- More general than ad-hoc media extensions (such as MMX, SSE)

Better match with compiler technology



Vector vs. Multimedia Extensions

- Vector instructions have a variable vector width, multimedia extensions have a fixed width
- Vector instructions support strided access, multimedia extensions do not
- Vector units can be combination of pipelined and arrayed functional units:



A[8]

A[4]

B[8]

B[4]

C[0]

A[9]

A[5]

B[9]

B[5]

C[1]

A[6]

Element aroup

C[2]

B[6]

A[7]

C[3]

B[7]



Chapter 6 — Parallel Processors from Client to Cloud — 15



Multithreading

One "thread of control" = one (traditional) sequential program. Multiple threads = parallel program.

and the Caches

- mimic multiple Performing multiple threads of execution in but Share the Functional Units
- cores, thus: Replicate registers, PC, etc.
 - Fast switching between threads
 - Fine-grain multithreading
 - Switch threads after each cycle
 - Interleave instruction execution
 - If one thread stalls, others are executed
 - Coarse-grain multithreading
 - Only switch on long stall (e.g., L2-cache miss)
 - Simplifies hardware, but doesn't hide short stalls (eg, data hazards)



Simultaneous Multithreading

- In multiple-issue dynamically scheduled processor
 - Schedule instructions from multiple threads
 - Instructions from independent threads execute when function units are available
 - Within threads, dependencies handled by scheduling and register renaming
- Example: Intel Pentium-4 HT
 - Two threads: duplicated registers, shared function units and caches



Multithreading Example





Future of Multithreading

- Will it survive? In what form?
- Power considerations ⇒ simplified microarchitectures
 - Simpler forms of multithreading
- Tolerating cache-miss latency
 - Thread switch may be most effective
- Multiple simple cores might share resources more effectively

Two different threads may have two different working sets of data/instructions; is it better to place them in a single cache, or in two different caches as two separate cores would do?



GPU Architectures

Graphics Processing Units

- Processing is highly data-parallel
 - GPUs are highly multithreaded
 - Use thread switching to hide memory latency
 - Less reliance on multi-level caches
 - Graphics memory is wide and high-bandwidth
- Trend toward general purpose GPUs
 - Heterogeneous CPU/GPU systems
 - CPU for sequential code, GPU for parallel code
- Programming languages/APIs
 - DirectX, OpenGL
 - C for Graphics (Cg), High Level Shader Language (HLSL)
 - Compute Unified Device Architecture (CUDA)



Example: NVIDIA Fermi

Multiple SIMD processors, each as shown:





Example: NVIDIA Fermi

- SIMD Processor: 16 SIMD lanes
- SIMD instruction
 - Operates on 32 element wide threads
 - Dynamically scheduled on 16-wide processor over 2 cycles
- 32K x 32-bit registers spread across lanes
 - 64 registers per thread context



GPU Memory Structures





Message Passing

- Each processor has private physical address space
- Hardware sends/receives messages between processors





Loosely Coupled Clusters

- Network of independent computers
 - Each has private memory and OS
 - Connected using I/O system
 - E.g., Ethernet/switch, Internet
- Suitable for applications with independent tasks
 - Web servers, databases, simulations, …
- High availability, scalable, affordable
- Problems
 - Administration cost (prefer virtual machines)
 - Low interconnect bandwidth
 - c.f. processor/memory bandwidth on an SMP

