

COMPUTER ORGANIZATION AND DESIGN

The Hardware/Software Interface



Large and Fast: Exploiting Memory Hierarchy

Principle of Locality

- §5.1 Introduction
- Programs access a small proportion of their address space at any time
 - Temporal locality Χρονική Τοπικότητα
 - Items accessed recently are likely to be accessed again soon
 - e.g., instructions in a loop, induction variables
 - Spatial locality

- Χωρική Τοπικότητα
- Items near those accessed recently are likely to be accessed soon
- E.g., sequential instruction access, array data

Hash Function: LS Block Address bits



Increased Line (Block) Size, to exploit Spatial Locality



'Vertical' Layout of the Words in a Line(Block)



Write-Through Ταυτόχρονη Εγγραφή

- On data-write hit, could just update the block in cache
 - But then cache and memory would be inconsistent
- Write through: also update memory
- But makes writes take longer
 - e.g., if base CPI = 1, 10% of instructions are stores, write to memory takes 100 cycles
 Write-Combined
 - Effective CPI = 1 + 0.1×100 = 11
- Solution: write buffer
 - Holds data waiting to be written to memory
 - CPU continues immediately
 - Only stalls on write if write buffer is already full



Write-Combining: sequential accesses to DRAM take shorter for subsequent words beyond the first one

Write-Back Ετερόχρονη Εγγραφή

- Alternative: On data-write hit, just update the block in cache
 - Keep track of whether each block is dirty
- When a dirty block is replaced
 - Write it back to memory
 - Can use a write buffer to allow replacing block to be read first

Main Memory is inconsistent with Cache

We will revisit this when talking about I/O, then about multicores...



Associative Caches

Μερικώς Προσεταιριστικές Κρυφές Μνήμες

Fully associative

- Allow a given block to go in any cache entry
- Requires all entries to be searched at once
- Comparator per entry (expensive)
- n-way set associative
 - Each set contains n entries Index portion of address
 - Block number determines which set

 (address)
 (Block number) modulo (#Sets in cache)
 - Search all entries in a given set at once
 - n comparators (less expensive)



Set Associative Cache Organization





Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 46

Replacement Policy

- Πώς να προβλέψουμε το μέλλον;; Συχνά, το πρόσφατο παρελθόν αποτελεί καλή ένδειξη γιά το προσεχές μέλλον!...
- Direct mapped: no choice
- Set associative
 - Prefer non-valid entry, if there is one
 - Otherwise, choose among entries in the set
- Least-recently used (LRU)
 - Choose the one unused for the longest time
 - Simple for 2-way, manageable for 4-way, too hard beyond that
- Random
 - Gives approximately the same performance as LRU for high associativity



Virtual Memory

- Use main memory as a "cache" for secondary (disk) storage
 - Managed jointly by CPU hardware and the operating system (OS) Virtualization
 - Programs share main memory & Protection
- Also solve the Each gets a private virtual address space Fragmentation holding its frequently used code and data available mem.
 - Protected from other programs
- process is fragmented

problem:

for new

- CPU and OS translate virtual addresses to physical addresses
 - VM "block" is called a page
 - VM translation "miss" is called a page fault





Translation Using a Page Table



Physical address





Replacement and Writes

- To reduce page fault rate, prefer leastrecently used (LRU) replacement
 - Reference bit (aka use bit) in PTE set to 1 on access to page
 - Periodically cleared to 0 by OS
 - A page with reference bit = 0 has not been used recently
- Disk writes take millions of cycles
 - Block at once, not individual locations
 - Write through is impractical
 - Use write-back
 - Dirty bit in PTE set when page is written



Fast Translation Using a TLB





TLB and Cache Interaction



- If cache tag uses physical address
 - Need to translate before cache lookup
- Alternative: use virtual address tag
 - Complications due to aliasing
 - Different virtual addresses for shared physical address

Often we want: physical addr. cache, and TLB access in parallel with tag read from cache. This requires cache index to be fully contained in page offset bits, which means:

Cache Way Size ≤ Page Size

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 80

Memory Protection

- Different tasks can share parts of their virtual address spaces
 - But need to protect against errant access
- **Requires OS assistance Supervisor** mode can

only be

(hardwired)

exception

handler

address,

ecall or

exception

- Hardware support for OS protection entered at
 - Privileged supervisor mode (aka kernel mode)
 - Privileged instructions only executable in supervisor mode
- Page tables and other state information only only through accessible in supervisor mode

System call exception (e.g., ecall in RISC-V)

Handling Exceptions

- Save PC of offending (or interrupted) instruction
 - In RISC-V: Supervisor Exception Program Counter (SEPC)
- Save indication of the problem
 - In RISC-V: Supervisor Exception Cause Register (SCAUSE)
 - 64 bits, but most bits unused
 - Exception code field: 2 for undefined opcode, 12 for hardware malfunction, ...
- Jump to handler
 - Assume at 0000 0000 1C09 0000_{hex}

Many RISC-V computers store the exception entry address in a special register named Supervisor Trap Vector (STVEC), which the OS can load with a value of its choosing.



Chapter 4 — The Processor — 87



e.g.: 100 ns per Bounk Access, Request every 25 ns "Split Transactions" on request/reply "bus" ... pipelining ...mulliple transactions interleaved in time



*

° Certain addresses are not regular memory

Instead, they correspond to registers in I/O devices



Processor Checks Status before Acting

° Path to device generally has 2 registers:

- 1 register says it's OK to read/write (I/O ready), often called <u>Control Register</u>
- 1 register that contains data, often called <u>Data Register</u>
- ^o Processor reads from Control Register "Polling" in loop, waiting for device to set Ready bit in Control reg to say its OK (0 P 1) "Busy wait" if done
- Processor then loads from (input) or writes to (output) data register

*

- "Busy wait" if done continuously; else, poll multiple devices on every interrupt from the real-time clock (usu. 50-120 Hz)



 traditional ("non-coherent"...) caching does <u>NOT</u> work when other devices (I/0, other proc. cores) access memory independently
 note: write-through is a "half-solution": works for output; but not for input...

An I/O interrupt is like an overflow exceptions except:

- An I/O interrupt is "asynchronous"
- More information needs to be conveyed
- An I/O interrupt is asynchronous with respect to instruction execution:
 - I/O interrupt is not associated with any instruction, but it can happen in the middle of any given instruction
 - I/O interrupt does not prevent any instruction from completion

Direct Memory Access (DMA)





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.



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)

Cache Line States: what do I know about it? "MDESI" I have modified M "Modified" A FI my copy, and I A fin responsible A for whiting 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!

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



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)

