

COMPUTER ORGANIZATION AND DESIGN

The Hardware/Software Interface



Large and Fast: Exploiting Memory Hierarchy

Principle of Locality

- Programs access a small proportion of their address space at any time
- Temporal locality
 - Items accessed recently are likely to be accessed again soon
 - e.g., instructions in a loop, induction variables
- Spatial locality
 - Items near those accessed recently are likely to be accessed soon
 - E.g., sequential instruction access, array data

Taking Advantage of Locality

- Memory hierarchy
- Store everything on disk
- Copy recently accessed (and nearby) items from disk to smaller DRAM memory
 - Main memory
- Copy more recently accessed (and nearby) items from DRAM to smaller SRAM memory
 - Cache memory attached to CPU



Memory Hierarchy Levels



- Block (aka line): unit of copying
 - May be multiple words
- If accessed data is present in upper level
 - Hit: access satisfied by upper level
 - Hit ratio: hits/accesses
- If accessed data is absent
 - Miss: block copied from lower level
 - Time taken: miss penalty
 - Miss ratio: misses/accesses
 = 1 hit ratio
 - Then accessed data supplied from upper level



Memory Technology

- Static RAM (SRAM)
 - 0.5ns 2.5ns, \$2000 \$5000 per GB
- Dynamic RAM (DRAM)
 - 50ns 70ns, \$20 \$75 per GB
- Magnetic disk
 - 5ms 20ms, \$0.20 \$2 per GB
- Ideal memory
 - Access time of SRAM
 - Capacity and cost/GB of disk



Cache Memory

Cache memory

- The level of the memory hierarchy closest to the CPU
- Given accesses X₁, ..., X_{n-1}, X_n





a. Before the reference to X_n

b. After the reference to X_n

- How do we know if the data is present?
- Where do we look?





Direct Mapped Cache

- Location determined by address
- Direct mapped: only one choice
 - (Block address) modulo (#Blocks in cache)



- #Blocks is a power of 2
- Use low-order address bits



Hash Function: LS Block Address bits



Tags and Valid Bits

- How do we know which particular block is stored in a cache location?
 - Store block address as well as the data
 - Actually, only need the high-order bits
 - Called the tag
- What if there is no data in a location?
 - Valid bit: 1 = present, 0 = not present
 - Initially 0



8-blocks, 1 word/block, direct mappedInitial state

Index	V	Tag	Data
000	N		
001	N		
010	N		
011	N		
100	N		
101	N		
110	N		
111	N		



Word addr	Binary addr	Hit/miss	Cache block
22	10 110	Miss	110

Index	V	Тад	Data
000	N		
001	N		
010	N		
011	N		
100	N		
101	N		
110	Υ	10	Mem[10110]
111	N		



Word addr	Binary addr	Hit/miss	Cache block
26	11 010	Miss	010

Index	V	Tag	Data
000	N		
001	N		
010	Υ	11	Mem[11010]
011	N		
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		



Word addr	Binary addr	Hit/miss	Cache block
22	10 110	Hit	110
26	11 010	Hit	010

Index	V	Тад	Data
000	N		
001	N		
010	Y	11	Mem[11010]
011	N		
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		



Word a	ddr	Binary ad	dr	Hit/miss	Cache block
16		10 000		Miss	000
3		00 011		Miss	011
16		10 000		Hit	000
Index	V	Tag	Dat	а	
000	Υ	10	Mem[10000]		
001	N				
010	Y	11	Mem[11010]		
011	Υ	00	Mem[00011]		
100	N				
101	N				
110	Y	10	Mem[10110]		
111	N				



Word addr	Binary addr	Hit/miss	Cache block
18	10 010	Miss	010

Index	V	Tag	Data
000	Y	10	Mem[10000]
001	N		
010	Υ	10	Mem[10010]
011	Y	00	Mem[00011]
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		



Address Subdivision





Average Access Time

- Hit time is also important for performance
- Average memory access time (AMAT)
 - AMAT = Hit time + Miss rate × Miss penalty
- Example
 - CPU with 1ns clock, hit time = 1 cycle, miss penalty = 20 cycles, I-cache miss rate = 5%
 - AMAT = 1 + 0.05 × 20 = 2ns
 - 2 cycles per instruction



Measuring Cache Performance

- Components of CPU time
 - Program execution cycles
 - Includes cache hit time
 - Memory stall cycles
 - Mainly from cache misses
- With simplifying assumptions:

Memory stall cycles

= Memory accesses Program × Miss rate × Miss penalty

 $= \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Misses}}{\text{Instruction}} \times \text{Miss penalty}$



Cache Performance Example

Given I-Cache I-cache miss rate = 2% Accesses D-cache miss rate = 4% per Instruction = 1 Miss penalty = 100 cycles Base CPI (ideal caché) = 2 Load & stores are 36% of instructions Data Cache Miss cycles per instruction Accesses I-cache: 0.02 × 100 = 2 per Instruction D-cache: 0.36 × 0.04 × 100 = 1.44 • Actual CPI = 2 + 2 + 1.44 = 5.44Ideal CPU is 5.44/2 = 2.72 times faster



Performance Summary

- When CPU performance increased
 - Miss penalty becomes more significant
- Decreasing base CPI
 - Greater proportion of time spent on memory stalls
- Increasing clock rate
 - Memory stalls account for more CPU cycles
- Can't neglect cache behavior when evaluating system performance



Increased Line (Block) Size, to exploit Spatial Locality



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



Block Size Considerations

- Larger blocks should reduce miss rate
 - Due to spatial locality
- But in a fixed-sized cache
 - Larger blocks \Rightarrow fewer of them
- They also reduce the number of Tags, hence speed up Tag look-up
- More competition \Rightarrow increased miss rate
- Larger blocks \Rightarrow pollution
- Larger miss penalty
 - Can override benefit of reduced miss rate
 - Early restart and critical-word-first can help





FIGURE 5.11 Miss rate versus block size. Note that the miss rate actually goes up if the block size is too large relative to the cache size. Each line represents a cache of different size. (This figure is independent of associativity, discussed soon.) Unfortunately, SPEC CPU2000 traces would take too long if block size were included, so these data are based on SPEC92.

Example: Intrinsity FastMATH

- Embedded MIPS processor
 - 12-stage pipeline
 - Instruction and data access on each cycle
- Split cache: separate I-cache and D-cache
 - Each 16KB: 256 blocks × 16 words/block
 - D-cache: write-through or write-back
- SPEC2000 miss rates
 - I-cache: 0.4%
 - D-cache: 11.4%
 - Weighted average: 3.2%



Example: Intrinsity FastMATH



Cache Misses

- On cache hit, CPU proceeds normally
- On cache miss (on I-cache miss, can also let the rest of the pipeline
 - Stall the CPU pipeline proceed to completion)
 - Fetch block from next level of hierarchy
 - Instruction cache miss
 - Restart instruction fetch
 - Data cache miss
 - Complete data access

Out-of-Order Pipelines do not stall the pipeline, but look for subsequent instructions that do not depend on the miss data; their D-cache must support one or more outstanding misses



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...



Write Allocation

- What should happen on a write miss?
- Alternatives for write-through
 - Allocate on miss: fetch the block
 - Write around: don't fetch the block
 - Since programs often write a whole block before reading it (e.g., initialization)
- For write-back
 - Usually fetch the block

Several modern processors allow software to control this policy on a per-page granularity (via page-table flag)



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)



Associative Cache Example





Spectrum of Associativity

For a cache with 8 entries



Two-way set associative

Four-way set associative



Eight-way set associative (fully associative)

Tag Data Tag Data



Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 42

Associativity Example

- Compare 4-block caches
 - Direct mapped, 2-way set associative, fully associative
 - Block access sequence: 0, 8, 0, 6, 8

Direct mapped

Block	Cache	Hit/miss	Cache content after access				
address	index		0	1	2	3	
0	0	miss	Mem[0]				
8	0	miss	Mem[8]				
0	0	miss	Mem[0]				
6	2	miss	Mem[0]		Mem[6]		
8	0	miss	Mem[8]		Mem[6]		



Associativity Example

Block access seq.: 0, 8, 0, 6, 8

2-way set associative

Block	Cache	Hit/miss	Cache content after access			
address	index		Set 0		Se	t 1
0	0	miss	Mem[0]			
8	0	miss	Mem[0]	Mem[8]		
0	0	hit	Mem[0]	Mem[8]		
6	0	miss	Mem[0]	Mem[6]		
8	0	miss	Mem[8]	Mem[6]		

Fully associative

Block address	Hit/miss	Cache content after access			
0	miss	Mem[0]			
8	miss	Mem[0]	Mem[8]		
0	hit	Mem[0]	Mem[8]		
6	miss	Mem[0]	Mem[8]	Mem[6]	
8	hit	Mem[0]	Mem[8]	Mem[6]	



How Much Associativity

- Increased associativity decreases miss rate
 - But with diminishing returns
- Simulation of a system with 64KB
 D-cache, 16-word blocks, SPEC2000
 - 1-way: 10.3%
 - 2-way: 8.6%
 - 4-way: 8.3%
 - 8-way: 8.1%



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



Multilevel Caches

- Primary cache attached to CPU
 - Small, but fast
- Level-2 cache services misses from primary cache
 - Larger, slower, but still faster than main memory
- Main memory services L-2 cache misses
- Some high-end systems include L-3 cache



Multilevel Cache Example

Given

- CPU base CPI = 1, clock rate = 4GHz
- Miss rate/instruction = 2%
- Main memory access time = 100ns
- With just primary cache
 - Miss penalty = 100ns/0.25ns = 400 cycles
 - Effective CPI = 1 + 0.02 × 400 = 9



Example (cont.)

- Now add L-2 cache
- Relative to ALL accesses to the entire cache hierarchy – not just the L2 cache – Penalty = 5ns/0.25ns = 20 cycles
 - Primary miss with L-2 miss 400
 Extra penalty = 500 cycles
 - CPI = 1 + 0.02 × 20 + 0.005 × 400 = 3.4
 - Performance ratio = 9/3.4 = 2.6



Multilevel Cache Considerations

- Primary cache
 - Focus on minimal hit time
- L-2 cache
 - Focus on low miss rate to avoid main memory access
 - Hit time has less overall impact
- Results

- & often fewer ways (smaller associativity)
- L-1 cache usually smaller than a single cache
- L-1 block size smaller than L-2 block size



Interactions with Advanced CPUs

- Out-of-order CPUs can execute instructions during cache miss
 - Pending store stays in load/store unit
 - Dependent instructions wait in reservation stations
 - Independent instructions continue
- Effect of miss depends on program data flow
 - Much harder to analyse
 - Use system simulation



Interactions with Software

- Misses depend on memory access patterns
 - Algorithm behavior
 - Compiler optimization for memory access



