



UC Berkeley Teaching Professor Dan Garcia

Great Ideas in Computer Architecture (a.k.a. Machine Structures)





cs61c.org



UC Berkeley Professor Bora Nikolić



Set-Associative Caches





N-Way Set Associative Cache (1/3)

- Memory address fields:
 - Tag: same as before
 - Offset: same as before
 - Index: points us to the correct "row" (called a set in this case)
- So what's the difference?
 - each set contains multiple blocks
 - once we've found correct set, must compare with all tags in that set to find our data



Caches III (3)

Size of \$ is # sets X N blocks/set X block sizeGarcia, Nikolić



Associative Cache Example







N-Way Set Associative Cache (2/3)

Basic Idea

- cache is direct-mapped w/respect to sets
- each set is fully associative with N blocks in it

Given memory address:

- Find correct set using Index value.
- Compare Tag with all Tag values in that set.
- If a match occurs, hit!, otherwise a miss.
- Finally, use the offset field as usual to find the desired data within the block.



et to sets I blocks in it

in that set a miss. al to find





N-Way Set Associative Cache (3/3)

- What's so great about this?
 - even a 2-way set assoc cache avoids a lot of conflict misses
 - hardware cost isn't that bad: only need N comparators
- In fact, for a cache with M blocks,
 - it's Direct-Mapped if it's 1-way set assoc
 - it's Fully Assoc if it's M-way set assoc
 - so these two are just special cases of the more general set associative design



Caches III (6)





4-Way Set Associative Cache Circuit









Block Replacement with Example





Block Replacement Policy

Direct-Mapped Cache

- index completely specifies position which position a block can go in on a miss
- N-Way Set Assoc
 - index specifies a set, but block can occupy any position within the set on a miss
- **Fully Associative**
 - block can be written into any position
- Question: if we have the choice, where should we write an incoming block?
 - If there's a valid bit off, write new block into first invalid.
 - If all are valid, pick a replacement policy
 - rule for which block gets "cached out" on a miss.



Caches III (9)





Block Replacement Policy

- LRU (Least Recently Used)
 - Idea: cache out block which has been accessed (read or write) least recently
 - Pro: temporal locality \rightarrow recent past use implies likely future use: in fact, this is a very effective policy
 - Con: with 2-way set assoc, easy to keep track (one LRU) bit); with 4-way or greater, requires complicated hardware and much time to keep track of this

HFO

Idea: ignores accesses, just tracks initial order

Random

If low temporal locality of workload, works ok



Caches III (10)





Block Replacement Example

Our same 2-way set associative cache with a four byte total capacity and one byte blocks. We perform the following byte accesses:

0, 2, 0, 1, 4, 0, 2, 3, 5, 4

How many hits and how many misses will there be for the LRU replacement policy?





Caches III (11)

loc 0 loc 1





Ży

Caches III (12)





Cache Simulator!





Caches III (13)





Average Memory Access Time (AMAT)







- How to choose between associativity, block size, replacement & write policy?
- Design against a performance model
 - Minimize: Average Memory Access Time
 = Hit Time
 - + Miss Penalty x Miss Rate
 - Influenced by technology & program behavior
- Create the illusion of a memory that is large, cheap, and fast - on average
- How can we improve miss penalty?



Caches III (15)

ivity, block odel

ehavio hat is ge ty?





Improving Miss Penalty

- When caches first became popular, Miss Penalty ~ 10 processor clock cycles
- Today 3 GHz Processor (1/3 ns per clock cycle) and 80 ns to go to DRAM ~200 processor clock cycles!



Solution: another cache between memory and the processor cache: <u>Second Level (L2) Cache</u>



Caches III (16)







Great Idea #3: Principle of Locality / Memory Hierarchy



Let's see Cache configuration on Dan's computer...



Caches III (17)

Extremely fast Extremely expensive Tiny capacity

Fast Priced reasonably Medium capacity





Analyzing Multi-level cache hierarchy



Avg Mem Access Time = L1 Hit Time + L1 Miss Pate * L1 Miss Penalty L1 Miss Penalty = L2 Hit Time + L2 Miss Pate * L2 Miss Penalty Avg Mem Access Time = L1 Hit Time + L1 Miss Rate * (L2 Hit Time + L2 Miss Rate * L2 Miss Penalty)



Caches III (18)







Assume

- Hit Time = 1 cycle
- Miss rate = 5%
- Miss penalty = 20 cycles

Calculate AMAT...

Avg mem access time
 = 1 + 0.05 × 20
 = 1 + 1 cycles
 = 2 cycles



Caches III (19)





Ways to reduce miss rate

Larger cache

- limited by cost and technology
- hit time of first level cache < cycle time (bigger caches are slower)
- More places in the cache to put each block of memory – associativity
 - fully-associative
 - any block any line
 - N-way set associated
 - N places for each block
 - direct map: N=1



Caches III (20)

(bigger each block







L1

- size: tens of KB
- hit time: complete in one clock cycle
- miss rates: 1-5%

L2:

- size: hundreds of KB
- hit time: few clock cycles
- miss rates: 10-20%
- L2 miss rate is fraction of L1 misses that also miss in L2 why so high?



Caches III (21)





Example: with L2 cache

Assume

- L1 Hit Time = 1 cycle
- L1 Miss rate = 5%
- L2 Hit Time = 5 cycles
- L2 Miss rate = 15% (% L1 misses that miss)
- L2 Miss Penalty = 200 cycles
- L1 miss penalty = 5 + 0.15 * 200 = 35

• Avg mem access time = $1 + 0.05 \times 35$

Caches III (22)

= 2.75 cycles





Example: without L2 cache

Assume

- L1 Hit Time = 1 cycle
- L1 Miss rate = 5%
- L1 Miss Penalty = 200 cycles

• Avg mem access time = $1 + 0.05 \times 200$ = 11 cycles

4x faster with L2 cache! (2.75 vs. 11)



Caches III (23)







Actual CPUs





An Actual CPU – Early PowerPC

Cache

- 32 KiB Instructions & 32 KiB Data L1 caches
- External L2 Cache interface with integrated controller and cache tags, supports up to 1 MiB external L2 cache
- Dual Memory Management Units (MMU) with Translation Lookaside Buffers (TLB)

Pipelining

- Superscalar (3 inst/cycle)
- 6 execution units (2 integer and 1 double precision IEEE floating point)





Caches III (25)





An Actual CPU – Pentium M





Caches III (26)



Streaming SIMD Extensions II compatible with Pentium® 4 Processor optimized software

<u>Dedicated Stack</u>
 <u>Management</u> –
 faster instruction
 at lower power
 levels

Enhanced Intel® SpeedStep® Technology - Multiple voltages & frequency

operating points

400 MHz Power Optimized System Bus

 faster system bus to enhance performance at lower power levels





An Actual CPU – Intel core i7





Caches III (27)





And in Conclusion...

- We've discussed memory caching in detail. Caching in general shows up over and over in computer systems
 - Filesystem cache, Web page cache, Game databases / tablebases, Software memoization, Others?
- Big idea: if something is expensive but we want to do it repeatedly, do it once and cache the result.
- Cache design choices:
 - Size of cache: speed v. capacity
 - Block size (i.e., cache aspect ratio)
 - Write Policy (Write through v. write back
 - Associativity choice of N (direct-mapped v. set v. fully associative)
 - Block replacement policy
 - 2nd level cache?
 - 3rd level cache?
- Use performance model to pick between choices, depending on programs, technology, budget,



Caches III (28)

