

Memory Hierarchy Basics

Introduction

- Six basic cache optimizations:
 - Larger block size
 - Reduces compulsory misses
 - Increases capacity and conflict misses, increases miss penalty
 - Larger total cache capacity to reduce miss rate
 - Increases hit time, increases power consumption
 - Higher associativity
 - Reduces conflict misses
 - Increases hit time, increases power consumption
 - Higher number of cache levels
 - Reduces overall memory access time
 - Giving priority to read misses over writes
 - Reduces miss penalty
 - Avoiding address translation in cache indexing
 - Reduces hit time
 - Victim cache



Ten Advanced Optimizations

Advanced Optimizations

- Small and simple first level caches
- Way Prediction
- Pipelined caches
- Non-blocking cache
- Multibanked cache
- Critical word first
- Merging write buffer
- Compiler optimization
- Hardware prefetching
- Compiler prefetching

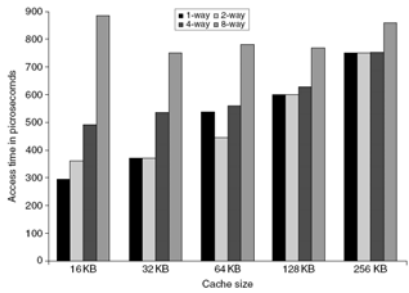


Small and Simple

- No mux in the critical path of a direct mapped cache.
- Bigger cache means more energy.
- CACTI – An idea for the project/paper review
- Many processors takes at least 2 clock cycles to access the cache, longer hit time may not be that critical
- The use of a virtual index cache, limits the cache size to page size \times associativity (recently a trend to increase associativity).

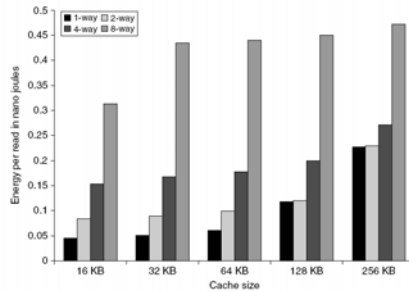


L1 Size and Associativity



Access time vs. size and associativity

L1 Size and Associativity



Energy per read vs. size and associativity

Way Prediction

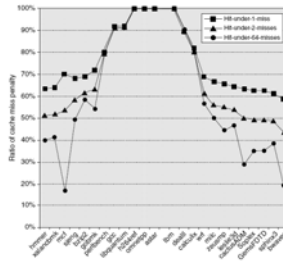
- To improve hit time, predict the way to pre-set mux
 - Mis-prediction gives longer hit time
 - Prediction accuracy
 - > 90% for two-way
 - > 80% for four-way
 - I-cache has better accuracy than D-cache
 - First used on MIPS R10000 in mid-90s
 - Used on ARM Cortex-A8
- Extend to predict block as well
 - "Way selection"
 - Increases mis-prediction penalty

Pipelining Cache

- Pipeline cache access to improve bandwidth
 - Examples:
 - Pentium: 1 cycle
 - Pentium Pro – Pentium III: 2 cycles
 - Pentium 4 – Core i7: 4 cycles
- Increases branch miss-prediction penalty (longer pipeline).
- Makes it easier to increase associativity

Nonblocking Caches

- For out-of-order execution (later on this point).
- Allow hits before previous misses complete
 - "Hit under miss"
 - "Hit under multiple miss"
- L2 must support this
- In general, processors can hide L1 miss penalty but not L2 miss penalty



Single core i7 using SPEC2006

Multibanked Caches

- Organize cache as independent banks to support simultaneous access
 - ARM Cortex-A8 supports 1-4 banks for L2
 - Intel i7 supports 4 banks for L1 and 8 banks for L2
- Interleave banks according to block address

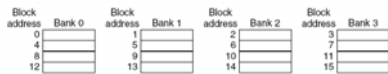


Figure 2.6 Four-way interleaved cache banks using block addressing. Assuming 64 bytes per blocks, each of these addresses would be multiplied by 64 to get byte addressing.

Critical Word First, Early Restart

- Critical word first
 - Request missed word from memory first
 - Send it to the processor as soon as it arrives
- Early restart
 - Request words in normal order
 - Send missed work to the processor as soon as it arrives
- Effectiveness of these strategies depends on block size and likelihood of another access to the portion of the block that has not yet been fetched

Merging Write Buffer

- When storing to a block that is already pending in the write buffer, update write buffer
- Reduces stalls due to full write buffer
- Do not apply to I/O addresses

Write address	V		V	V	V
100	1	Mem[100]	0	0	0
108	1	Mem[108]	0	0	0
116	1	Mem[116]	0	0	0
124	1	Mem[124]	0	0	0

No write buffering

Write address	V		V	V	V
100	1	Mem[100]	1	Mem[108]	1
	0		0	0	0
	0		0	0	0
	0		0	0	0

Write buffering

Compiler Optimizations

- Loop Interchange
 - Swap nested loops to access memory in sequential order (row major access)
- Blocking
 - Instead of accessing entire rows or columns, subdivide matrices into blocks
 - Requires more memory accesses but improves locality of accesses
