























### **Reducing Cache Miss Rates #1**

- 1. Allow more flexible block placement
- In a direct mapped cache a memory block maps to exactly one cache block
- At the other extreme, could allow a memory block to be mapped to any cache block – fully associative cache
- A compromise is to divide the cache into sets each of which consists of n "ways" (n-way set associative). A memory block maps to a unique set (specified by the index field) and can be placed in any way of that set (so there are n choices)

(block address) modulo (# sets in the cache)



















- 2. Use multiple levels of caches
- With advancing technology have more than enough room on the die for bigger L1 caches or for a second level of caches – normally a unified L2 cache – and in some cases even a unified L3 cache
- For our example, CPI<sub>ideal</sub> of 2, 100 cycle miss penalty (to main memory), 36% load/stores, a 2% (4%) L1 I\$ (D\$) miss rate, add a UL2\$ that has a 25 cycle miss penalty and a 0.5% miss rate

CPI<sub>stalls</sub> = 2 + .02×25 + .36×.04×25 + .005×100 + .36×.005×100 = 3.54 (as compared to 5.44 with no L2\$)



• L2\$ local miss rate >> than the global miss rate

# L1 typical L2 typical Total size (blocks) 255 to 2000 4000 to

I OTAI SIZE (DIOCKS)	250 to 2000	4000 to 250,000
Total size (KB)	16 to 64	500 to 8000
Block size (B)	32 to 64	32 to 128
Miss penalty (clocks)	10 to 25	100 to 1000
Miss rates (global for L2)	2% to 5%	0.1% to 2%

	Intel P4	AMD Opteron	
L1 organization	Split I\$ and D\$	Split I\$ and D\$	
L1 cache size	8KB for D\$, 96KB for trace cache (~I\$)	64KB for each of I\$ and D\$	
L1 block size	64 bytes	64 bytes	
L1 associativity	4-way set assoc.	2-way set assoc.	
L1 replacement	~ LRU	LRU	
L1 write policy	write-through	write-back	
L2 organization	Unified	Unified	
L2 cache size	512KB	1024KB (1MB)	
L2 block size	128 bytes	64 bytes	
L2 associativity	8-way set assoc.	16-way set assoc.	
L2 replacement	~LRU	~LRU	
L2 write policy	write-back	write-back	

# 4 Questions for the Memory Hierarchy

- Q1: Where can a block be placed in the upper level? (Block placement)
- Q2: How is a block found if it is in the upper level? (Block identification)
- Q3: Which block should be replaced on a miss? (Block replacement)
- Q4: What happens on a write? (Write strategy)

# Q1&Q2: Where can a block be placed/found?

		<b>B</b> 1 1 1
	# of sets	Blocks per set
Direct mapped	# of blocks in cache	1
Set associative	(# of blocks in cache)/ associativity	Associativity (typically 2 to 16)
Fully associative	1	# of blocks in cache

	Location method	# of comparisons
Direct mapped	Index	1
Set associative	Index the set; compare set's tags	Degree of associativity
Fully associative	Compare all blocks tags	# of blocks

# Q3: Which block should be replaced on a miss?

- Easy for direct mapped only one choice
- Set associative or fully associative
  - Random
  - LRU (Least Recently Used)
- For a 2-way set associative cache, random replacement has a miss rate about 1.1 times higher than LRU.
- LRU is too costly to implement for high levels of associativity (> 4-way) since tracking the usage information is costly

### Q4: What happens on a write?

- <u>Write-through</u> The information is written to both the block in the cache and to the block in the next lower level of the memory hierarchy
  - Write-through is always combined with a write buffer so write waits to lower level memory can be eliminated (as long as the write buffer doesn't fill)
- <u>Write-back</u> The information is written only to the block in the cache. The modified cache block is written to main memory only when it is replaced.
  - Need a dirty bit to keep track of whether the block is clean or dirty
- Pros and cons of each?
  - Write-through: read misses don't result in writes (so are simpler and cheaper)
  - Write-back: repeated writes require only one write to lower level

# Improving Cache Performance

- 1. Reduce the miss rate
  - bigger cache
  - more flexible placement (increase associativity)
  - larger blocks (16 to 64 bytes typical)
  - victim cache small buffer holding most recently discarded blocks

# 2. Reduce the miss penalty

- smaller blocks
- use a write buffer to hold dirty blocks being replaced so don't have to wait for the write to complete before reading
- check write buffer on read miss may get lucky
- for large blocks fetch critical word first
- use multiple cache levels L2 cache not tied to CPU clock rate
- faster backing store/improved memory bandwidth
- wider buses
   memory interleaving, page mode DRAMs

### Improving Cache Performance

- 3. Reduce the time to hit in the cache
  - smaller cache
  - direct mapped cache
  - smaller blocks
  - for writes
    - no write allocate no "hit" on cache, just write to write buffer
       write allocate to avoid two cycles (first check for hit, then write) pipeline writes via a delayed write buffer to cache







# Recap Q3: Which block should be replaced on a miss?

- Easy for Direct Mapped
- Set Associative or Fully Associative:
  - Random
  - LRU (Least Recently Used)

Associativity	2 way	2 way	4 way	4 way	8 way	8 way
Size	LRU	Random	LRU	Random	LRU	Random
16 KB	5.2%	5.7%	4.7%	5.3%	4.4%	5.0%
64 KB	1.9%	2.0%	1.5%	1.7%	1.4%	1.5%
256 KB	1.15%	1.17%	1.13%	1.13%	1.12%	1.12%

# Know This!

- Calculate runtime given cache statistics (miss rate, miss penalty, etc.)
- □ Calculate Average Memory Access Time (AMAT)
- Understand direct mapped, set-associative, fully associative
  - Comparison between them
  - Sources of cache misses
  - Architecture
  - Addressing into them (tag, index, byte)
  - Cache behavior with them

### Sample Problem: Impact on Performance

- Suppose a processor executes at
  - Clock Rate = 2 GHz (.5 ns per cycle)
  - Base CPI = 1.1 (assuming 1-cycle cache hits)
  - 50% arith/logic, 30% ld/st, 20% control
- Suppose that 10% of data memory operations (lw, sw) get 50 cycle miss penalty
- Suppose that 1% of instructions get same miss penalty
- What is CPI?
- □ What is AMAT?

### Sample Problem: Impact on Performance

- CPI = Base CPI + average stalls per instruction
- 1.1 (cycles/ins) +
- [ 0.30 (DataMops/ins) x 0.10 (miss/DataMop) x 50(cycle/miss)] + [ 1 (InstMop/ins)x 0.01 (miss/InstMop) x 50 (cycle/miss)] = (1.1 + 1.5 + .5) cycle/ins = 3.1
- □ 64.5% of the time the proc is stalled waiting for memory!
- AMAT=(1/1.3)x[1+0.01x50]+(0.3/1.3)x[1+0.1x50]=2.54

### Sample Problem: Cache Performance

### Processor:

- CPI = 2
- Icache miss rate = 2%, miss penalty = 100 cycles
- Dcache miss rate = 4%, miss penalty = 100 cycles
- Using SPECint2000 load/store percentage of 36%:
  - Speedup from this processor to one that never missed?
    What if CPI = 1?

  - What if we doubled clock rate of computer without changing memory speed (CPI = 2)?

# Calculating Cache Performance, CPI=2

I miss cycles = I \* 2% \* 100 = 2.00 \* I
D miss cycles = I \* 4% \* 36% \* 100 = 1.44 \* I
So memory stalls/instr = 2.00 + 1.44 = 3.44

□ <u>CPU time with stalls</u> = <u>I \* CPI<sub>stall</sub> \* clkcycle</u>
 □ CPU time, no stalls
 I \* CPI<sub>perf</sub> \* clkcycle

io i or iperi one

□ CPI<sub>stall</sub> = 2 + 3.44; CPI<sub>perf</sub> = 2 □ CPI<sub>stall</sub> / CPI<sub>perf</sub> = 5.44 / 2 = 2.72

# Calculating Cache Performance, CPI=1

□ I miss cycles = I \* 2% \* 100 = 2.00 \* I

D miss cycles = I \* 4% \* 36% \* 100 = 1.44 \* I

□ So memory stalls/instr = 2.00 + 1.44 = 3.44

□ <u>CPU time with stalls</u> = <u>I \* CPI<sub>stall</sub> \* clkcycle</u> □ CPU time, no stalls I \* CPI<sub>perf</sub> \* clkcycle

CPI<sub>stall</sub> = 1 + 3.44; CPI<sub>perf</sub> = 1
 CPI<sub>stall</sub> / CPI<sub>perf</sub> = 4.44 / 4 = 4.44

 Calculating Cache Performance, 2x

 I miss cycles = 1 \* 2% \* 200 = 4.00 \* 1

 D miss cycles = 1 \* 4% \* 36% \* 200 = 2.88 \* 1

 So memory stalls/instr = 4.00 + 2.88 = 6.88

 CPU time (slow clk)

 = CPU time (fast clk)

 I \* CPI<sub>slow</sub> \* clkcycle

 CPU time (fast clk)

 I \* CPI<sub>fast</sub> \* clkcycle/2

 CPI<sub>fast</sub> = 2 + 6.88; CPI<sub>slow</sub> = 5.44

 Speedup = 5.44 / (8.88 \* 0.5) = 1.23

 Ideal machine is 2x faster