

Who Cares about Memory Hierarchy?

• Processor Only Thus Far in Course

Dean Tullsen

Locality and cacheing

- Memory hierarchies exploit locality by *cacheing* (keeping close to the processor) data likely to be used again.
- ٠ This is done because we can build large, slow memories and small, fast memories, but we can't build large, fast memories.
- If it works, we get the illusion of SRAM access time with disk capacity

SRAM (static RAM) -- 5-20 ns access time, very expensive (onchip faster) DRAM (dynamic RAM) -- 60-100 ns, cheaper disk -- access time measured in milliseconds, very cheap

Dean Tullsen

Cache Fundamentals



is found in the cache. • *cache miss* -- an access which isn't

- *hit time* -- time to access the higher cache *miss penalty* -- time to move data from
- *miss penalty* -- time to move data from lower level to upper, then to cpu
- *hit ratio* -- percentage of time the data is found in the higher cache
- *miss ratio* -- (1 hit ratio)



Cache Fundamentals, cont.

• *cache block size* or *cache line size--* the amount of data that gets transferred on a cache miss.



instructions.

• *instruction cache* -- cache that only holds

- *data cache* -- cache that only caches data.
- *unified cache* -- cache that holds both.



A typical memory hierarchy





Cache Organization: Where can a block be placed in the cache?



Cache Access: How Is a Block Found In the Cache?

- Tag on each block
 - No need to check index or block offset
- · Increasing associativity shrinks index, expands tag

Address



•

Dean Tullsen

Cache Organization – Overview



Cache Access

- 16 KB, 4-way set-associative cache, 32-bit address, byteaddressable memory, 32-byte cache blocks/lines
- how many tag bits?
- Where would you find the word at addres 0x200356A4?



Which Block Should be Replaced on a Miss?

- Direct Mapped is Easy
- Set associative or fully associative:
 - longest till next use (ideal, impossible)
 - least recently used (best practical approximation)
 - pseudo-LRU (e.g., NMRU, NRU)
 - random (easy)
 - how many bits for LRU?

Associativity:	2-way		4-way		8-way	
Size	LRU	Random	LRU	Random	LRU	Random
16 KB	5.18%	5.69%	4.67%	5.29%	4.39%	4.96%
64 KB	1.88%	2.01%	1.54%	1.66%	1.39%	1.53%
256 KB	1.15%	1.17%	1.13%	1.13%	1.12%	1.12%
64 KB 256 KB	1.88% 1.15%	2.01% 1.17%	1.54% 1.13%	1.66% 1.13%	1.39% 1.12%	1.53% 1.12%

Dean Tullsen

A set-associative cache

4	00000100						
8	00001000	00000100	taa	data	taa	data	
12	00001100	I F	tag	uata	lag	uata	
4	00000100	L,					
8	00001000	-					
20	00010100	ŀ					
4	00000100	L			I		
8	00001000						
20	00010100	4 entries, each block holds one word, each word in memory maps to one of a set of <i>n</i> cache lines					
24	00011000						
36	00100100						
	00000100						
4	00000100						

- A cache that can put a line of data in exactly *n* places is called *n*-way set-associative.
- The cache lines that share the same index are a cache *set*.

Cache Access Cache Access • 32-bit address, and 128 KB cache with 64-byte blocks, 4-• 48-bit address, and 1 MB cache with 64-byte blocks, 8way set associative. way set associative. 4 blocks/set 8 blocks/set 32 48 ?? sets sets 2 CSE 240A Dean Tullsen CSE 240A Dean Tullsen **Cache Access** What Happens on a Write? • 64-bit address, and 32 KB cache with 32-byte blocks, • Write through: The information is written to both the block in direct-mapped. the cache and to the block in the lower-level memory. 4 blocks/set 64 • Write back: 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. - is block clean or dirty? sets • Pros and Cons of each: 2 - WT: read misses cannot result in writes (because of replacements) - WB: no writes of repeated writes • WT always combined with *write buffers* so that don't wait for lower level memory CSE 240A Dean Tullsen CSE 240A Dean Tullsen

What happens on a write miss?

- *write-allocate* -- make room for the cache line in the cache, fetch rest of line from memory.
- *no-write-allocate* (write-around) -- write to lower levels of memory hierarchy, ignoring this cache.
- Tradeoffs?

CSE 240A

- Which makes most sense for write-back?
- Which makes most sense for write-through?

21264 L1 Data Cache

• 64 KB, 64-byte blocks, 2-way set associative, ? blocks, ? sets



Cache Organization: Separate Instruction and Data Caches?

Size	Instruction Cache	Data Cache	Unified Cache
1 KB	3.06%	24.61%	13.34%
2 KB	2.26%	20.57%	9.78%
4 KB	1.78%	15.94%	7.24%
8 KB	1.10%	10.19%	4.57%
16 KB	0.64%	6.47%	2.87%
32 KB	0.39%	4.82%	1.99%
64 KB	0.15% —	3.77%	1.35%
128 KB	0.02%	2.88%	→0.95%

if 75% of accesses are instructions? Other reasons to separate?

Cache Performance

- CPU time = (CPU execution clock cycles + Memory stall clock cycles) x clock cycle time
- Memory stall clock cycles = (Reads x Read miss rate x Read miss penalty + Writes x Write miss rate x Write miss penalty)
- Memory stall clock cycles = Memory accesses x Miss rate • x Miss penalty

Dean Tullsen

Dean Tullsen

CSE 240A

CSE 240A

Cache Performance

CPUtime = IC x (CPI_{execution} + Memory stalls per instruction) Average memory-access time = Hit time + Miss rate x Miss penalty (ns or x Clock cycle time clocks) How are we going to improve cache performance?? CPUtime = IC x ($CPI_{execution}$ + Mem accesses per instruction x 1. Miss rate x Miss penalty) x Clock cycle time (includes hit time as part of CPI) 2. (Alternate view of memory performance) 3. Average memory-access time = Hit time + Miss rate x Miss penalty (ns or clocks) CSE 240A Dean Tullsen CSE 240A Dean Tullsen

Improving Cache Performance

Caches, pt I: Key Points

- CPU-Memory gap is a major performance obstacle
- Caches take advantage of program behavior: locality
- Designer has lots of choices -> cache size, block size, associativity, replacement policy, write policy, ...
- Time of program still only reliable performance measure