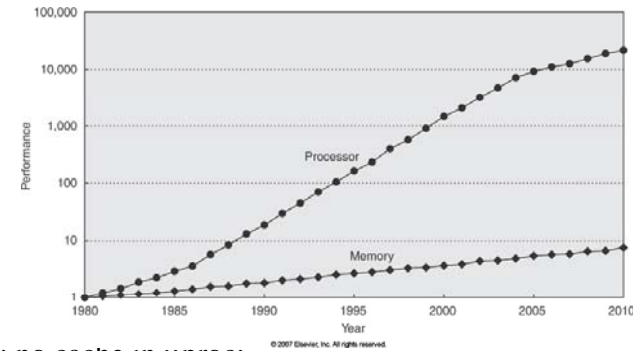


Memory Subsystem Design

OR
Nothing Beats Cold, Hard Cache

Who Cares about Memory Hierarchy?

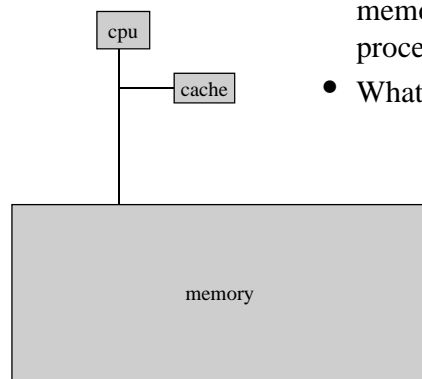
- Processor Only Thus Far in Course



1980: no cache in μ proc;

1995 2-level cache, 60% trans. on Alpha 21164 μ proc

Memory Cache



- Can put **small, fast** memory close to processor.
- What do we put there?

Memory Locality

- Memory hierarchies take advantage of **memory locality**.
- **Memory locality** is the principle that future memory accesses are *near* past accesses.
- Memory hierarchies take advantage of two types of locality
 - **Temporal locality** -- near in time => we will often access the same data again very soon
 - **Spatial locality** -- near in space/distance => our next access is often very close to our last access (or recent accesses).

1,2,3,1,2,3,8,8,47,9,10,8,8...

Locality and cacheing

- Memory hierarchies exploit locality by *cacheing* (keeping close to the processor) data likely to be used again.
- 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)

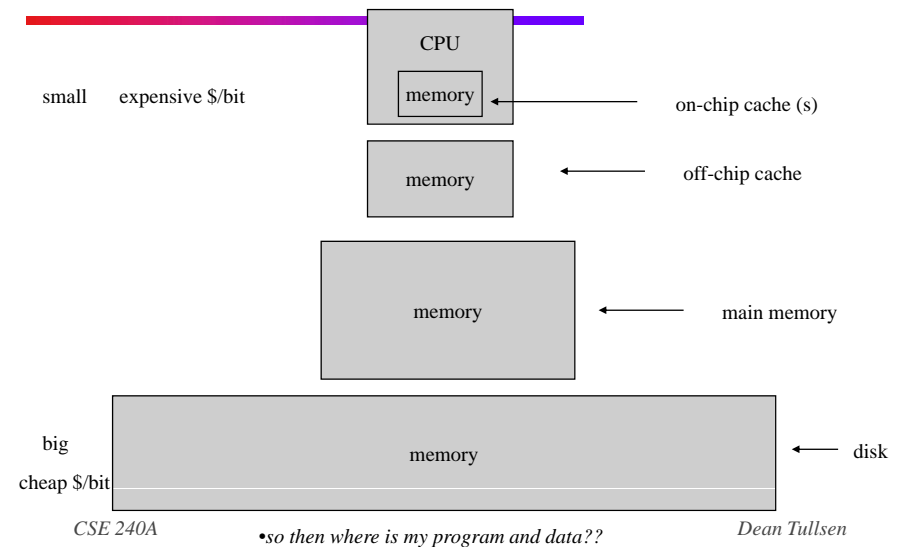
DRAM (dynamic RAM) -- 60-100 ns, cheaper

disk -- access time measured in milliseconds, very cheap

CSE 240A

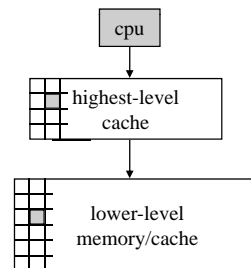
Dean Tullsen

A typical memory hierarchy



Cache Fundamentals

- *cache hit* -- an access where the data 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 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)

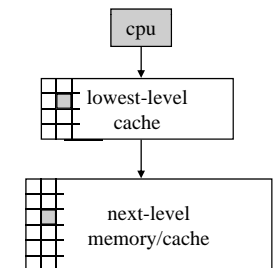


CSE 240A

Dean Tullsen

Cache Fundamentals, cont.

- *cache block size* or *cache line size*-- the amount of data that gets transferred on a cache miss.
- *instruction cache* -- cache that only holds instructions.
- *data cache* -- cache that only caches data.
- *unified cache* -- cache that holds both.

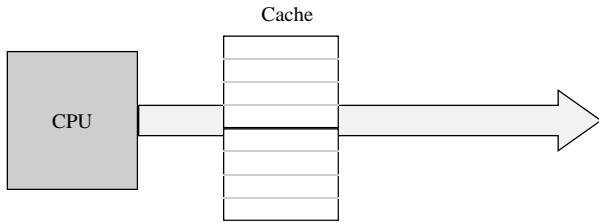


CSE 240A

Dean Tullsen

Accessing a simple cache

- blocksize = 4 words (16 bytes), cache size = 2 blocks (32 bytes), associativity = full



Memory	
0	63
4	37
8	4
12	149
16	12
20	82
24	2
28	3
32	4
36	5
40	17
44	3245
48	63
52	37
56	4
60	149
64	12
68	82

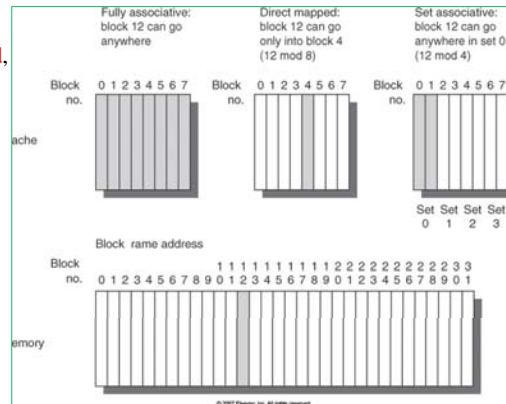
CPU reads addresses 8, 0, 4, 0, 4, 8, 12, 16, 20, writes 12, reads 8, 20, 28, 56, 20, 60, 12

Cache Characteristics

- Cache Organization (size, associativity, block size)
- Cache Access
- Cache Replacement
- Write Policy

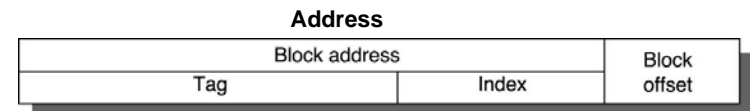
Cache Organization: Where can a block be placed in the cache?

- Block 12 placed in 8-block cache:
 - Fully associative, direct mapped, n-way set associative
 - index - pointer to the set in the cache where a memory location might be cached
 - (associativity = degree of freedom in placing a particular block of memory)
 - (set = a collection of cache blocks with the same cache index)

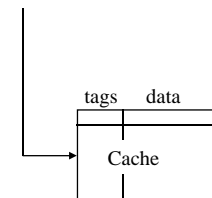


Cache Access: How Is a Block Found In the Cache?

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

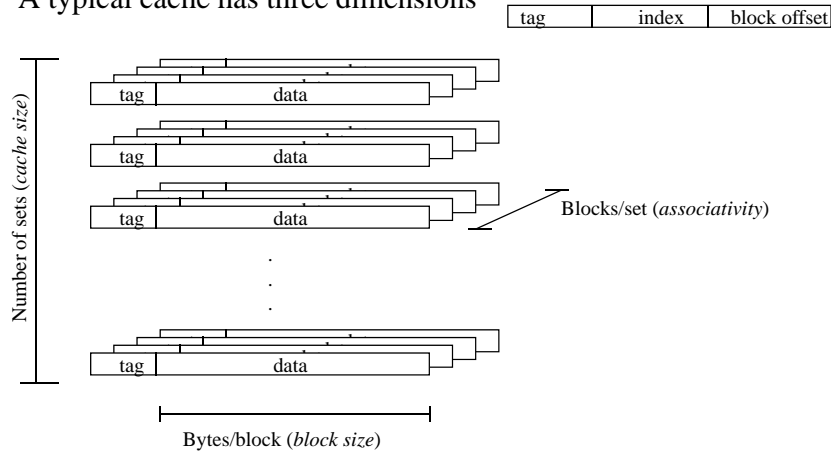


FA: No index, large tags
DM: Large index, smaller tags



Cache Organization – Overview

- A typical cache has three dimensions



CSE 240A

Dean Tullsen

Cache Access

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

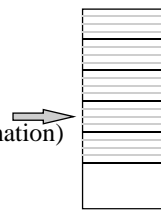


CSE 240A

Dean Tullsen

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%

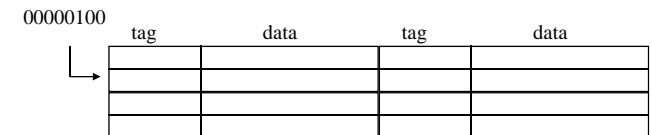
CSE 240A

Dean Tullsen

A set-associative cache

address string:

- 4 00000100
- 8 00001000
- 12 00001100
- 4 00000100
- 8 00001000
- 20 00010100
- 4 00000100
- 8 00001000
- 20 00010100
- 24 00011000
- 36 00100100
- 4 00000100
- 20 00010100



4 entries, each block holds one word, each word in memory maps to one of a set of n cache lines

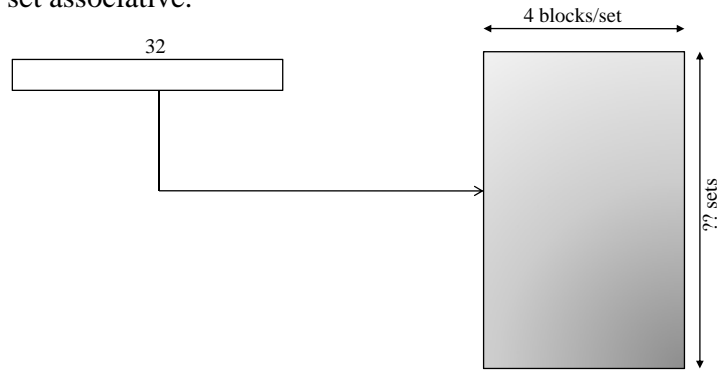
- A cache that can put a line of data in exactly n places is called n -way set-associative.
- The cache lines that share the same index are a cache set.

CSE 240A

Dean Tullsen

Cache Access

- 32-bit address, and 128 KB cache with 64-byte blocks, 4-way set associative.

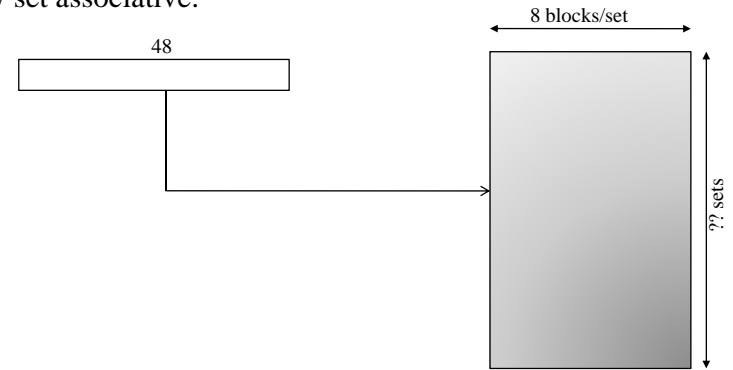


CSE 240A

Dean Tullsen

Cache Access

- 48-bit address, and 1 MB cache with 64-byte blocks, 8-way set associative.

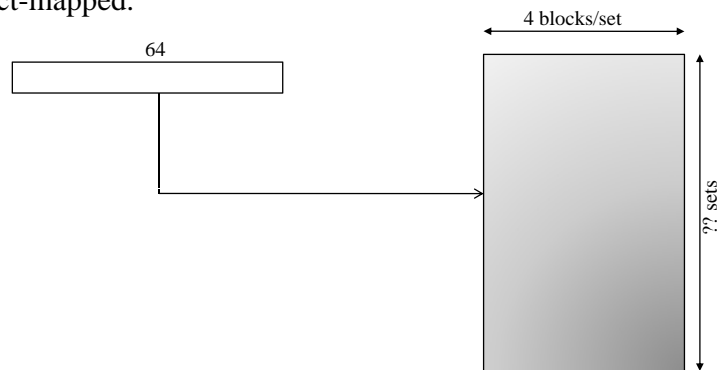


CSE 240A

Dean Tullsen

Cache Access

- 64-bit address, and 32 KB cache with 32-byte blocks, direct-mapped.



CSE 240A

Dean Tullsen

What Happens on a Write?

- *Write through*: The information is written to both the block in the cache and to the block in the lower-level memory.
- *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?
- Pros and Cons of each:
 - 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

What happens on a write miss?

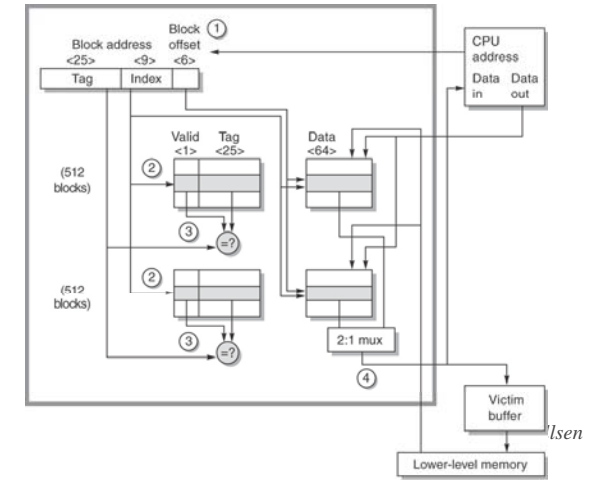
- *write-allocate* -- make room for the cache line in the cache, fetch rest of line from memory.
- *no-write-allocate* (write-around) -- write to lower levels of memory hierarchy, ignoring this cache.
- Tradeoffs?
- Which makes most sense for write-back?
- Which makes most sense for write-through?

CSE 240A

Dean Tullsen

21264 L1 Data Cache

- 64 KB, 64-byte blocks, 2-way set associative, 256 blocks, 2 sets
- write-back



CSE 240A

© 2007 Altera, Inc. All rights reserved.

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?

CSE 240A

Dean Tullsen

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

CSE 240A

Dean Tullsen

Cache Performance

$$\text{CPUtime} = \text{IC} \times (\text{CPI}_{\text{execution}} + \text{Memory stalls per instruction}) \times \text{Clock cycle time}$$

$$\text{CPUtime} = \text{IC} \times (\text{CPI}_{\text{execution}} + \text{Mem accesses per instruction} \times \text{Miss rate} \times \text{Miss penalty}) \times \text{Clock cycle time}$$

(includes hit time as part of CPI)

(Alternate view of memory performance)

$$\text{Average memory-access time} = \text{Hit time} + \text{Miss rate} \times \text{Miss penalty (ns or clocks)}$$

Improving Cache Performance

$$\text{Average memory-access time} = \text{Hit time} + \text{Miss rate} \times \text{Miss penalty (ns or clocks)}$$

How are we going to improve cache performance??

- 1.
- 2.
- 3.

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