


Cache Principles



What did you get for 7?

Hit.

8?

Miss.

14?

Hit.

19?

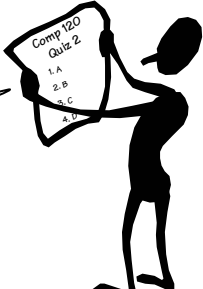
Miss!

FIFO?

LRU.

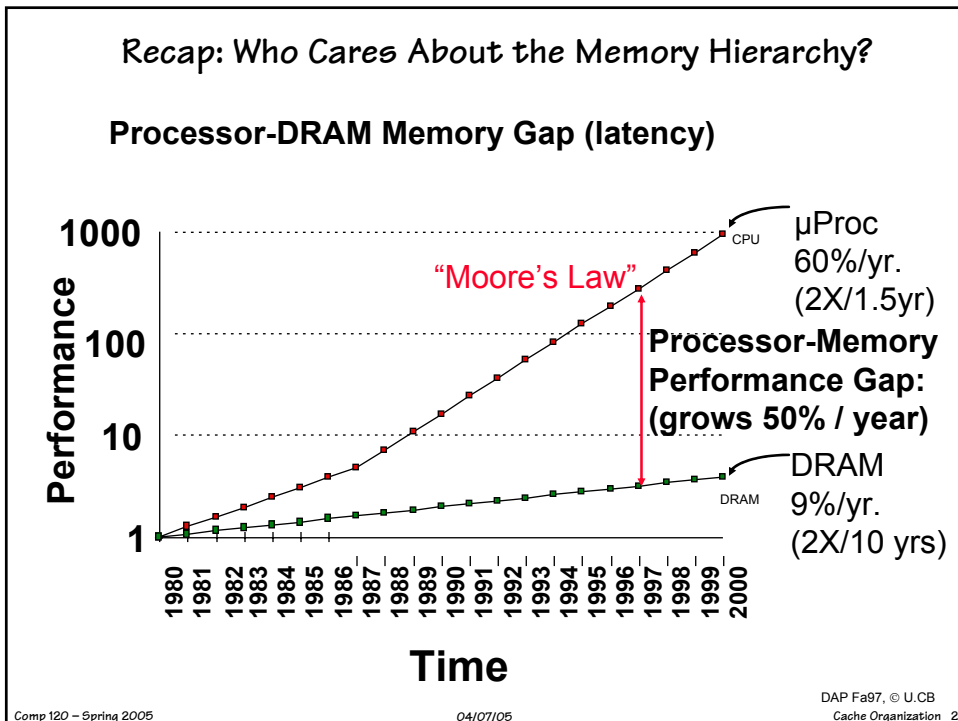
Write-back?

No, through



Study Chapter 7.1-7.3

Comp 120 – Spring 2005
04/07/05
Cache Organization 1



DRAM v. Desktop Microprocessors Cultures

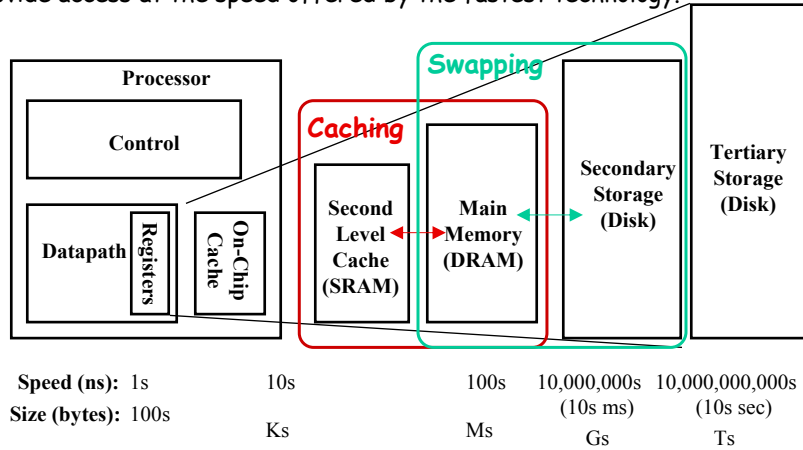
Standards	pinout, package, refresh rate, capacity, ...	binary compatibility, IEEE 754, I/O bus
Sources	Multiple	Single
Figures of Merit	1) capacity, 1a) \$/bit 2) BW, 3) latency	1) SPEC speed 2) cost
Improve Rate/year	1) 60%, 1a) 25%, 2) 20%, 3) 7%	1) 60%, 2) little change

Recap: Memory Hierarchy of a Modern Computer System

By taking advantage of the principle of locality:

Present the user with as much memory as is available in the cheapest technology.

Provide access at the speed offered by the fastest technology.



Recap:

Two Different Types of Locality:

Temporal Locality (Locality in Time): If an item is referenced, it will tend to be referenced again soon.

Spatial Locality (Locality in Space): If an item is referenced, items whose addresses are close by tend to be referenced soon.

By taking advantage of the principle of locality:

Present the user with as much memory as is available in the cheapest technology.

Provide access at the speed offered by the fastest technology.

DRAM is slow but cheap and dense:

Good choice for presenting the user with a BIG memory system

SRAM is fast but expensive and not very dense:

Good choice for providing the user FAST access time.

Cache Review

Caches improve AVERAGE memory access time

$$t_{ave} = \alpha t_c + (1 - \alpha)(t_c + t_m) = t_c + (1 - \alpha)t_m$$

where t_c = access time of a Cache HIT
 t_m = additional access time for a Cache MISS
 α = HIT ratio (ratio of references found in cache)
 $(1 - \alpha)$ = MISS ratio

e.g. $10ns + 0.02 \times 100ns = 12ns$ (98%) (caching)

e.g. $100ns + 0.000002 \times 10ms = 120ns$ (99.9998%) (swapping)

Cache Principle

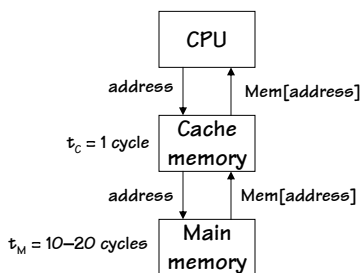
IDEA: Suppose you're making a sequence of data requests. For each request you supply some sort of tag (or key) and the cache give you back the data associated with that tag. If there is a **predictable pattern** to the sequence, it's possible to build a **cache** that quickly responds to many (or, if we're lucky, most) requests.

A cache is usually some sort of local storage that can be accessed very quickly because of its small size and/or high-bandwidth connection, i.e., much more quickly than the usual mechanism for responding to requests.

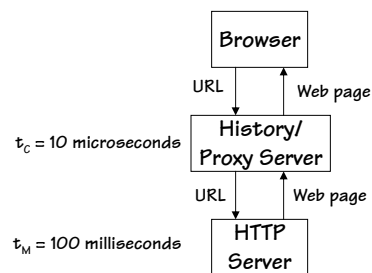
This request is made only when tag can't be found in cache, i.e., with probability $1-\alpha$, where α is the probability that the cache has the data (a "hit")



Real-world Cache Applications

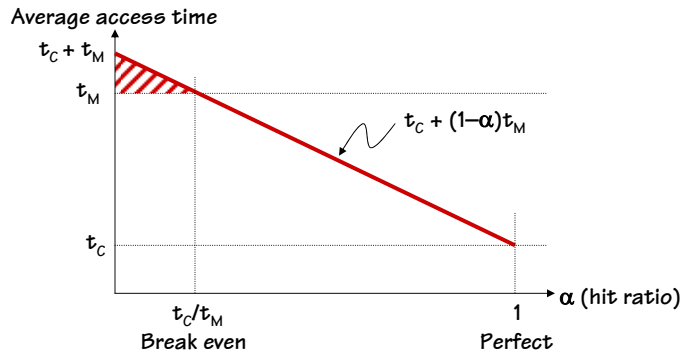


Memory system: assuming a hit ratio of 95% and $t_M=20$ cycles, average access time is $t_{ave} = 1 + (1 - .95) 20 = 2$ cycles. Cache is worthwhile, but still not as fast as we'd like (need high hit rates).



World-Wide Web: links, back/forward navigation aid prediction. Just fair cache performance reduces effective t_M (reduces latency). $t_{ave} = .01 + (1 - .2) 100 = 80$ ms. BONUS: frees up bandwidth to server by reducing contention (win-win).

Cache Economics



When the desired t_{ave} is close to t_C , an excellent hit ratio is a must. This is the regime where CPU caches operate... effectively they would like every access to come from cache.

When the desired t_{ave} is closer to t_M , and $t_C \ll t_M$, then even a modest hit ratio can make a significant difference in access time.

Taking Advantage of Access Patterns

Spatial locality:

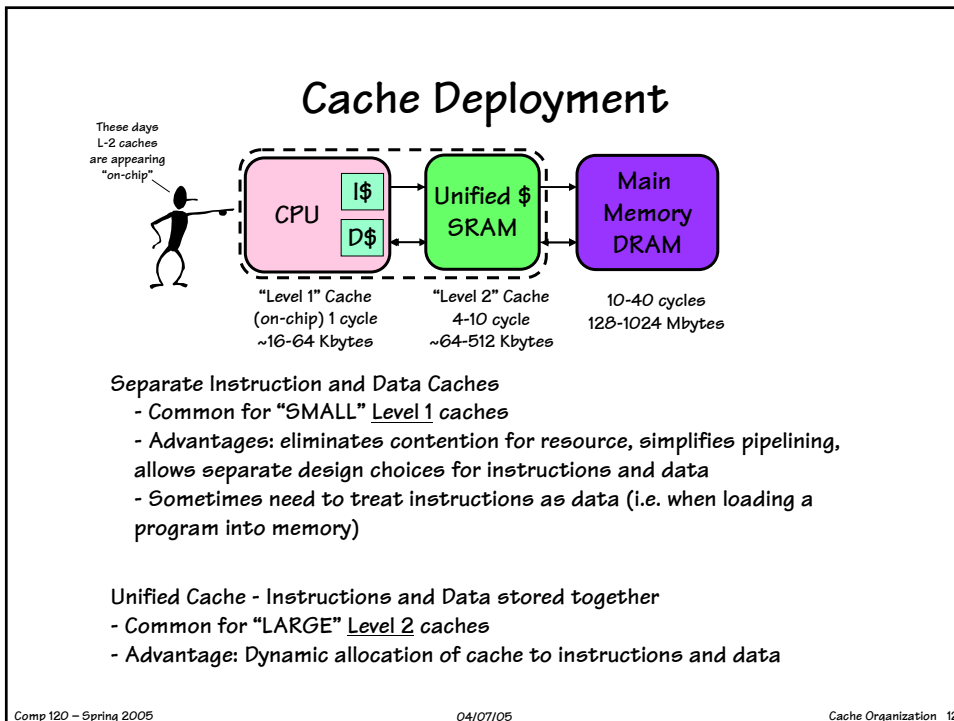
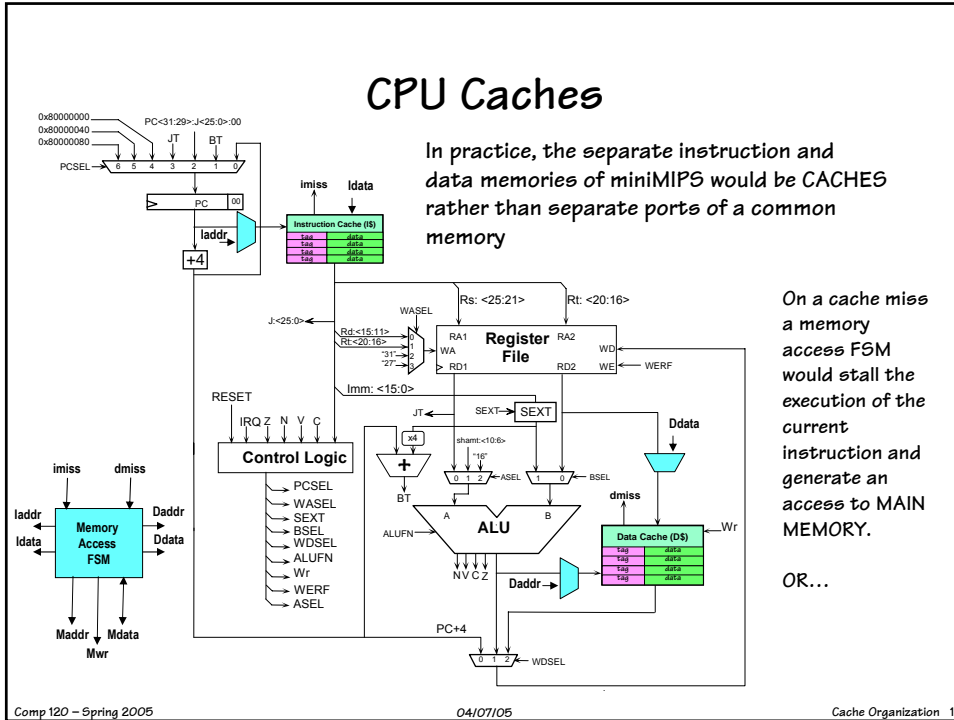
Access to X suggests that locations “near” X might soon be referenced.

- Return a **block** of “adjacent” locations for each request, e.g., a block of consecutive memory locations, several linked web pages or a JAR file of related java classes. Amortize long latencies, leverage good throughput
- Speculatively **prefetch** other nearby blocks (do this as background task so as not to preempt useful work)

Temporal locality:

Access to X suggests subsequent accesses to X.

- Keep recent accesses in cache (Obvious, but effective).
- Impacts what we throw out of cache. According to temporal locality, the **least recently used** cache element is the least likely to be accessed in the future. This observation works well when the cache is large compared to the number of requests issued before revisiting a location.



Level-2 Cache Reward

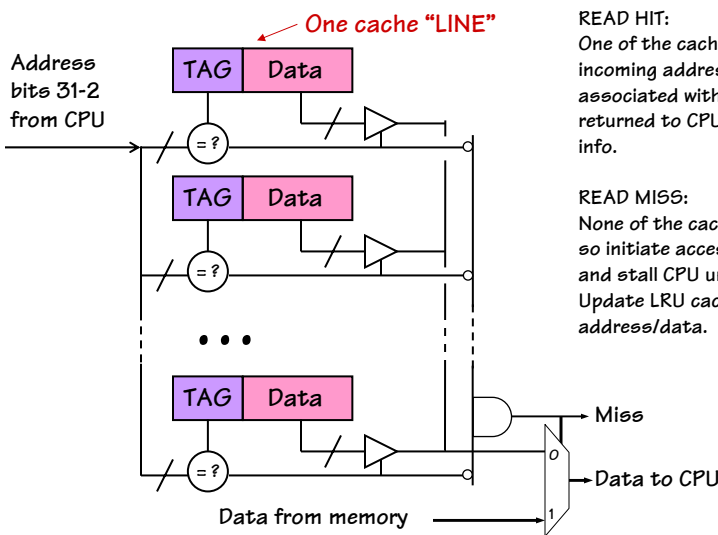
Does Level-2 Cache help?

$$t_{ave} = t_{c1} + \beta t_{c2} + \beta\beta t_m$$

$$t_{ave} = t_{c1} + \beta t_m$$

e.g. 80% of Level-2 request hit (99% of original)
 $1 + 0.05 \times 10 + 0.05 \times 0.2 \times 100 = 2.5$
 $1 + 0.05 \times 100 = 6$

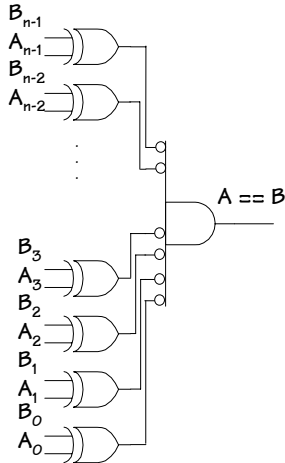
Fully Associative Cache (from last Lecture)



READ HIT:
 One of the cache tags matches the incoming address; the data associated with that tag is returned to CPU. Update usage info.

READ MISS:
 None of the cache tags matched, so initiate access to main memory and stall CPU until complete. Update LRU cache entry with address/data.

Tags Are Expensive!



- Tag comparison logic is **LARGE**
An XOR gate is as large as a one memory bit.
- Tag comparison logic is **SLOW**
High-Fan-In NOR gate
- Tag storage overhead is high

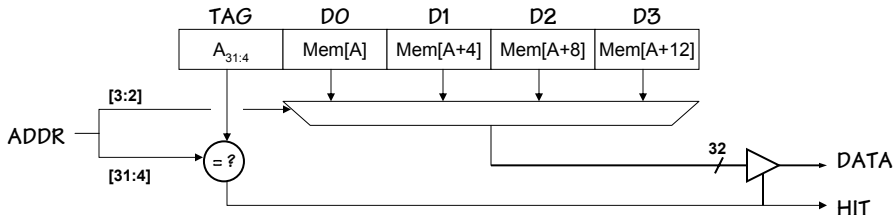
e.g. For a Fully Associative Cache

Tag bits/cache entry = 30 bits
Data bits/cache entry = 32 bits

48% of cache's memory is devoted to tag storage!

Amortize Tag Costs: More Data/Tag

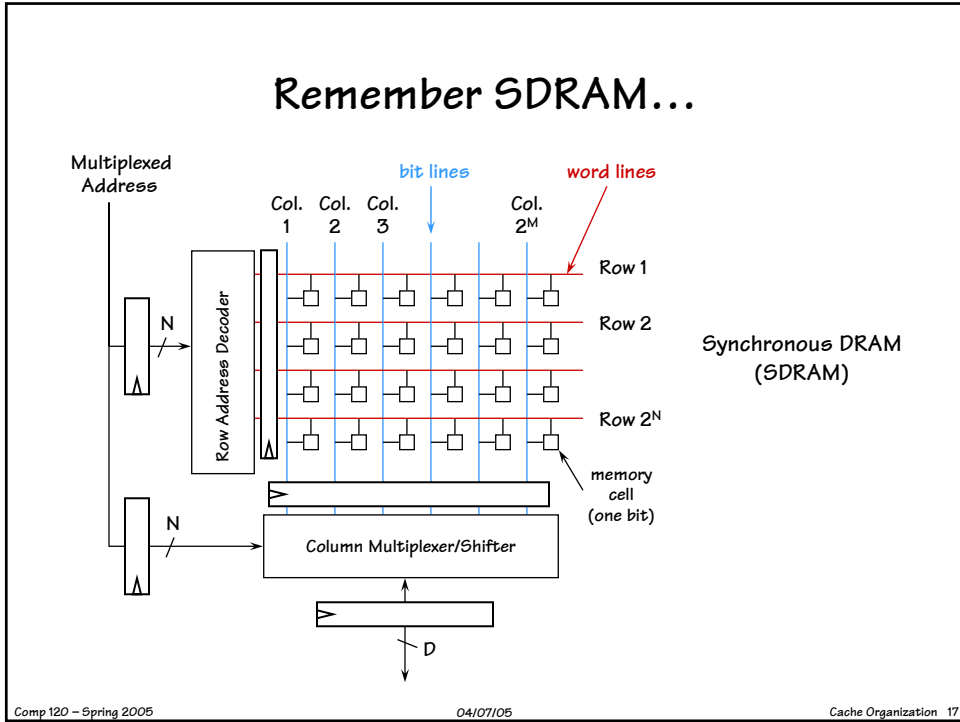
Enlarge each line in fully-associative cache:



- **blocks** of 2^B words, on 2^B word boundaries
- always read/write 2^B word **BLOCK** from/to memory
- exploits *spatial locality*: nearby words in block, likely to accessed
- *cost*: some fetches of unaccessed words
- **BIG WIN** if path to memory is wide, or *sequential accesses* are fast (SDRAM, see next slide)

Tag bits/cache entry = $(30 - 2)$ bits

Data bits/cache entry = $4 * 32$ bits Only 18% of cache's memory used for tags



Block Size Tradeoff

In general, larger block size take advantage of spatial locality

BUT:

Larger block size means larger miss penalty:

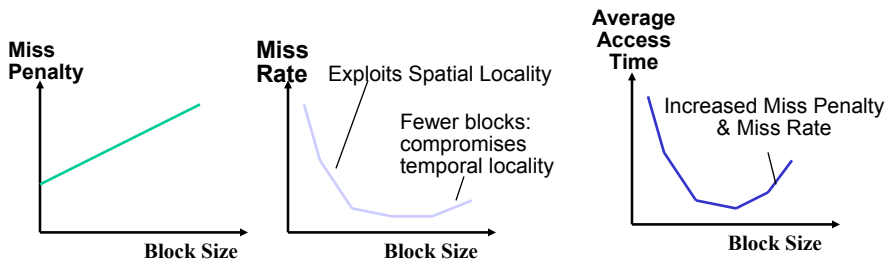
Takes longer time to fill up the block

If block size is too big relative to cache size, miss rate will go up

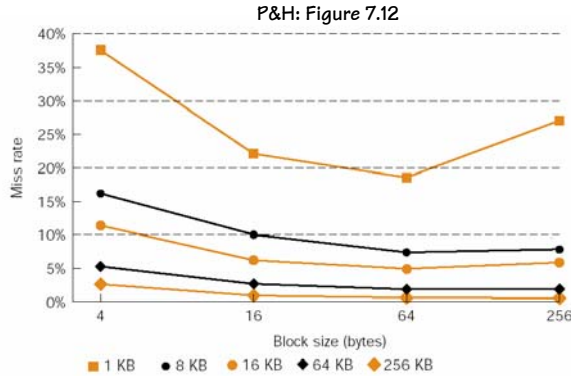
Too few cache blocks

In general, Average Access Time:

$$= \text{Hit Time} \times (1 - \text{Miss Rate}) + \text{Miss Penalty} \times \text{Miss Rate}$$

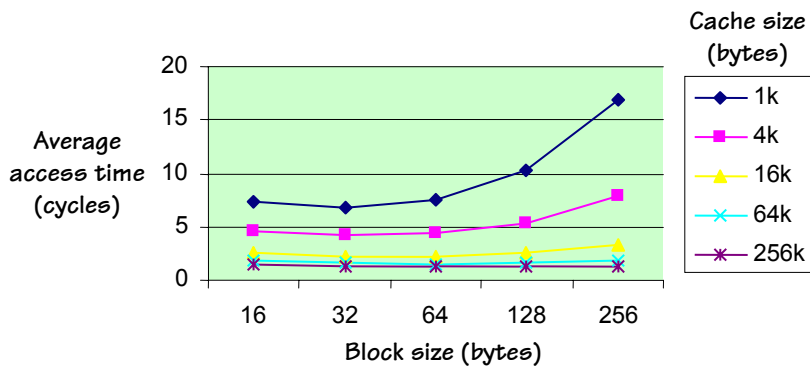


Block Size vs. Miss Rate

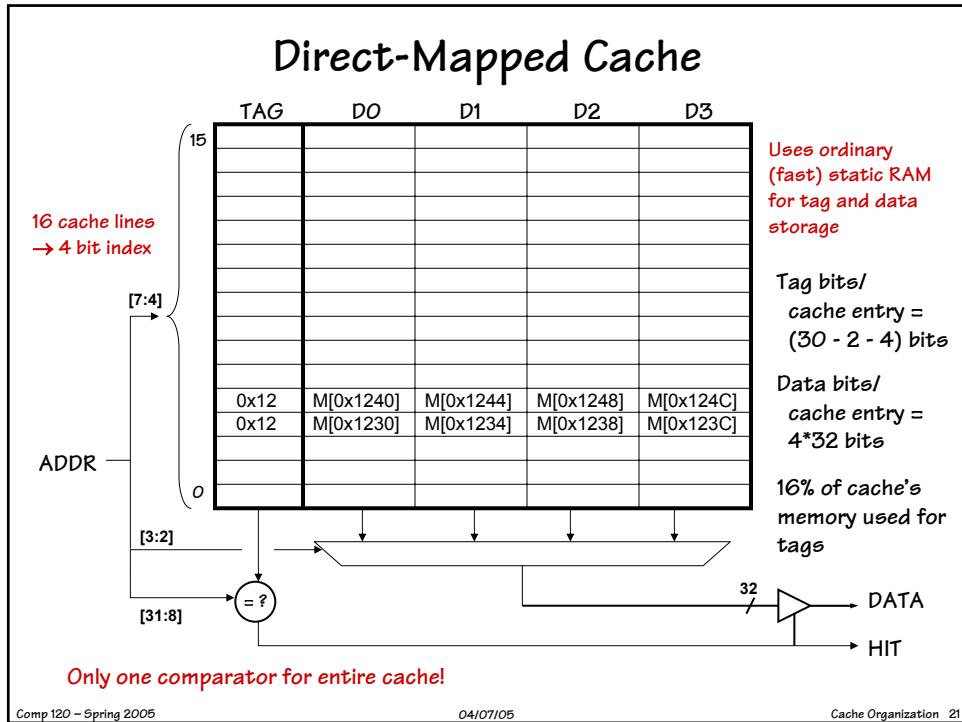


- spatial locality: larger blocks → reduce miss rate
- fixed cache size: larger blocks
 - fewer lines in cache
 - higher miss rate, especially in small caches

Block Size vs. Access Time



- Average access time = 1 + (miss rate)*(miss penalty)
- miss penalty: time it takes to read block from memory:
 - 40 cycles latency, then 16 bytes every 2 cycles



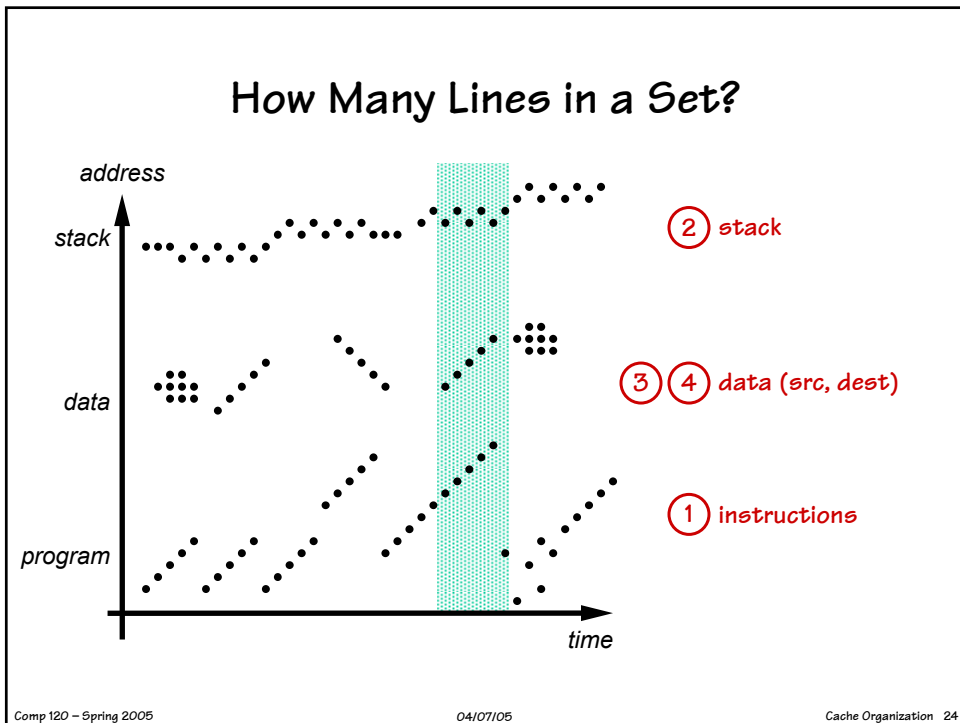
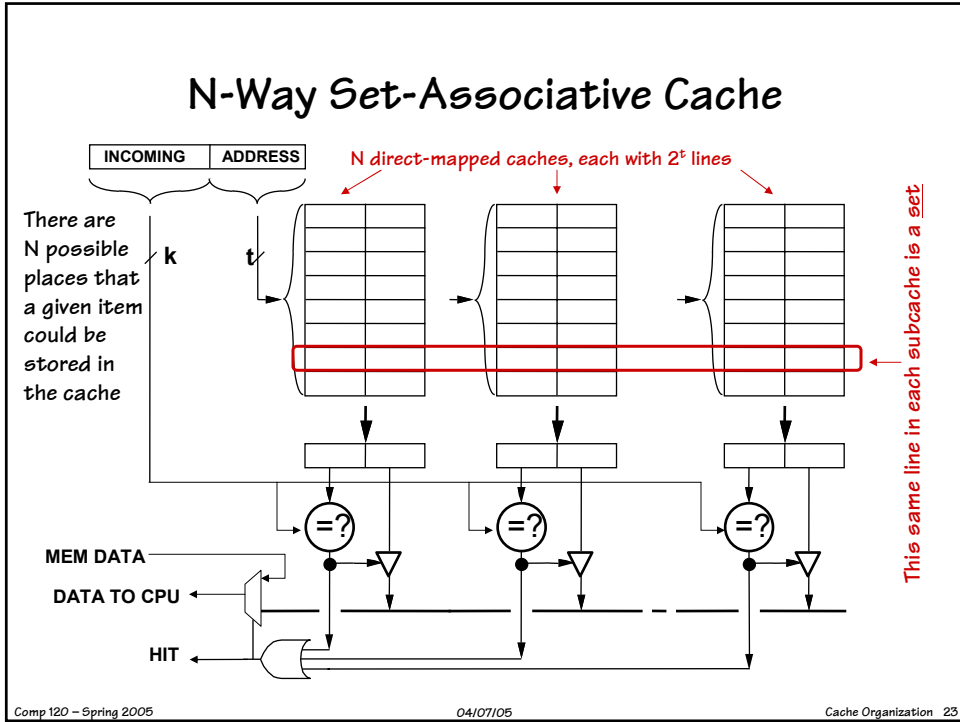
Fully-Assoc. vs. Direct-mapped

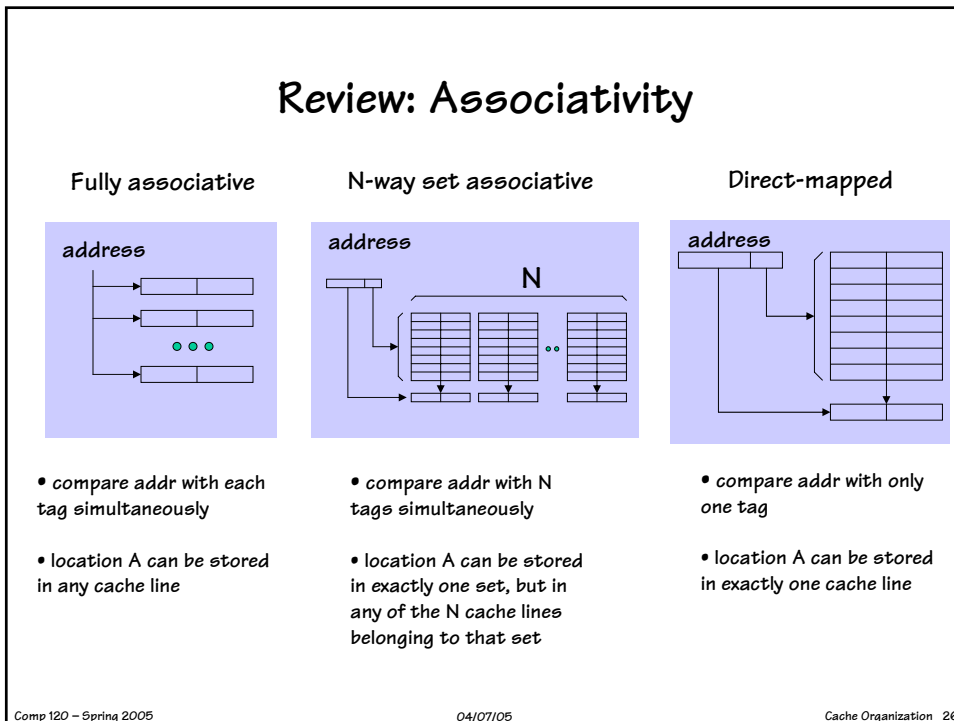
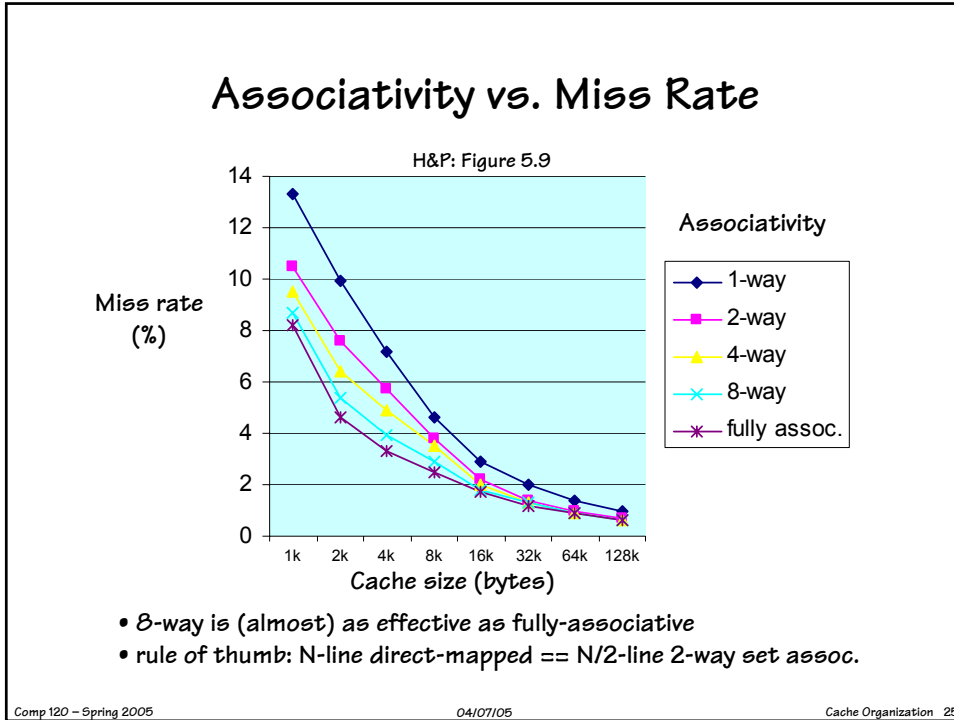
Fully-associative N-line cache:

- N tag comparators, registers used for tag/data storage (\$\$\$)
- Location A can be stored in ANY of the N cache lines; no "collisions"
- Replacement strategy (e.g., LRU=Least Recently Used) used to pick which line to use when loading new word(s) into cache

Direct-mapped N-line cache:

- 1 tag comparator, SRAM used for tag/data storage (\$)
- Location A is stored in a SPECIFIC line of the cache determined by its address; address "collisions" possible
- Replacement strategy not needed: each word can only be cached in one specific cache line





A Summary on Sources of Cache Misses

Compulsory (cold start or process migration, first reference): first access to a block

“Cold” fact of life: not a whole lot you can do about it

Note: If you are going to run “billions” of instruction, Compulsory Misses are insignificant

Conflict (collision):

Multiple memory locations mapped to the same cache location

Solution 1: increase cache size

Solution 2: increase associativity

Capacity:

Cache cannot contain all blocks accessed by the program

Solution: increase cache size

Invalidation: other process (e.g., I/O) updates memory

Sources of Cache Misses Answer

	Direct Mapped	N-way Set Associative	Fully Associative
Cache Size	Big	Medium	Small
Compulsory Miss	Same	Same	Same
Conflict Miss	High	Medium	Zero
Capacity Miss	Low	Medium	High
Invalidation Miss	Same	Same	Same

Note:

If you are going to run “billions” of instruction, Compulsory Misses are insignificant.

More \$ Next Time

- 1) *Replacement Strategies* – what cache line do we replace on a miss? What bookkeeping logic is required?
- 2) *Writes to Cache* – should we update memory on writes or only the cache?
- 3) *Cache bookkeeping* – overhead for keeping track of least-recently used cache items, “dirty” cache lines, etc.
- 4) *Cache Simulations* – How do the choices of block-size, associativity, replacement strategy, and writing policy effect cache performance