

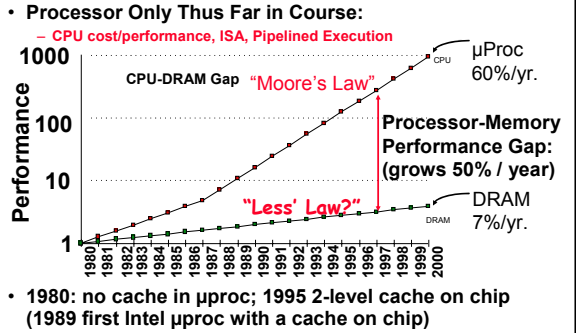
EECS 252 Graduate Computer Architecture

Lec 12 - Caches

David Culler
Electrical Engineering and Computer Sciences
University of California, Berkeley

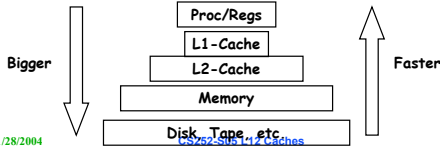
<http://www.eecs.berkeley.edu/~culler>
<http://www-inst.eecs.berkeley.edu/~cs252>

Review: Who Cares About the Memory Hierarchy?



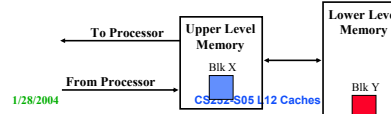
Review: What is a cache?

- Small, fast storage used to improve average access time to slow memory.
- Exploits spacial and temporal locality
- In computer architecture, almost everything is a cache!
 - Registers a cache on variables
 - First-level cache a cache on second-level cache
 - Second-level cache a cache on memory
 - Memory a cache on disk (virtual memory)
 - TLB a cache on page table
 - Branch-prediction a cache on prediction information?



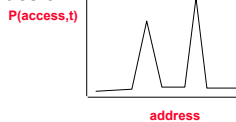
Review: Terminology

- Hit:** data appears in some block in the upper level (example: Block X)
 - Hit Rate: the fraction of memory access found in the upper level
 - Hit Time: Time to access the upper level which consists of RAM access time + Time to determine hit/miss
- Miss:** data needs to be retrieve from a block in the lower level (Block Y)
 - Miss Rate = 1 - (Hit Rate)
 - Miss Penalty: Time to replace a block in the upper level + Time to deliver the block the processor
- Hit Time << Miss Penalty (500 instructions on 21264!)



Why it works

- Exploit the statistical properties of programs
- Locality of reference
 - Temporal
 - Spatial



Average Memory Access Time

$$AMAT = HitTime + MissRate \times MissPenalty$$

$$= (HitTime_{Inst} + MissRate_{Inst} \times MissPenalty_{Inst}) + (HitTime_{Data} + MissRate_{Data} \times MissPenalty_{Data})$$

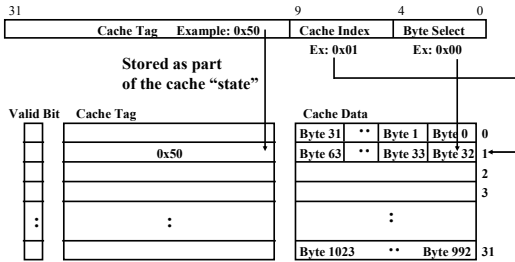
- Simple hardware structure that observes program behavior and reacts to improve future performance
- Is the cache visible in the ISA?

Block Placement

- Q1: Where can a block be placed in the upper level?
 - Fully Associative,
 - Set Associative,
 - Direct Mapped

1 KB Direct Mapped Cache, 32B blocks

- For a 2^N byte cache:
 - The uppermost $(32 - N)$ bits are always the Cache Tag
 - The lowest N bits are the Byte Select (Block Size = 2^M)



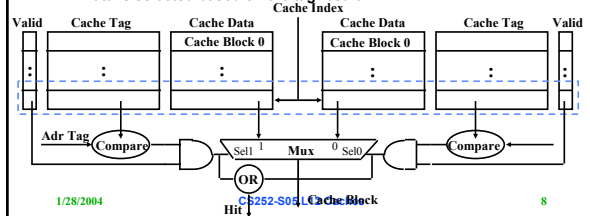
1/28/2004

CS252-S05 L12 Caches

7

Review: Set Associative Cache

- N-way set associative:** N entries for each Cache Index
 - N direct mapped caches operates in parallel
 - How big is the tag?
- Example: Two-way set associative cache**
 - Cache Index selects a "set" from the cache
 - The two tags in the set are compared to the input in parallel
 - Data is selected based on the tag result



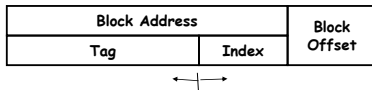
1/28/2004

CS252-S05 L12 Caches

8

Q2: How is a block found if it is in the upper level?

- Index identifies set of possibilities
- Tag on each block
 - No need to check index or block offset
- Increasing associativity shrinks index, expands tag



$$\text{Cache size} = \text{Associativity} * 2^{\text{index_size}} * 2^{\text{offset_size}}$$

1/28/2004

CS252-S05 L12 Caches

9

Q3: Which block should be replaced on a miss?

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

Assoc:	2-way		4-way		8-way	
Size	LRU	Ran	LRU	Ran	LRU	Ran
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%

1/28/2004

CS252-S05 L12 Caches

10

Q4: What happens on a write?

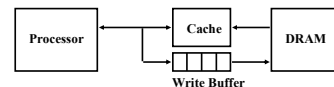
- Write through**—The information is written to both the block in the cache and to the block in the lower-level memory.
- Write back**—The information is written only to the block in the cache. The modified cache block is written to main memory only when it is replaced.
 - is block clean or dirty?
- Pros and Cons of each?
 - WT: read misses cannot result in writes
 - WB: no repeated writes to same location
- WT always combined with write buffers so that don't wait for lower level memory
- What about on a miss?
 - Write_no_allocate vs write_allocate

1/28/2004

CS252-S05 L12 Caches

11

Write Buffer for Write Through



- A Write Buffer is needed between the Cache and Memory
 - Processor: writes data into the cache and the write buffer
 - Memory controller: write contents of the buffer to memory
- Write buffer is just a FIFO:
 - Typical number of entries: 4
 - Works fine if: Store frequency (w.r.t. time) \ll 1 / DRAM write cycle

1/28/2004

CS252-S05 L12 Caches

12

Review: Cache performance

Miss-oriented Approach to Memory Access:

$$CPU\ time = IC \times \left(CPI_{Execution} + \frac{Mem\ Access}{Inst} \times Miss\ Rate \times Miss\ Penalty \right) \times Cycle\ Time$$

Separating out Memory component entirely

– AMAT = Average Memory Access Time

$$CPU\ time = IC \times \left(CPI_{ArithOps} + \frac{Mem\ Access}{Inst} \times AMAT \right) \times Cycle\ Time$$

– Effective CPI = $CPI_{ideal_mem} + P_{mem} * AMAT$

1/28/2004

CS252-S05 L12 Caches

13

Impact on Performance

- Suppose a processor executes at
 - Clock Rate = 200 MHz (5 ns per cycle), Ideal (no misses) CPI = 1.1
 - 50% arith/logic, 30% ld/st, 20% control
- Suppose that 10% of memory operations get 50 cycle miss penalty
- Suppose that 1% of instructions get same miss penalty
- CPI = ideal 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
- 58% of the time the proc is stalled waiting for memory!
- AMAT = $(1/1.3) \times [1 + 0.01 \times 50] + (0.3/1.3) \times [1 + 0.1 \times 50] = 2.54$

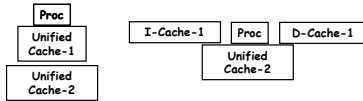
1/28/2004

CS252-S05 L12 Caches

14

Example: Harvard Architecture

Unified vs Separate I&D (Harvard)



Statistics (given in H&P):

- 16KB I&D: Inst miss rate=0.64%, Data miss rate=6.47%
- 32KB unified: Aggregate miss rate=1.99%

Which is better (ignore L2 cache)?

- Assume 33% data ops \Rightarrow 75% accesses from instructions (1.0/1.33)
- hit time=1, miss time=50
- Note that *data* hit has 1 stall for unified cache (only one port)

$$AMAT_{Harvard} = 75\% \times (1 + 0.64\% \times 50) + 25\% \times (1 + 6.47\% \times 50) = 2.05$$

$$AMAT_{Unified} = 75\% \times (1 + 1.99\% \times 50) + 25\% \times (1 + 1 + 1.99\% \times 50) = 2.24$$

1/28/2004

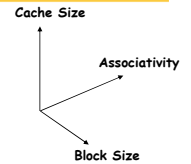
CS252-S05 L12 Caches

15

The Cache Design Space

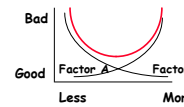
Several interacting dimensions

- cache size
- block size
- associativity
- replacement policy
- write-through vs write-back



The optimal choice is a compromise

- depends on access characteristics
 - » workload
 - » use (I-cache, D-cache, TLB)
- depends on technology / cost



Simplicity often wins

1/28/2004

CS252-S05 L12 Caches

16

Review: Improving Cache Performance

$$CPU\ time = IC \times \left(CPI_{ArithOps} + \frac{Memory\ accesses}{Instruction} \times Miss\ rate \times Miss\ penalty \right) \times Clock\ cycle\ time$$

1. Reduce the miss rate.

2. Reduce the miss penalty, or

3. Reduce the time to hit in the cache.

1/28/2004

CS252-S05 L12 Caches

17

Reducing Misses

Classifying Misses: 3 Cs

- **Compulsory**—The first access to a block is not in the cache, so the block must be brought into the cache. Also called **cold start misses** or **first reference misses**.
(Misses in even an Infinite Cache)
- **Capacity**—If the cache cannot contain all the blocks needed during execution of a program, **capacity misses** will occur due to blocks being discarded and later retrieved.
(Misses in Fully Associative Size X Cache)
- **Conflict**—If block-placement strategy is set associative or direct mapped, conflict misses (in addition to compulsory & capacity misses) will occur because a block can be discarded and later retrieved if too many blocks map to its set. Also called **collision misses** or **interference misses**.
(Misses in N-way Associative, Size X Cache)

More recent, 4th “C”:

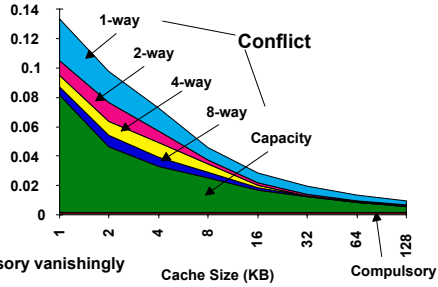
- **Coherence** - Misses caused by cache coherence.

1/28/2004

CS252-S05 L12 Caches

18

3Cs Absolute Miss Rate (SPEC92)



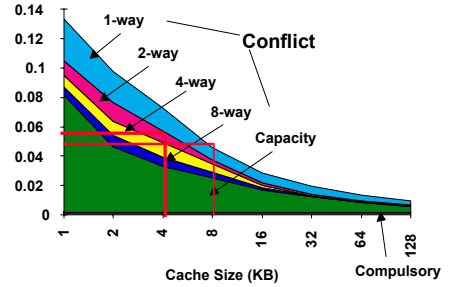
1/28/2004

CS252-S05 L12 Caches

19

2:1 Cache Rule

miss rate 1-way associative cache size X
 ~ miss rate 2-way associative cache size X/2

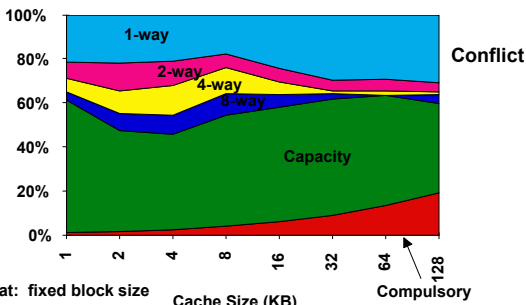


1/28/2004

CS252-S05 L12 Caches

20

3Cs Relative Miss Rate



1/28/2004

CS252-S05 L12 Caches

21

How Can Reduce Misses?

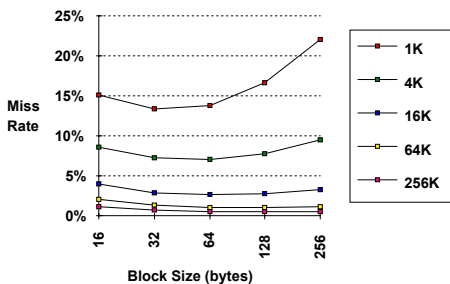
- 3 Cs: **Compulsory, Capacity, Conflict**
- In all cases, assume total cache size not changed:
- What happens if:
 - 1) Change Block Size:
Which of 3Cs is obviously affected?
 - 2) Change Associativity:
Which of 3Cs is obviously affected?
 - 3) Change Algorithm / Compiler:
Which of 3Cs is obviously affected?

1/28/2004

CS252-S05 L12 Caches

22

1. Reduce Misses via Larger Block Size



1/28/2004

CS252-S05 L12 Caches

23

2. Reduce Misses via Higher Associativity

- **2:1 Cache Rule:**
 - Miss Rate DM cache size N ~ Miss Rate 2-way cache size N/2
- **Beware: Execution time is only final measure!**
 - Will Clock Cycle time increase?
 - Hill [1988] suggested hit time for 2-way vs. 1-way external cache +10%, internal + 2%

1/28/2004

CS252-S05 L12 Caches

24

Example: Avg. Memory Access Time vs. Miss Rate

- assume CCT = 1.10 for 2-way, 1.12 for 4-way, 1.14 for 8-way vs. CCT direct mapped

Cache Size (KB)	Associativity			
	1-way	2-way	4-way	8-way
1	2.33	2.15	2.07	2.01
2	1.98	1.86	1.76	1.68
4	1.72	1.67	1.61	1.53
8	1.46	1.48	1.47	1.43
16	1.29	1.32	1.32	1.32
32	1.20	1.24	1.25	1.27
64	1.14	1.20	1.21	1.23
128	1.10	1.17	1.18	1.20

(Red means A.M.A.T. not improved by more associativity)

1/28/2004

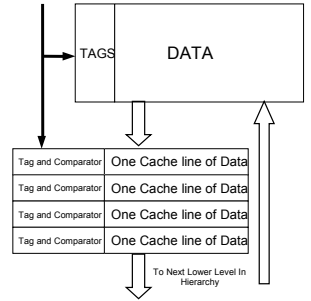
CS252-S05 L12 Caches

25

3. Reducing Misses via a "Victim Cache"

- How to combine fast hit time of direct mapped yet still avoid conflict misses?

- Add buffer to place data discarded from cache
- Jouppi [1990]: 4-entry victim cache removed 20% to 95% of conflicts for a 4 KB direct mapped data cache
- Used in Alpha, HP machines



1/28/2004

CS252-S05 L12 Caches

26

4. Reducing Misses via "Pseudo-Associativity"

- How to combine fast hit time of Direct Mapped and have the lower conflict misses of 2-way SA cache?
- Divide cache: on a miss, check other half of cache to see if there, if so have a **pseudo-hit** (slow hit)



- Drawback: CPU pipeline is hard if hit takes 1 or 2 cycles
 - Better for caches not tied directly to processor (L2)
 - Used in MIPS R1000 L2 cache, similar in UltraSPARC

1/28/2004

CS252-S05 L12 Caches

27

5. Reducing Misses by Hardware Prefetching of Instructions & Data

- E.g., Instruction Prefetching
 - Alpha 21064 fetches 2 blocks on a miss
 - Extra block placed in "stream buffer"
 - On miss check stream buffer
- Works with data blocks too:
 - Jouppi [1990] 1 data stream buffer got 25% misses from 4KB cache; 4 streams got 43%
 - Palacharla & Kessler [1994] for scientific programs for 8 streams got 50% to 70% of misses from 2 64KB, 4-way set associative caches
- Prefetching relies on having extra memory bandwidth that can be used without penalty

1/28/2004

CS252-S05 L12 Caches

28

6. Reducing Misses by Software Prefetching Data

- Data Prefetch
 - Load data into register (HP PA-RISC loads)
 - Cache Prefetch: load into cache (MIPS IV, PowerPC, SPARC v. 9)
 - Special prefetching instructions cannot cause faults; a form of speculative execution
- Issuing Prefetch Instructions takes time
 - Is cost of prefetch issues < savings in reduced misses?
 - Higher superscalar reduces difficulty of issue bandwidth

1/28/2004

CS252-S05 L12 Caches

29

7. Reducing Misses by Compiler Optimizations

- McFarling [1989] reduced caches misses by 75% on 8KB direct mapped cache, 4 byte blocks **in software**
- Instructions
 - Reorder procedures in memory so as to reduce conflict misses
 - Profiling to look at conflicts (using tools they developed)
- Data
 - Merging Arrays**: improve spatial locality by single array of compound elements vs. 2 arrays
 - Loop Interchange**: change nesting of loops to access data in order stored in memory
 - Loop Fusion**: Combine 2 independent loops that have same looping and some variables overlap
 - Blocking**: Improve temporal locality by accessing "blocks" of data repeatedly vs. going down whole columns or rows

1/28/2004

CS252-S05 L12 Caches

30

Merging Arrays Example

```

/* Before: 2 sequential arrays */
int val[SIZE];
int key[SIZE];

/* After: 1 array of structures */
struct merge {
    int val;
    int key;
};
struct merge merged_array[SIZE];
    
```

Reducing conflicts between val & key;
improve spatial locality

1/28/2004

CS252-S05 L12 Caches

31

Loop Interchange Example

```

/* Before */
for (k = 0; k < 100; k = k+1)
    for (j = 0; j < 100; j = j+1)
        for (i = 0; i < 5000; i = i+1)
            x[i][j] = 2 * x[i][j];

/* After */
for (k = 0; k < 100; k = k+1)
    for (i = 0; i < 5000; i = i+1)
        for (j = 0; j < 100; j = j+1)
            x[i][j] = 2 * x[i][j];
    
```

Sequential accesses instead of striding through
memory every 100 words; improved spatial
locality

1/28/2004

CS252-S05 L12 Caches

32

Loop Fusion Example

```

/* Before */
for (i = 0; i < N; i = i+1)
    for (j = 0; j < N; j = j+1)
        a[i][j] = 1/b[i][j] * c[i][j];
for (i = 0; i < N; i = i+1)
    for (j = 0; j < N; j = j+1)
        d[i][j] = a[i][j] + c[i][j];

/* After */
for (i = 0; i < N; i = i+1)
    for (j = 0; j < N; j = j+1)
        {
            a[i][j] = 1/b[i][j] * c[i][j];
            d[i][j] = a[i][j] + c[i][j];
        }
    
```

2 misses per access to a & c vs. one miss per access;
improve spatial locality

1/28/2004

CS252-S05 L12 Caches

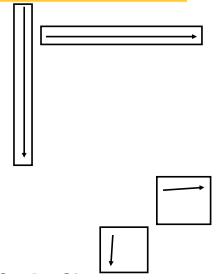
33

Blocking Example

```

/* Before */
for (i = 0; i < N; i = i+1)
    for (j = 0; j < N; j = j+1)
        {
            r = 0;
            for (k = 0; k < N; k = k+1) {
                r = r + y[i][k]*z[k][j];
            }
            x[i][j] = r;
        }
    
```

- Two Inner Loops:
 - Read all NxN elements of z[]
 - Read N elements of 1 row of y[] repeatedly
 - Write N elements of 1 row of x[]
- Capacity Misses a function of N & Cache Size:
 - $2N^3 + N^2 \Rightarrow$ (assuming no conflict; otherwise ...)
- Idea: compute on BxB submatrix that fits



1/28/2004

CS252-S05 L12 Caches

34

Blocking Example

```

/* After */
for (jj = 0; jj < N; jj = jj+B)
    for (kk = 0; kk < N; kk = kk+B)
        for (i = 0; i < N; i = i+1)
            for (j = jj; j < min(jj+B-1, N); j = j+1)
                {
                    r = 0;
                    for (k = kk; k < min(kk+B-1, N); k = k+1) {
                        r = r + y[i][k]*z[k][j];
                    }
                    x[i][j] = x[i][j] + r;
                }
    
```

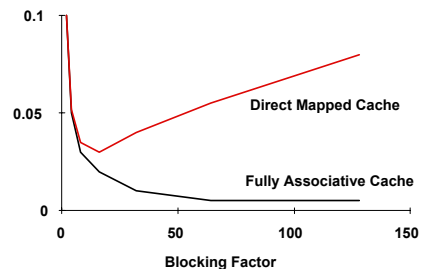
- B called **Blocking Factor**
- Capacity Misses from $2N^3 + N^2$ to $2N^3/B + N^2$
- Conflict Misses Too?

1/28/2004

CS252-S05 L12 Caches

35

Reducing Conflict Misses by Blocking



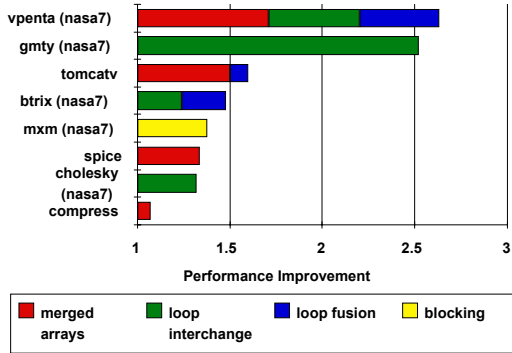
- Conflict misses in caches not FA vs. Blocking size
 - Lam et al [1991] a blocking factor of 24 had a fifth the misses vs. 48 despite both fit in cache

1/28/2004

CS252-S05 L12 Caches

36

Summary of Compiler Optimizations to Reduce Cache Misses (by hand)



Impact of Memory Hierarchy on Algorithms

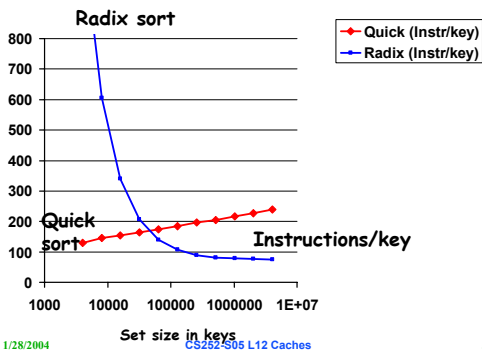
- Today CPU time is a function of (ops, cache misses) vs. just f(ops): What does this mean to Compilers, Data structures, Algorithms?
- "The Influence of Caches on the Performance of Sorting" by A. LaMarca and R.E. Ladner. *Proceedings of the Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*, January, 1997, 370-379.
- Quicksort: fastest comparison based sorting algorithm when all keys fit in memory
- Radix sort: also called "linear time" sort because for keys of fixed length and fixed radix a constant number of passes over the data is sufficient independent of the number of keys
- For Alphastation 250, 32 byte blocks, direct mapped L2 2MB cache, 8 byte keys, from 4000 to 4000000

1/28/2004

CS252-S05 L12 Caches

38

Quicksort vs. Radix as vary number keys: Instructions

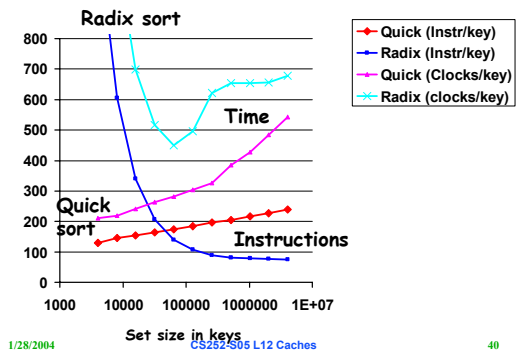


1/28/2004

CS252-S05 L12 Caches

39

Quicksort vs. Radix as vary number key Instrs & Time

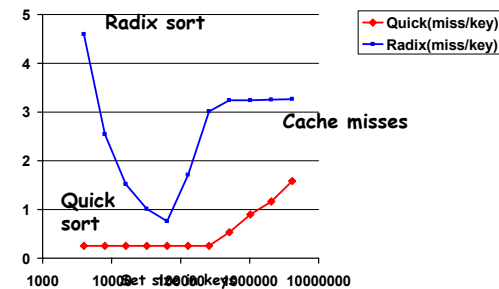


1/28/2004

CS252-S05 L12 Caches

40

Quicksort vs. Radix as vary number keys Cache misses



What is proper approach to fast algorithms?

1/28/2004

CS252-S05 L12 Caches

41

Review: What happens on Cache miss?

- For in-order pipeline, 2 options:
 - Freeze pipeline in Mem stage (popular early on: Sparc, R4000)


```
IF ID EX Mem stall stall stall ... stall Mem Wr
IF ID EX stall stall stall ... stall Ex Mem Wr
```

 - » Stall, Load cache line, Restart mem stage
 - » This is why cost on CM = Penalty + Hit Time
 - Use Full/Empty bits in registers + MSHR queue
 - » MSHR = "Miss Status/Handler Registers" (Kroft)
 - Each entry in this queue keeps track of status of outstanding memory requests to one complete memory line.
 - Per cache-line: keep info about memory address.
 - For each word: register (if any) that is waiting for result.
 - Used to "merge" multiple requests to one memory line
 - » New load creates MSHR entry and sets destination register to "Empty". Load is "released" from pipeline.
 - » Attempt to use register before result returns causes instruction to block in decode stage.
 - » Limited "out-of-order" execution with respect to loads. Popular with in-order superscalar architectures.
- Out-of-order pipelines already have this functionality built in (load queues, etc)

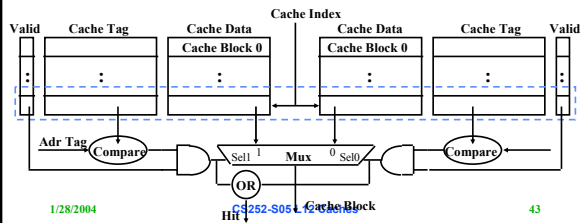
1/28/2004

CS252-S05 L12 Caches

42

Disadvantage of Set Associative Cache

- **N-way Set Associative Cache v. Direct Mapped Cache:**
 - N comparators vs. 1
 - Extra MUX delay for the data
 - Data comes AFTER Hit/Miss
- **In a direct mapped cache, Cache Block is available BEFORE Hit/Miss:**
 - Possible to assume a hit and continue. Recover later if miss.



Review: Four Questions for Memory Hierarchy Designers

- **Q1: Where can a block be placed in the upper level?**
(Block placement)
 - Fully Associative, Set Associative, Direct Mapped
- **Q2: How is a block found if it is in the upper level?**
(Block identification)
 - Tag/Block
- **Q3: Which block should be replaced on a miss?**
(Block replacement)
 - Random, LRU
- **Q4: What happens on a write?**
(Write strategy)
 - Write Back or Write Through (with Write Buffer)

1/28/2004 CS252-S05 L12 Caches 44

Summary

$$CPUtime = IC \times \left(CPI_{memory} + \frac{Memory\ accesses}{Instruction} \times Miss\ rate \times Miss\ penalty \right) \times Clock\ cycle\ time$$

- **3 Cs: Compulsory, Capacity, Conflict**
 1. Reduce Misses via Larger Block Size
 2. Reduce Misses via Higher Associativity
 3. Reducing Misses via Victim Cache
 4. Reducing Misses via Pseudo-Associativity
 5. Reducing Misses by HW Prefetching Instr, Data
 6. Reducing Misses by SW Prefetching Data
 7. Reducing Misses by Compiler Optimizations
- **Remember danger of concentrating on just one parameter when evaluating performance**

1/28/2004 CS252-S05 L12 Caches 45