

Lecture 17: Memory Caching

COS 471a, COS 471b / ELE 375

Computer Architecture and Organization

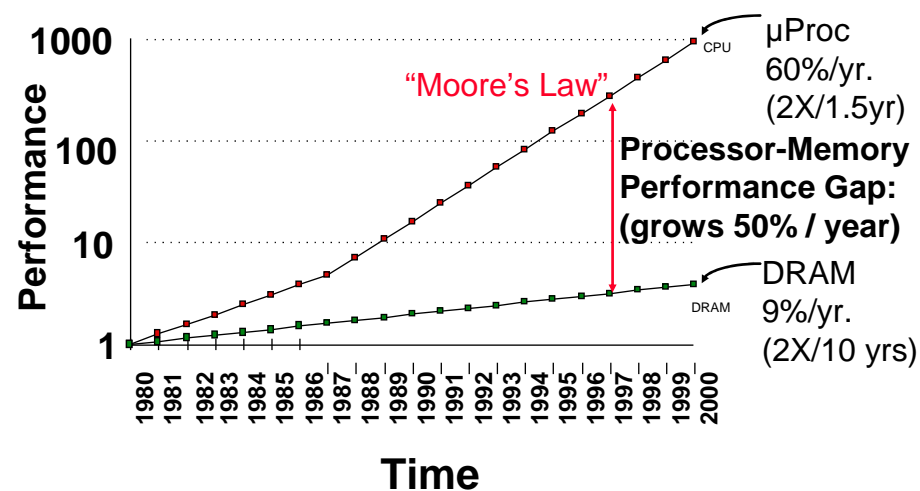
Princeton University
Fall 2004

Prof. David August

1

We still have a problem!

Processor-DRAM Performance Gap (latency)



5

Caching and The Principle of Locality

- Program access a relatively small portion of the address space at any instant of time. (90-10 rule)

Temporal Locality

- If an item is referenced, it will tend to be referenced again soon

Spatial Locality

- If an item is referenced, nearby items will tend to be referenced soon

USE CACHES!

6

Spatial Locality in Instruction & Data

Instruction and Data References have distinct behavior:

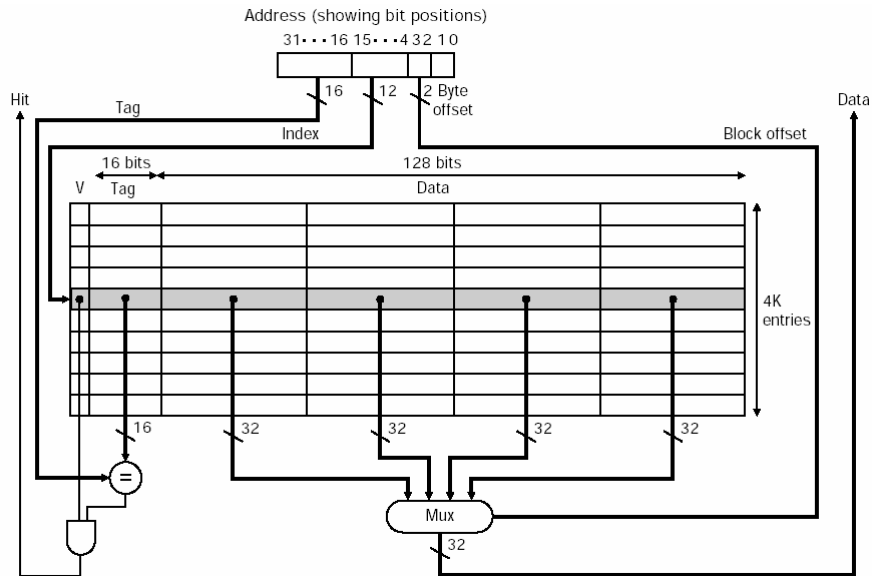
Program	Block size in words	Instruction miss rate	Data miss rate	Effective combined miss rate
gcc	1	6.1%	2.1%	5.4%
	4	2.0%	1.7%	1.9%
spice	1	1.2%	1.3%	1.2%
	4	0.3%	0.6%	0.4%

Split Instruction and Data Caches

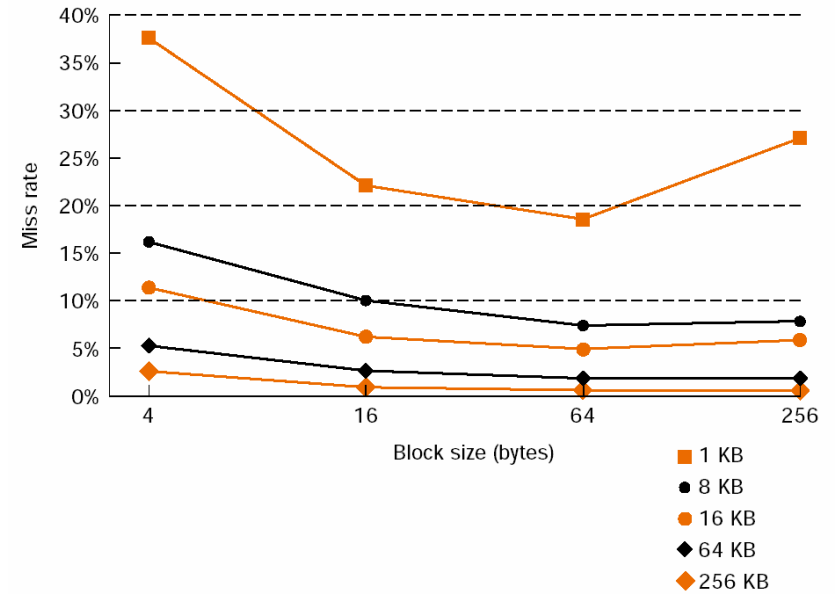
- Optimize for behavior
- Smaller caches are faster
- Problem - when data is code or code is data

7

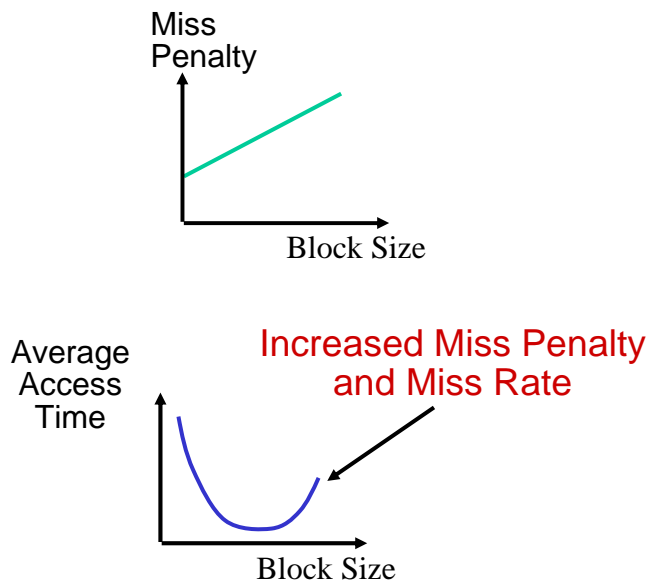
Direct Mapped Cache, Increased Block Size Capture Spatial Locality



Block Size Increase: Miss Rate



Block Size Increase: Overall Performance



Block Size Increase: Fill Time

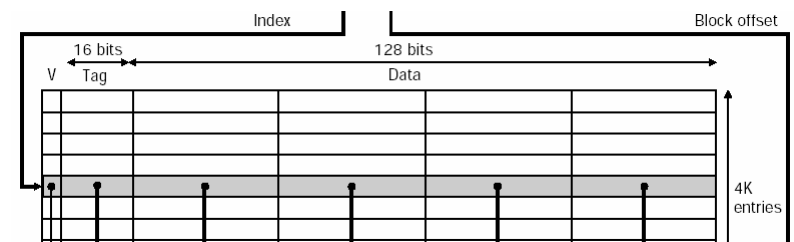
Larger Block Size à Must Wait for Block to Fill

Early Restart

- Deliver word to process/continue execution when word requested is delivered.

Critical Word First

- Early Restart and Fetch the requested word first.



The Three Types of Cache Misses

1. Conflict Misses

- Two distinct memory addresses map to the same cache location
- Big problem in direct-mapped caches

How do we reduce these?

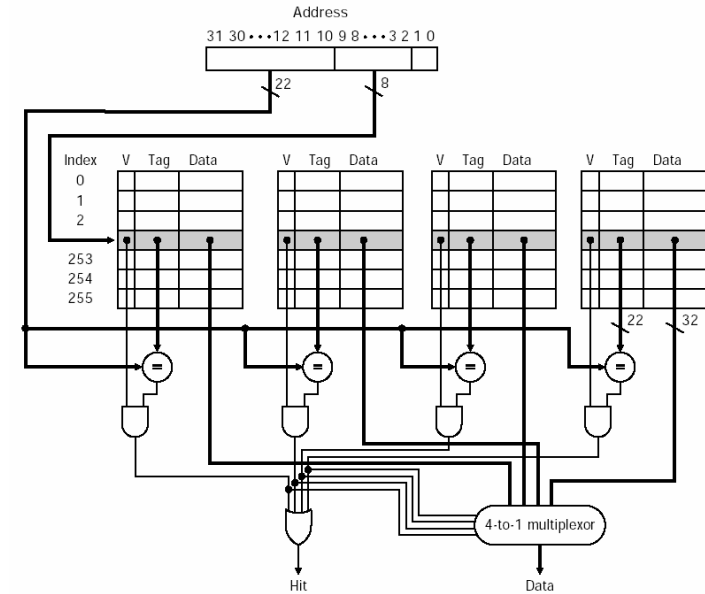
Solution 1: Make cache bigger (limits)

Solution 2: ...

12

Four-Way Set Associative Cache

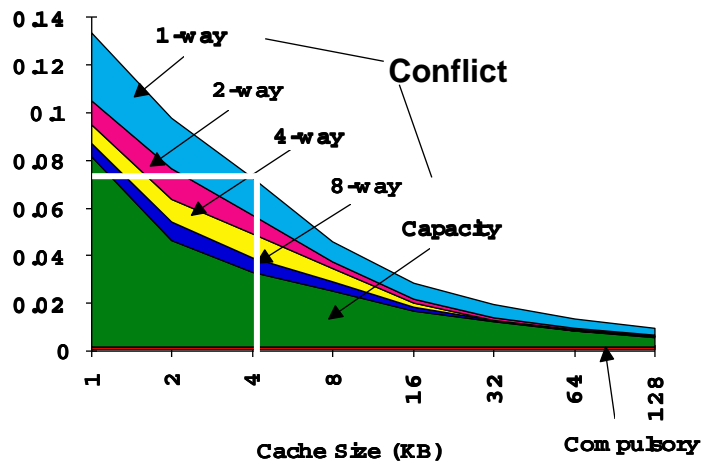
Avoid Conflicts



13

2:1 Cache Rule

Rule of Thumb: a direct-mapped cache of size N has about the same miss rate as a 2-way set associative cache of size $N/2$.



14

The Three Types of Cache Misses

2. Capacity Misses

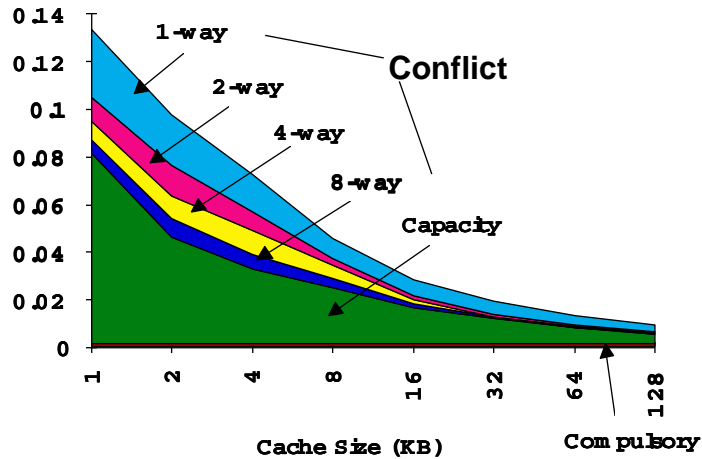
- Occurs because the cache has a limited size
- Increase the size of the cache, it goes away
- Sketchy definition, so just get the general idea
- Easy to understand in Fully Associative Caches.

How do we reduce these?

15

Capacity Misses

Fully Associative Cache yields no conflict misses.



16

The Three Types of Cache Misses

3. Compulsory Misses

- Occur when a program is first started
- Cache does not contain any of program's data yet

How do we reduce these?

17

Prefetching!

Reduces all types of misses, including "compulsory"!

Original Code:

```
for(y = 0; y < SIZE_Y; y++)
  for(x = 0; x < SIZE_X; x++)
    sum += Array[x][y];
```

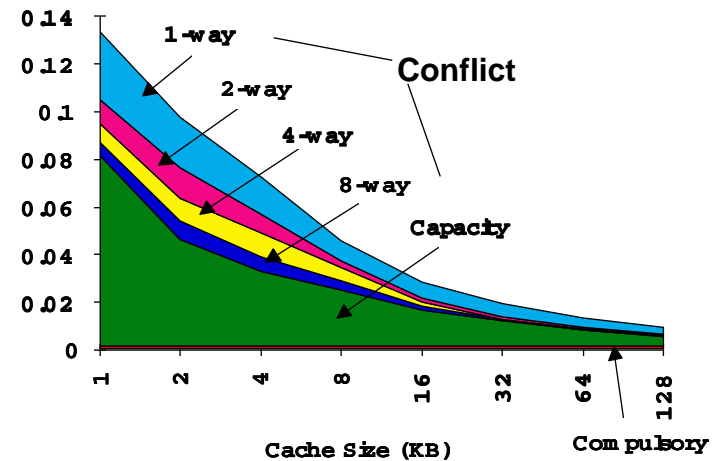
Code with Prefetching (ignoring boundary condition):

```
for(y = 0; y < SIZE_Y; y++)
  for(x = 0; x < SIZE_X; x++) {
    junk = Array[x+16][y];
    sum += Array[x][y];
  }
```

18

Compulsory Misses

Fully Associative Cache yields no conflict misses.



19

3C Summary

Compulsory misses (cold start)

- Cold fact of life
- First time data is referenced
- Run billions of instructions, become insignificant

Capacity misses

- Working set is larger than cache size
- Solution: increase cache size

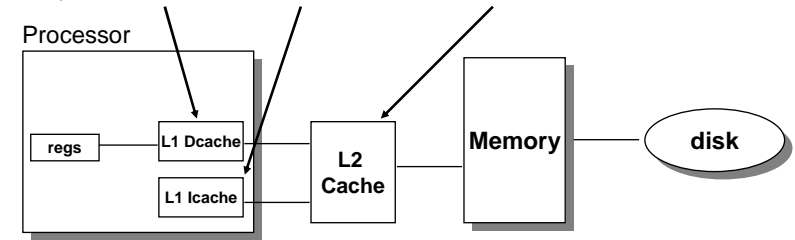
Conflict misses

- Multiple memory locations mapped to the same location
- One set fills up, but space in other cache sets
- Solution 1: increase cache size
- Solution 2: increase associative indexes

20

Multi-Level Caches

Options: *separate* data and instruction caches, or a *unified* cache



Inclusive vs. Exclusive

Sample Sizes:

- L1: 32KB, 32 Byte Lines, 4-Way Set Associative
- L2: 256KB, 128 Byte Lines, 8-Way Set Associative
- L3: 4MB, 256 Byte Lines, Direct Mapped

22

Split Instruction and Data Caches

Self-Modifying Code!?!?

- Ignore problem, software must flush cache
- Permit duplicate lines: invalidate I-cache line on write
- Do not permit duplicate lines: data is exclusive to D- or I-Cache

- Page Faults - More next week

23

Handling Writes in Caches

First, Two Observations:

1. Writes change state à wait until exceptions are cleared
2. Stores aren't the source of a dependence - latency tolerant

Typical Implementation Decisions:

- Cache write policy?
 - Write-Through
 - Write-Back
 - Write-Around
- Include a Write buffer?
 - Small pseudo-FIFO buffer alongside cache

25

Write-Back vs. Write-Through Caches

Write back

- Writes only go into top level of hierarchy
- Maintain a record of "dirty" lines
- Faster write speed (only has to go to top level to be considered complete)

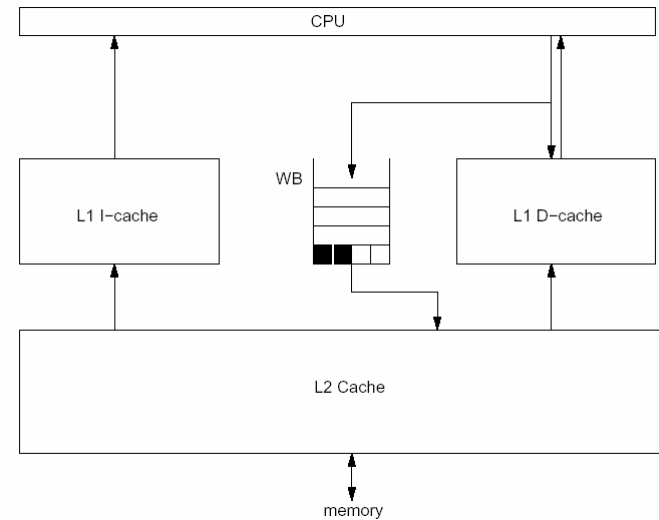
Write through

- All writes go into L1 cache and then also write through into subsequent levels of hierarchy
- Better for "cache coherence" issues
- No dirty/clean bit records required
- Faster evictions

Write Around?

26

Write Buffer



Source: Skadron/Clark

27

Cache Summary

- Two types of locality: spatial and temporal
- Spatial locality: larger block sizes
- Cache contents include data, tags, and valid bits

- Miss penalty is increasing (processor vs. memory)
- Modern processors use set-associative caches worth the cost

- Multi-level caches used to reduce miss penalty

- Variations: Victim Caches, Trace Caches

28