# Lecture 20: More caches

Monday, February 25, 2019    10:59 AM

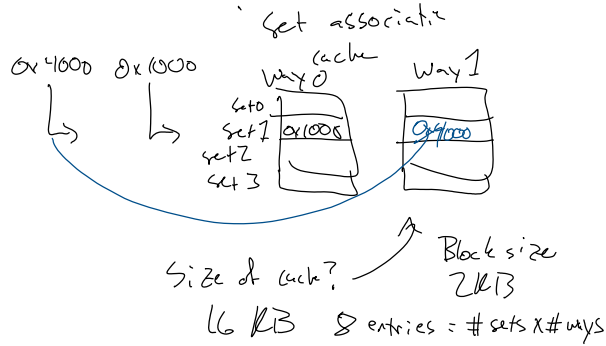Set associativity → multiple ways for data to be placed

Direct-mapped cache → each address only maps to a single location



Cache

0x1000

0x4000
miss
↳ conflict miss

A miss that wouldn't happen if we had more associativity

0x1000

0x4000

0x4000  0x1000

Set associative cache

way 0          Way 1

set 0
set 1   0x1000        0x4000
set 2
set 3

Size of cache?          Block size
                          2 KB
16 KB   8 entries = # sets × # ways

What is better?

Performance → reduce miss ratio → More associative

Fully associative   # of ways = lines in cache → any address can map to any location
  ↳ minimize conflict misses

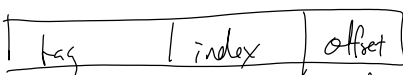Downside of associativity is power and area

Direct mapped
index ─────→

- 1 block + 1 tag
- compare 1 tag

Set associative
index ──→ way 0    way 1    way 2

- # ways (3) blocks + # ways tags
- compare many tags
- usually in parallel

Want high associativity when hit rate is very important → if memory is far away or high miss latency

address                              accessing SA cache

| tag | index | offset |

# bits                    # of bits → size of
whatever is left          cache line
32 bit → 15   direct-mapped    16 B cache line → 4 bits
address bits   # of bits
              log of
              # lines → 13 bits

# of
8192 lines → 15 n.

## Changes for SA?

now index bits → log(# of lines in a way or # of sets)
↳ log(# of unique addressable locations in cache)

# of index bits direct-mapped ⇒ log(# of lines)
One way                # sets = $\frac{\text{# lines}}{1}$

# of index bits for fully associative cache ⇒ $\emptyset$ index bits

### Example
64 kB cache

| 43 | 11 | 4 |
|---|---|---|
| tag | index | Offset |

**Equations**

Blocks/lines = sets × ways

Sets = $2^{\text{index bits}}$

Size = block size × sets × ways

Offset bits = log (block size)

**What is the associativity?**

# sets = 2048 ($2^{11}$)      # of lines = $\frac{64kB}{16B}$ = 4096

line size = 16 ($2^4$)

# ways = $\frac{\text{# lines}}{\text{#sets}}$ = $\frac{4096}{2048}$ = 2 → 2-way set associative

lines in cache = sets × ways

**RISC-V 32 bit**

32 bit addresses

$2^{32}$ → 4GB

most one app can address

---

## Kinds of misses

✗ **Conflict misses** result from low associativity
↳ reduce by incr. set associativity

✗ **Cold misses / Compulsory misses** → first access to an address
↳ reduce by prefetching
↳ generally hard to reduce these misses

✗ **Capacity misses** → an address was evicted before it was reused
↳ reduce these by making cache bigger
↳ smarter replacement policy

---

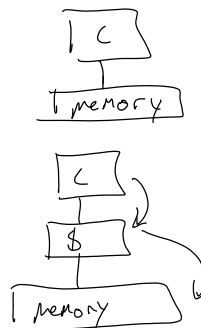## Predict Performance for caches
↳ Average memory access time (AMAT)

Memory takes 50 cycles

AMAT = memory time = 50 cycles

Cache → 2 cycle hit time + 95% hit ratio

AMAT = 0.95×2 + $\underbrace{0.05×2 + 0.05×50}_{\text{miss}}$

→ 2 + 0.05 50

$- 2 + 0.05 \cdot 50$

$= 4.5 \text{ cycles}$

AMAT = <u>hit time for cache</u> + miss ratio × (miss time)

10 cycle hit time and 80% hit ratio for L2

AMAT = $2 + 0.05[10 + 0.2 \times 50]$

$= 3$