

CS2100 - L14 - Caches (Set Associative & Fully Associative)

Week
8+9

14.1 - Cache performance

14.2 - Types of cache misses

14.3 - Set associative cache

14.4 - Fully associative cache

14.5 - Block replacement policy

14.6 - Cache framework

14.1 - Cache performance

Memory Access Time: Formula

Average Access Time

$$= \text{Hit rate} \times \text{Hit Time} + (1 - \text{Hit rate}) \times \text{Miss penalty}$$

■ **Hit:** Data is in cache (e.g., X)

- **Hit rate:** Fraction of memory accesses that hit
- **Hit time:** Time to access cache

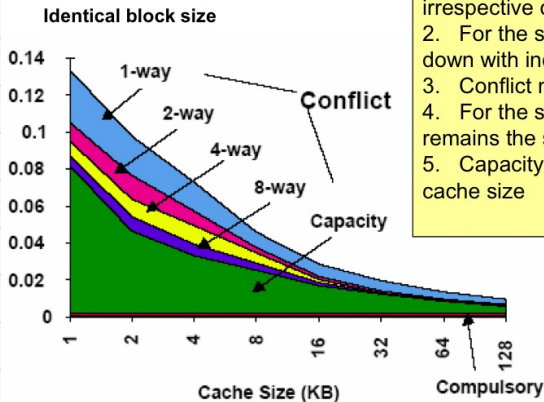
■ **Miss:** Data is not in cache (e.g., Y)

- **Miss rate** = $1 - \text{Hit rate}$
- **Miss penalty:** Time to replace cache block + deliver data

■ **Hit time < Miss penalty**

[L14 - AY2021S1]

Cache Performance



Observations:

1. Cold/compulsory miss remains the same irrespective of cache size/associativity
2. For the same cache size, conflict miss goes down with increasing associativity
3. Conflict miss is 0 for FA caches
4. For the same cache size, capacity miss remains the same irrespective of associativity
5. Capacity miss decreases with increasing cache size

Total Miss = Cold miss + Conflict miss + Capacity miss
Capacity miss (FA) = Total miss (FA) - Cold miss (FA), when Conflict Miss $\rightarrow 0$

[L14 - AY2021S1]

well defined for FA

Cache Misses: Classifications

Compulsory / Cold Miss

- First time a **memory block** is accessed
- Cold fact of life: Not much can be done
- **Solution**: Increase cache block size

Conflict Miss

- Two or more distinct memory blocks map to the same cache block
- Big problem in direct-mapped caches
- **Solution 1**: Increase cache size
 - Inherent restriction on cache size due to SRAM technology
- **Solution 2**: **Set-Associative caches** (coming next ..)

Capacity Miss

- Due to limited cache size
- Will be further clarified in "fully associative caches" later

Set Associative (SA) Cache

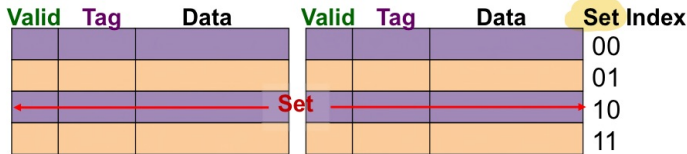
■ N-way Set Associative Cache

- A memory block can be placed in a fixed number of locations ($N > 1$) in the cache

■ Key Idea:

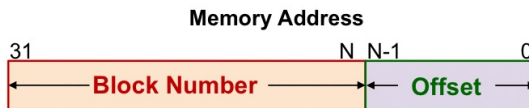
- Cache consists of a number of sets:
 - Each set contains N cache blocks
- Each memory block maps to a unique cache set
- Within the set, a memory block can be placed in **any** element of the set

SA Cache Structure



2-way Set Associative Cache

Set-Associative Cache: Mapping



Cache Block size = 2^N bytes

Cache Set Index
= (BlockNumber) modulo (NumberOfCacheSets)



Cache Block size = 2^N bytes

Number of cache sets = 2^M

Offset = N bits

Set Index = M bits

Tag = $32 - (N + M)$ bits

Observation:
It is essentially unchanged from the direct-mapping formula

[L14 - AY2021S1]

Advantage of Associativity (3/3)

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 $N / 2$

Fully Associative (FA) Cache

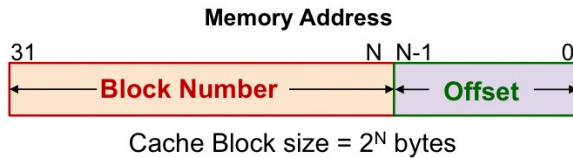
■ Fully Associative Cache

- A memory block can be placed in any location in the cache

■ Key Idea:

- Memory block placement is no longer restricted by cache index / cache set index
- ++ Can be placed in any location, **BUT**
- Need to search all cache blocks for memory access
→ lots of hardware needed

Fully Associative Cache: Mapping




Cache Block size = 2^N bytes
Number of cache blocks = 2^M
Offset = **N bits**
Tag = **32 - N bits**

Observation:
The block number
serves as the tag in FA
cache

14.5 - Block replacement policy

Block Replacement Policy (3/3)

- Drawback for LRU ^{least Recently Used}  ^{Difficult to shuffle blocks within set}
 - ❑ Hard to keep track if there are many choices
- Other replacement policies:
 - ❑ First in first out (FIFO)
 - Second chance variant ^(gauges whether block has been accessed after being placed)
 - ❑ Random replacement (RR)
 - ❑ Least frequently used (LFU)

Cache Organizations: Summary

One-way set associative (direct mapped)

Block	Tag	Data
0		
1		
2		
3		
4		
5		
6		
7		

$N=1$

Two-way set associative

Set	Tag	Data	Tag	Data
0				
1				
2				
3				

$N=2$

Four-way set associative

$N=4$

Set	Tag	Data	Tag	Data	Tag	Data	Tag	Data
0								
1								

Eight-way set associative (fully associative)

Tag	Data	Tag	Data	Tag	Data	Tag	Data	Tag	Data	Tag	Data	Tag	Data	Tag	Data

14.6 - Cache Framework

Cache Framework (1/2)

Block Placement: Where can a block be placed in cache?

Direct Mapped:

- Only one block defined by index

N-way Set-Associative:

- Any one of the **N** blocks within the set defined by index

Fully Associative:

- Any cache block

Block Identification: How is a block found if it is in the cache?

Direct Mapped:

- Tag match with only one block

N-way Set Associative:

- Tag match for all the blocks within the set

Fully Associative:

- Tag match for all the blocks within the cache

[L14 - AY2021S1]

Cache Framework (2/2)

Block Replacement: Which block should be replaced on a cache miss?

Direct Mapped:

- No Choice

n-way Set-Associative:

- Based on replacement policy

Fully Associative:

- Based on replacement policy

Write Strategy: What happens on a write?

Write Policy: Write-through vs. write-back

Write Miss Policy: Write allocate vs. write no allocate



L13 - DM Cache

[L14 - AY2021S1]