

Cache II

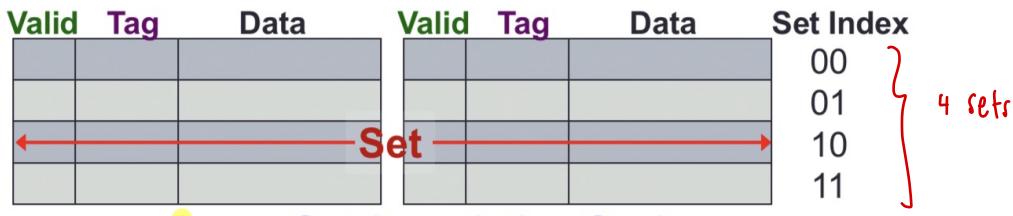
Block Size Trade-off

- Average Access Time = Hit rate \times Hit Time + (1 - Hit rate) \times Miss penalty.
- Larger block size:
 - + takes advantage of spatial locality
 - larger miss penalty: takes longer time to fill up the block
 - If block size is too big relative to cache size
 - Too few cache blocks → miss rate will go up.

Set Associative (SA) Cache

- Conflict misses (several blocks mapped to same block / set)
 - ↳ Solution via SA cache.
- N-way Set Associative Cache.
 - ↖ number of blocks in one set.
- 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 of the N cache blocks in the set.

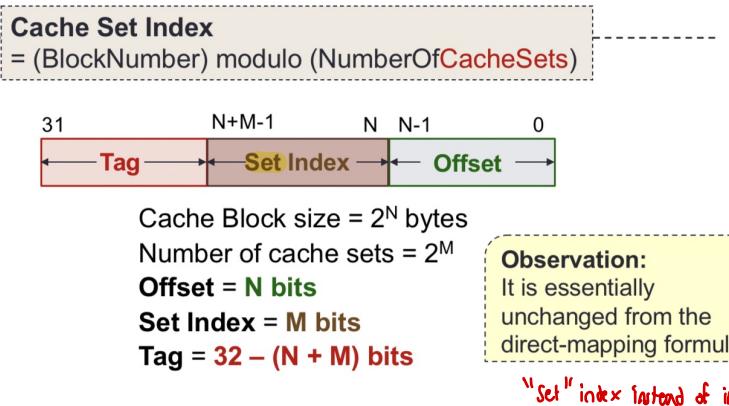
Structure of SA Cache



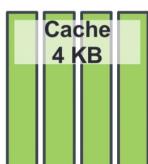
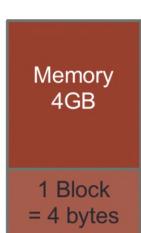
2-way Set Associative Cache

- each set has 2 cache blocks
- A memory block maps to a unique set.
 - In the set, the memory block can be placed in either of the cache blocks.
 - └ need to search both to look for the memory block.

Mapping:

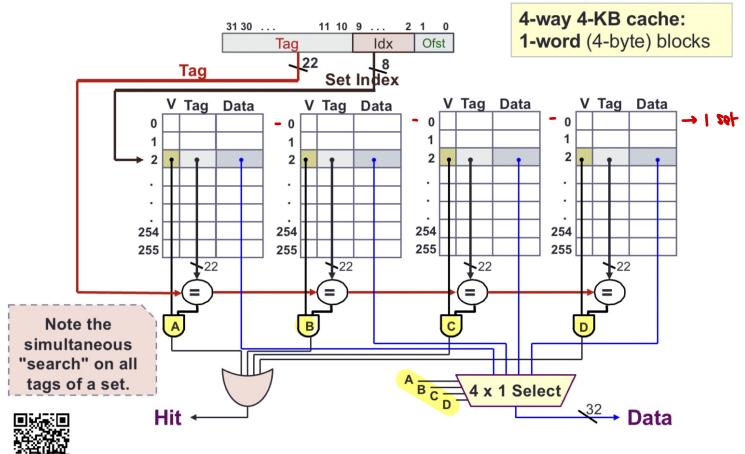


Example:

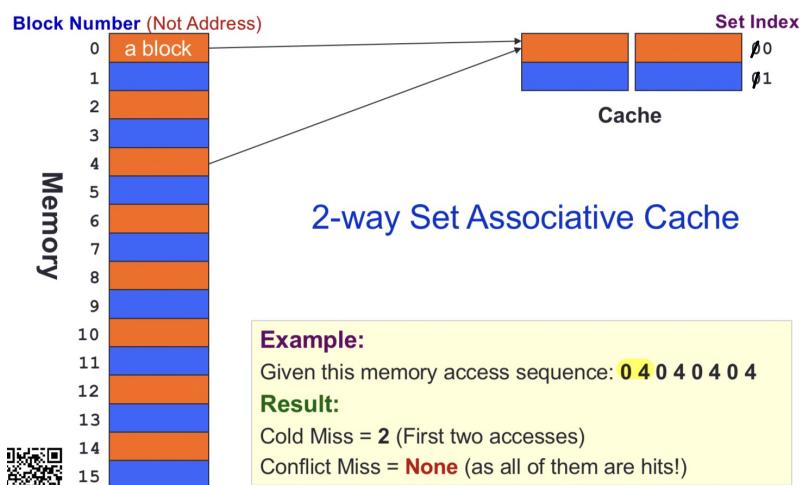
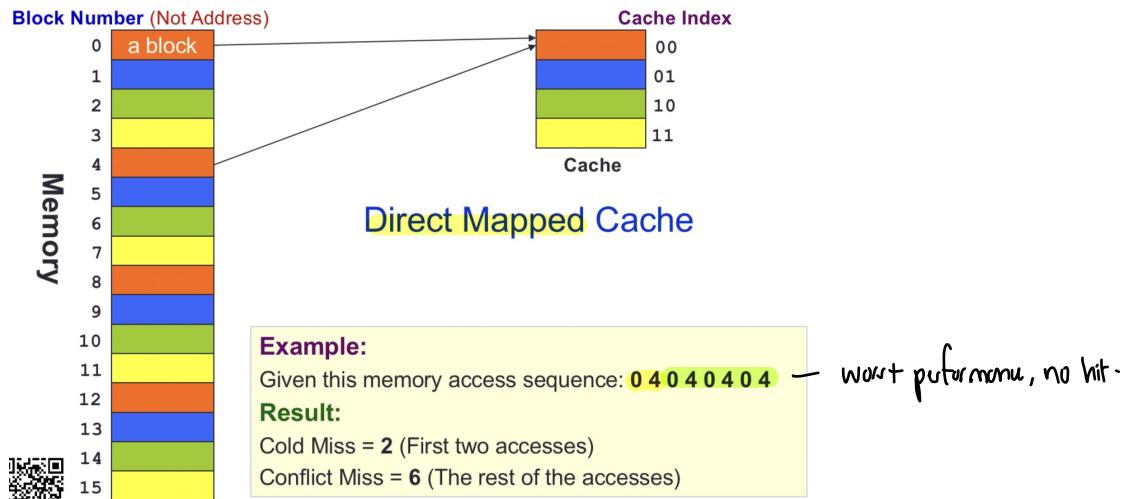


4-way associative, number of sets = $1024 / 4 = 256 = 2^8$
Set Index, M = 8 bits
Cache Tag = $32 - 8 - 2 = 22$ bits

Every set has 4 blocks.

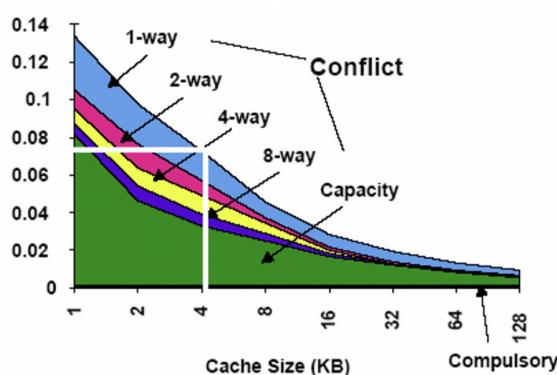


Advantage of Associativity.



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**



Given:

- Memory access sequence: 4, 0, 8, 36, 0
- 2-way set-associative cache with a total of four 8-byte blocks → **total of 2 sets**
- Indicate hit/miss for each access



Offset, N = 3 bits

Block Number = $32 - 3 = 29$ bits

2-way associative, number of sets = 2^1

Set Index, M = 1 bits

Cache Tag = $32 - 3 - 1 = 28$ bits

3. SA Cache Example: Load #1

Miss

Tag Index Offset

- Load from 4 → 0 | 100

Check: Both blocks in Set 0 are invalid [**Cold Miss**]

Result: Load from memory and place in Set 0 - Block 0

Set Index	Block 0				Block 1			
	Valid	Tag	W0	W1	Valid	Tag	W0	W1
0	✓ 1	0	M[0]	M[4]	0			
1	0				0			

3. SA Cache Example: Load #2

Miss Hit

Tag Index Offset

- Load from 0 → 0 | 000

Result:

[Valid and Tags match] in Set 0-Block 0 [**Spatial Locality**]

Set Index	Block 0				Block 1			
	Valid	Tag	W0	W1	Valid	Tag	W0	W1
0	1	0	M[0]	M[4]	0			

3. SA Cache Example: Load #3

Miss Hit Miss

Tag Index Offset

- Load from 8 → 1 | 000

Check: Both blocks in Set 1 are invalid [**Cold Miss**]

Result: Load from memory and place in Set 1 - Block 0

Set Index	Block 0				Block 1			
	Valid	Tag	W0	W1	Valid	Tag	W0	W1
0	1	0	M[0]	M[4]	0			
1	✓ 1	0	M[8]	M[12]	0			

3. SA Cache Example: Load #4

Miss Hit Miss Miss

Tag Index Offset

- Load from 36 → 0 | 100

Check: [Valid but tag mismatch] Set 0 - Block 0

[Invalid] Set 0 - Block1 [**Cold Miss**]

Result: Load from memory and place in Set 0 - Block 1

Set Index	Block 0				Block 1			
	Valid	Tag	W0	W1	Valid	Tag	W0	W1
0	1	0	M[0]	M[4]	✓ 1	2	M[32]	M[36]
1	1	0	M[8]	M[12]	0			

3. SA Cache Example: Load #5

Miss Hit Miss Miss Hit

Tag Index Offset

- Load from 0 → 0 | 000

Check: [Valid and tags match] Set 0-Block 0

[Valid but tags mismatch] Set 0-Block1

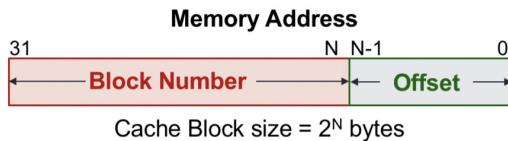
[**Temporal Locality**]

Set Index	Block 0				Block 1			
	Valid	Tag	W0	W1	Valid	Tag	W0	W1
0	1	0	M[0]	M[4]	1	2	M[32]	M[36]
1	1	0	M[8]	M[12]	0			

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 or cache set index.
 - Can be placed in **any location**
 - But need to **search all cache blocks** for memory access.
 - Inurs capacity misses.

Mapping:



Cache Block size = 2^N bytes

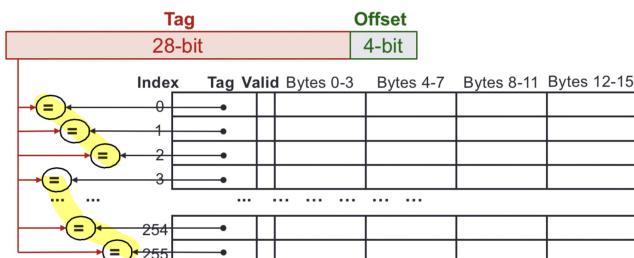
Number of cache blocks = 2^M

Offset = **N bits**

Tag = **32 - N bits**

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

- Example: $2^{12} \times 2^4$
 - 4KB cache size and 16-Byte block size $\Rightarrow 2^8 = 256$ blocks
 - Compare tags and valid bit in parallel



No Conflict Miss (since data can go anywhere)

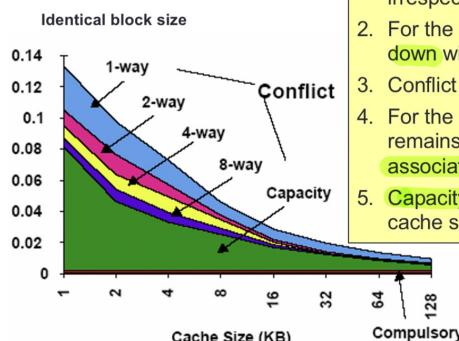
have to compare every Valid bit (no index)

a lot

of comparators!

4. Cache Performance

\Rightarrow Useful for small cache.



Observations:

- Cold/compulsory miss remains the same irrespective of cache size/associativity.
- For the same cache size, conflict miss goes down with increasing associativity.
- Conflict miss is 0 for FA caches.
- For the same cache size, capacity miss remains the same irrespective of associativity.
- Capacity miss decreases with increasing cache size. (FA)

Total Miss = Cold miss + Conflict miss + Capacity miss

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

5. Block Replacement Policy (1/3)

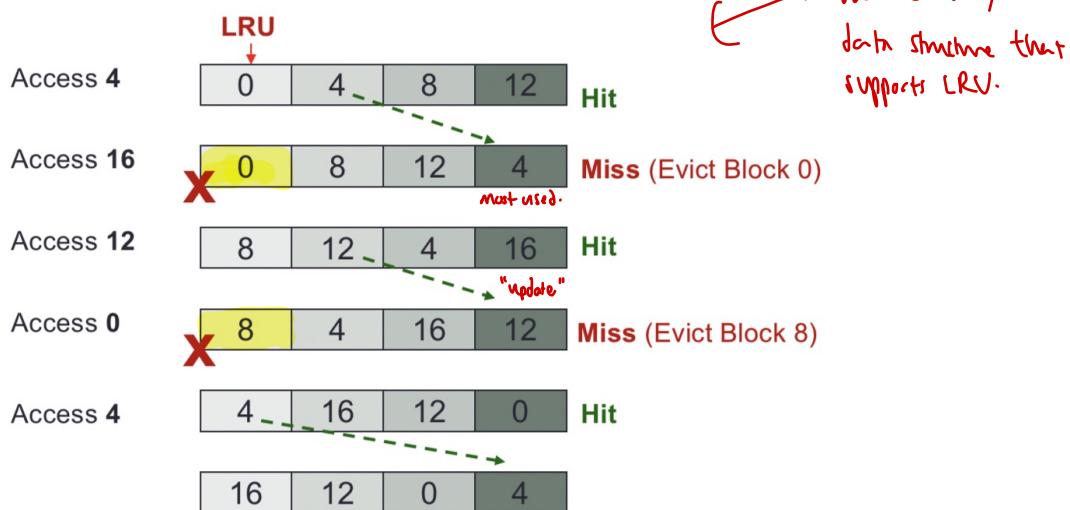
- Set Associative or Fully Associative Cache: (not Direct Mapped Cache)
 - Can choose where to place a memory block
 - Potentially replacing another cache block if full
 - Need **block replacement policy**

▪ Least Recently Used (LRU)

- **How:** For cache hit, record the cache block that was accessed
 - When replacing a block, choose one which has not been accessed for the longest time
- **Why:** Temporal locality – Least used less likely to be used again.

Least Recently Used policy in action:

- 4-way SA cache
- Memory accesses: 0 4 8 12 4 16 12 0 4



- Drawback for LRU
 - Hard to keep track if there are many choices
- Other replacement policies:
 - First in first out (FIFO)
 - Random replacement (RR)
 - Least frequently used (LFU)

Example (for all 3 types):

6. Additional Examples #1

- **Direct-Mapped Cache:**
 - Four 8-byte blocks
- Memory accesses:
4, 8, 36, 48, 68, 0, 32

Addr:	Tag	Index	Offset
4:	00...0000	00	100
8:	00...0000	01	000
36:	00...0001	00	100
48:	00...0001	10	000
68:	00...0010	00	100
0:	00...0000	00	000
32:	00...0001	00	000

Index	Valid	Tag	Word0	Word1
0	0 ¹	0 ¹ 1 ² 0	M[0] M[32] M[64]	M[36] M[36] M[48]
1	0 ¹	0	M[8]	M[12]
2	0 ¹	1	M[48]	M[52]
3	0			



6. Additional Examples #2

- **Fully-Associative Cache:**
 - Four 8-byte blocks
 - LRU Replacement Policy
- Memory accesses:
4, 8, 36, 48, 68, 0, 32

Addr:	Tag	Offset
4:	00...000000	100
8:	00...000001	000
36:	00...001000	100
48:	00...001100	000
68:	00...010000	100
0:	00...000000	000
32:	00...001000	000

Index	Valid	Tag	Word0	Word1
0	0 ¹	0 ⁸	M[0] M[64]	M[4] M[68]
1	0 ¹	1 ⁰	M[8] M[0]	M[12] M[4]
2	0 ¹	4	M[32]	M[36]
3	0 ¹	6	M[48]	M[52]



replace at current cache for next iteration.

6. Additional Examples #3

- **2-way Set-Associative Cache:**
 - Four 8-byte blocks
 - LRU Replacement Policy
- Memory accesses:
4, 8, 36, 48, 68, 0, 32

Addr:	Tag	Index	Offset
4:	00...000000	0	100
8:	00...000001	1	000
36:	00...001000	0	100
48:	00...001100	0	000
68:	00...010000	0	100
0:	00...000000	0	000
32:	00...001000	0	000

Set Index	Block 0				Block 1			
	Valid	Tag	Word0	Word1	Valid	Tag	Word0	Word1
0	0 ¹	0 3 0	M[0] M[48] M[0]	M[4] M[52] M[4]	0 ¹	2 4 2	M[32] M[64] M[32]	M[36] M[68] M[36]
1	0 ¹	0	M[8]	M[12]	0			

7. Summary: Cache Organizations

One-way set associative
(direct mapped)

Block	Tag	Data
0		
1		
2		
3		
4		
5		
6		
7		

Two-way set associative

Set	Tag	Data	Tag	Data
0				
1				
2				
3				

Four-way set associative

Set	Tag	Data	Tag	Data	Tag	Data	Tag	Data
0								
1								

Eight-way set associative (fully associative)

Tag	Data														

offset
= What word
in the block.

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

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