



# Lecture 11: Cache II

## **CS10014 Computer Organization**

Department of Computer Science

Tsung Tai Yeh

Thursday: 1:20 pm– 3:10 pm

Classroom: EC-022



# Acknowledgements and Disclaimer

- Slides were developed in the reference with
  - CS 61C at UC Berkeley
    - <https://inst.eecs.berkeley.edu/~cs61c/sp23/>
  - CS252 at ETHZ
    - <https://safari.ethz.ch/digitaltechnik/spring2023>
  - CIS510 at Upenn
    - <https://www.cis.upenn.edu/~cis5710/spring2019/>



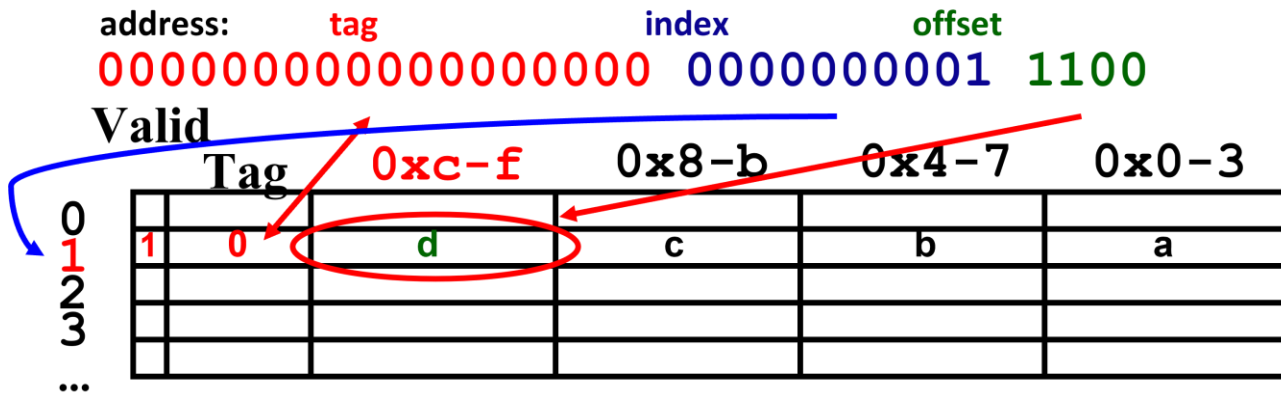
# Outline

- Associativity
- Fully Associative Cache
- N-way Set Associative Cache
- Cache Write Issue
- Block Replacement Policy
- Multi-level Caches



# Review: Cache Basics (1/2)

- How to transparently move data among levels of a storage hierarchy
  - Address => index to set of candidates
  - Compare desired address with tag
  - Service hit or miss -> load new block and binding on miss





## Review: Cache Basics (2/2)

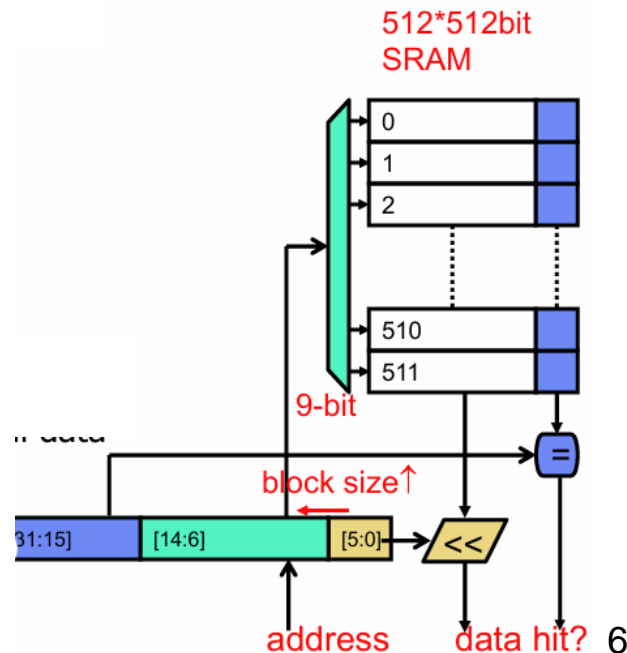
- 4-bit addresses -> 16 bytes memory
- 8 bytes cache, 2 bytes blocks
  - The number of sets: 4 (capacity / block size)
  - How address splits into offset/index/tag bits
    - **Offset:** least-significant  $\log_2(\text{block size}) = \log_2(2) = 1$  -> 000**0**
    - **Index:** next  $\log_2(\text{number-of-sets}) = \log_2(4) = 2$  -> 00**00**
    - **Tag:** rest =  $4 - 1 - 2 = 1$  -> **0**000





# Block Size Tradeoff (1/7)

- Given capacity, manipulate miss rate by changing cache organization
- One option: increase **block size**
  - Exploit spatial locality
  - Notice index/offset bits change
  - Tag remain the same
- Increasing cache block size
  - + reduce miss rate (up to a point)
  - + reduce tag overhead (why?)
  - - potentially useless data transfer





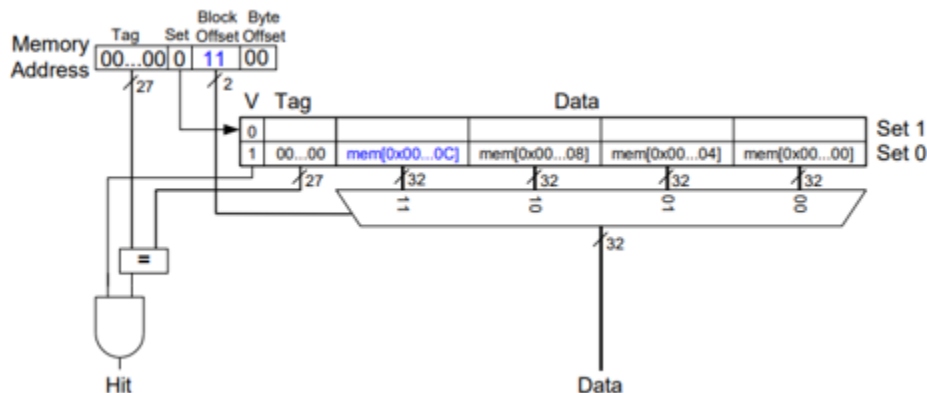
## Block Size Tradeoff (2/7)

- 4KB cache with 1024 4B blocks?
  - 4B blocks -> 2-bit offset, 1024 frames -> 10-bit index
  - 32-bit address - 2-bit offset – 10-bit index = 20-bit tag
  - 20-bit tag / 32-bit block = 63% overhead
- 4KB cache with 512 8B blocks?
  - 8B blocks -> 3-bit offset, 512 frames -> 9-bit index
  - 32-bit address – 3-bit offset – 9-bit index = 20-bit tag
  - 20-bit tag / 64-bit block = 32% overhead
- A realistic example: 64KB cache with 64B blocks
  - 16-bit tag / 512-bit block = ~2% overhead



# Block Size Tradeoff (3/7)

- Benefits of larger block size
  - **Spatial locality**
    - If we access a given word, we're likely to access other nearby words soon
    - Works nicely in sequential array accesses too







## Block Size Tradeoff (4/7)

- **Drawbacks of larger block size**
  - **Larger block size** means **larger miss penalty**
    - On a miss, takes longer time to load a new block from the next level memory
  - If block size is too big relative to cache size, then there are too few blocks
    - Result: miss rate goes up
  - In general, minimize **Average Memory Access Time (AMAT)**
    - = Hit time + Miss Penalty x Miss Rate



## Block Size Tradeoff (5/7)

- **Hit Time**

- Time to find and retrieve data from current level cache

- **Miss Penalty**

- Average time to retrieve data on a current level miss (includes the possibility of misses on successive levels of memory hierarchy)

- **Hit Rate**

- % of requests that are found in current level cache

- **Miss Rate**

- $1 - \text{Hit Rate}$



## Block Size Tradeoff (6/7)

**Valid Bit**



**Tag**



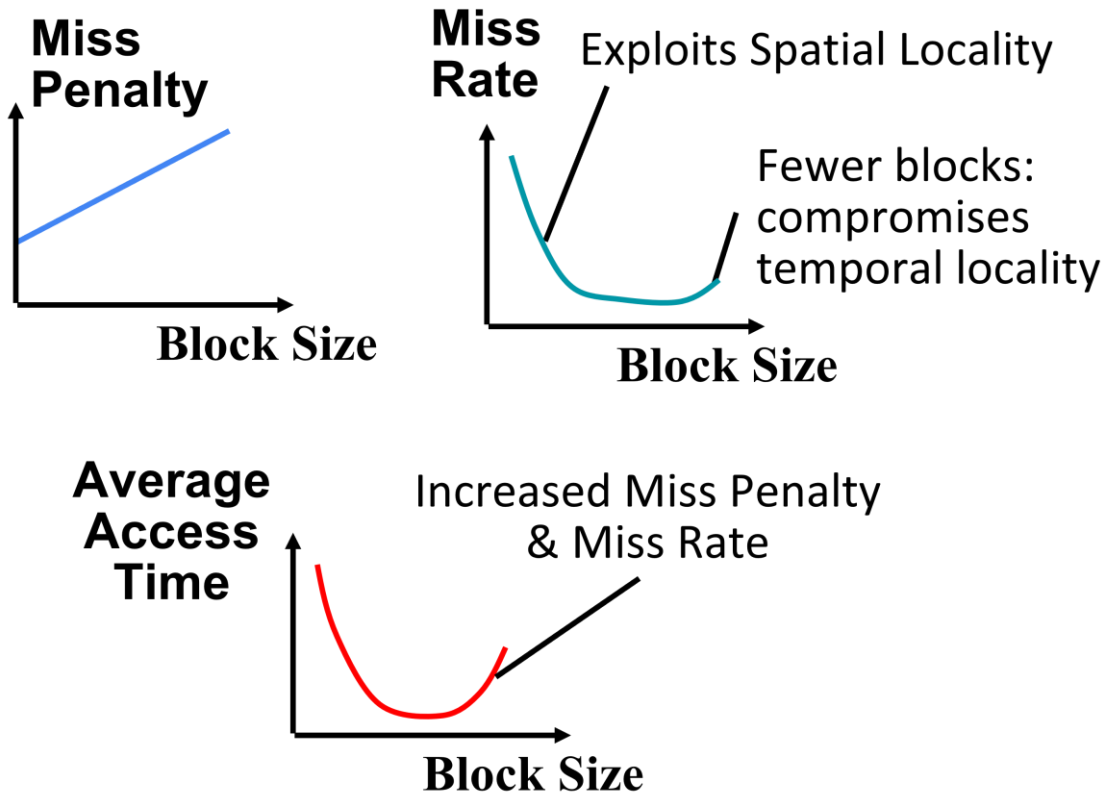
**Cache Data**



- Cache Size = 4 bytes, Block Size = 4 bytes
  - Only ONE entry (row) in the cache
- If item accessed, likely accessed again soon
  - But unlikely will be accessed again immediately!
- The next access will likely to be a miss again
  - Continually loading data into the cache but discard data before using it again
  - Nightmare for cache designer: **Ping Pong Effect**



# Block Size Tradeoff (7/7)





# Takeaway Question

- Parameters
  - Simple pipeline with base CPI of 1
  - Instruction mix: 30% loads/stores
  - I\$: %miss = 2%,  $t_{\text{miss}} = 10$  cycles
  - D\$: %miss = 10%,  $t_{\text{miss}} = 10$  cycles
- What is the new CPI?



# Takeaway Question

- Parameters

- Simple pipeline with base CPI of 1
- Instruction mix: 30% loads/stores
- I\$: %miss = 2%,  $t_{\text{miss}} = 10$  cycles
- D\$: %miss = 10%,  $t_{\text{miss}} = 10$  cycles

- What is the new CPI?

- $\text{CPI}_{\text{I\$}} = \% \text{miss I\$} * t_{\text{miss}} = 0.02 * 10 \text{ cycles} = 0.2 \text{ cycle}$
- $\text{CPI}_{\text{D\$}} = \% \text{ load/store} * \% \text{miss D\$} * t_{\text{miss}} = 0.3 * 0.02 * 10 \text{ cycles} = 0.3 \text{ cycle}$
- $\text{CPI}_{\text{new}} = \text{CPI} + \text{CPI}_{\text{I\$}} + \text{CPI}_{\text{D\$}} = 1 + 0.2 + 0.3 = 1.5$



# Types of Cache Misses (1/3)

- **“Three Cs” Model of Misses**
- **1<sup>st</sup> C: Compulsory misses**
  - Occur when a program is first started
  - Cache does not contain any of that program’s data yet, so misses are bound to occur
  - Can’t be avoided easily



## Types of Cache Misses (2/3)

- **“Three Cs” Model of Misses**

- **2<sup>nd</sup> C: Conflict Misses**

- Miss that occurs because two distinct memory addresses map to the same cache location
- Two blocks (which happen to map to the same location) can keep overwriting each other
- Big problem in **direct-mapped caches**

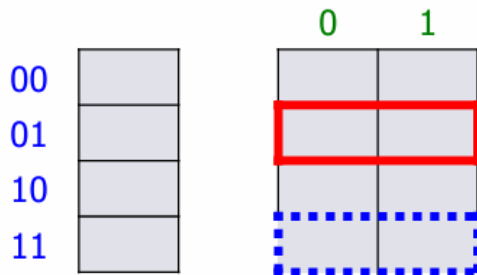
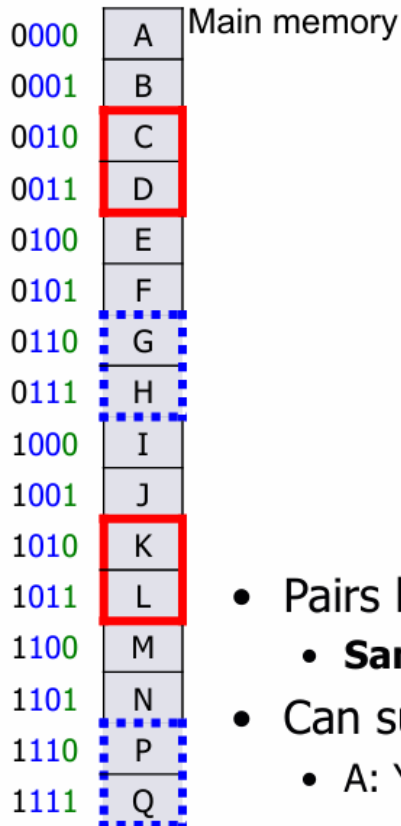
- Dealing with conflict misses

- Solution 1: Make the cache size bigger
- Solution 2: Multiple distinct blocks can fit in the same cache index<sub>16</sub>





# Types of Cache Misses (3/3)



- Pairs like “0010” and “1010” **conflict**
  - **Same index!**
- Can such pairs to simultaneously reside in cache?
  - A: Yes, if we reorganize cache to do so

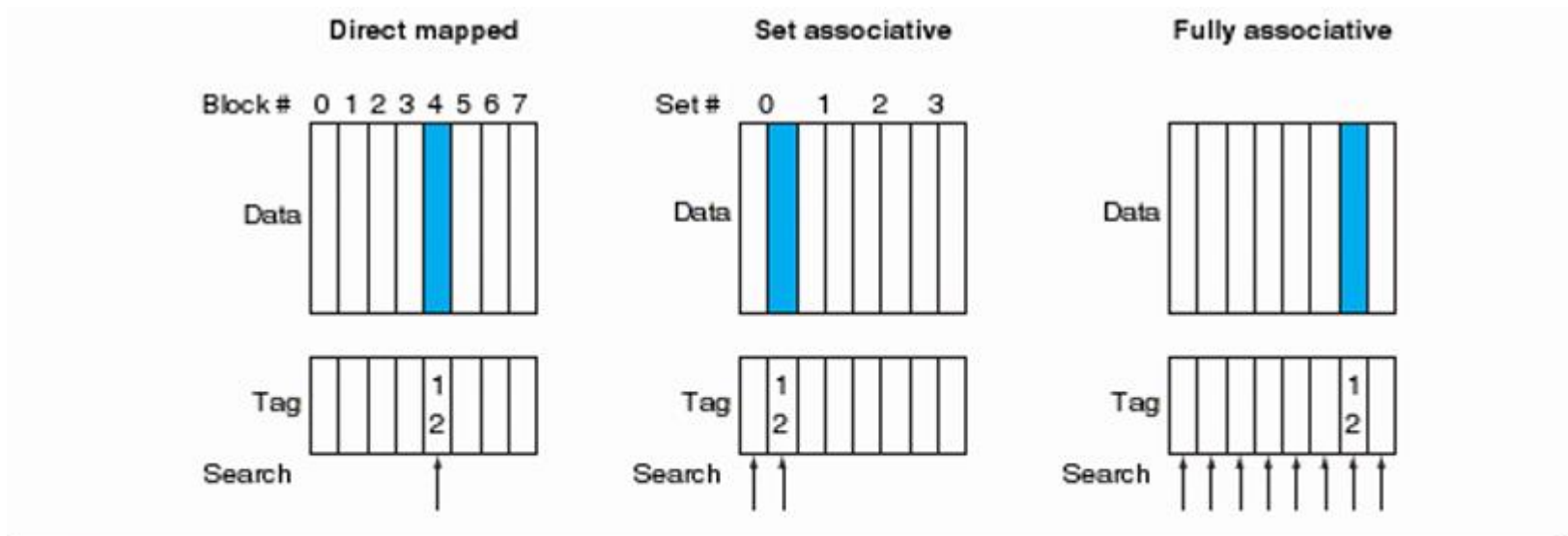


# Outline

- **Associativity**
- Fully Associative Cache
- N-way Set Associative Cache
- Cache Write Issue
- Block Replacement Policy
- Multi-level Caches



# Associativity (1/6)



**FIGURE 5.13** The location of a memory block whose address is 12 in a cache with eight blocks varies for direct-mapped, set-associative, and fully associative placement. In direct-mapped placement, there is only one cache block where memory block 12 can be found, and that block is given by  $(12 \bmod 8) = 4$ . In a two-way set-associative cache, there would be four sets, and memory block 12 must be in set  $(12 \bmod 4) = 0$ ; the memory block could be in either element of the set. In a fully associative placement, the memory block for block address 12 can appear in any of the eight cache blocks.



## Associativity (2/6)

- **Direct-mapped cache**

- Index completely specifies position which position a block can go in on a miss

- **N-Way Set associative**

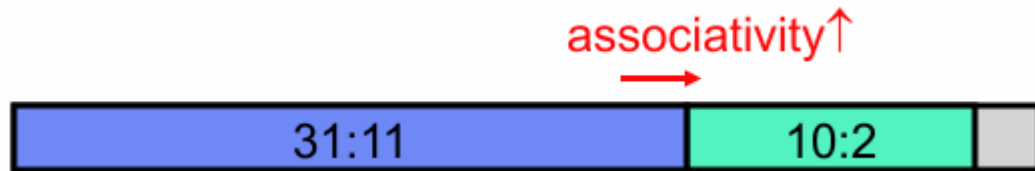
- Index specifies a set, but block can occupy any position within the set on a miss

- **Fully associative**

- Block can be written into any position

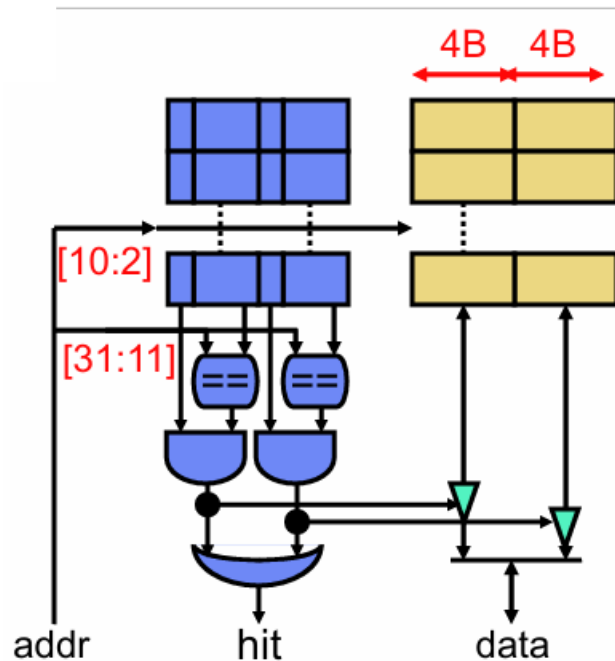


# Associativity (3/6)



- Set-associativity

- Block can reside in one of few frames
- Frame groups called sets
- Each frame in set called a way
- This is 2-way set-associative (SA)
- 1-way  $\rightarrow$  directed-mapped (DM)
- 1-set  $\rightarrow$  fully-associative (FA)
- + Reduce conflicts
- - Increase latency<sub>hit</sub>
  - Additional tag match & muxing



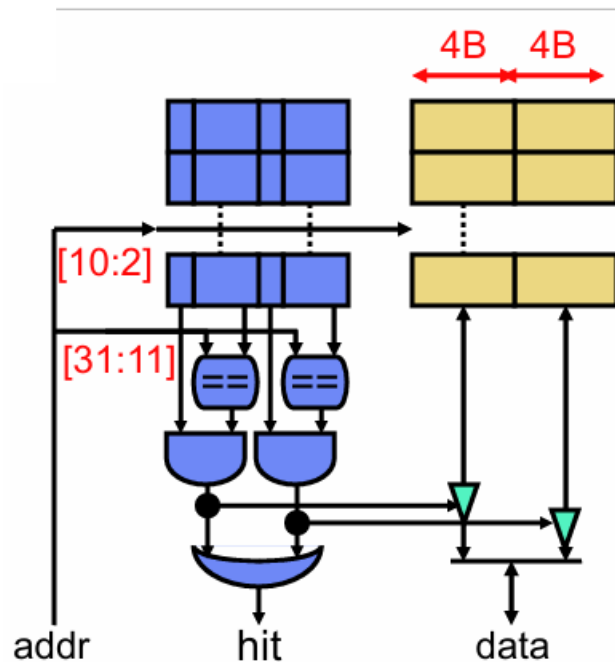


# Associativity (4/6)



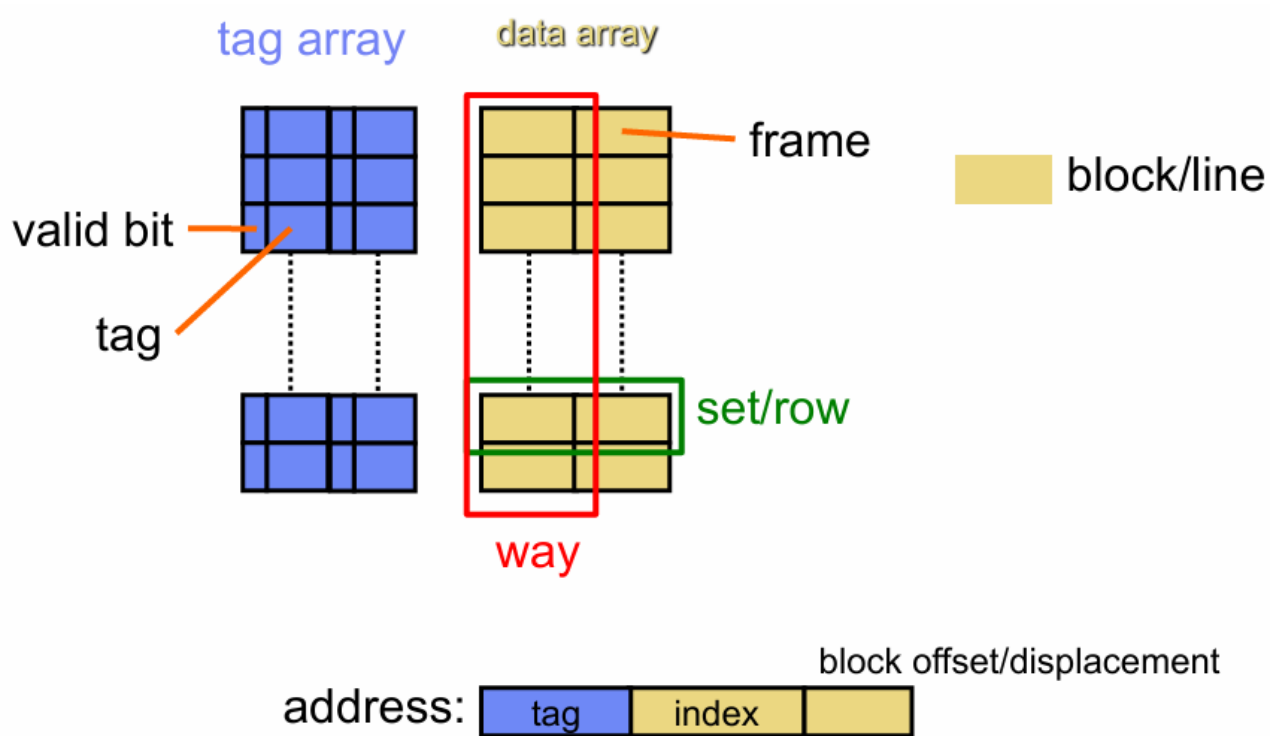
## • Lookup algorithm

- Use index bits to find set
- Read data/tags in all frames in parallel
- **Any** (match and valid bit), Hit
- Notice tag/index/offset bits
  - Only 9-bit index (versus 10-bit for direct mapped)





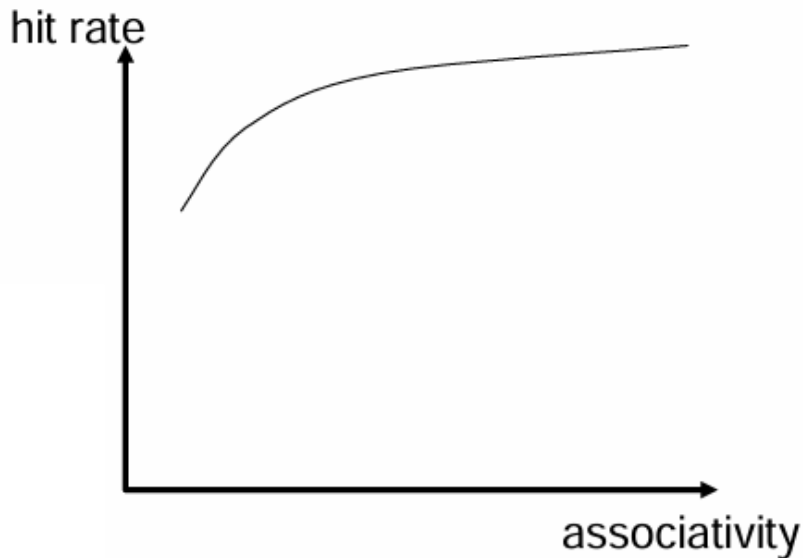
# Associativity (5/6)





## Associativity (6/6)

- How many blocks can be present in the same index (i.e., set)?
- Larger associativity
  - Lower miss rate (reduced conflict)
  - Higher hit latency and area cost
- Smaller associativity
  - Lower cost
  - Lower hit latency
  - Especially important for L1 caches







# Outline

- Associativity
- **Fully Associative Cache**
- N-way Set Associative Cache
- Cache Write Issue
- Block Replacement Policy
- Multi-level Caches



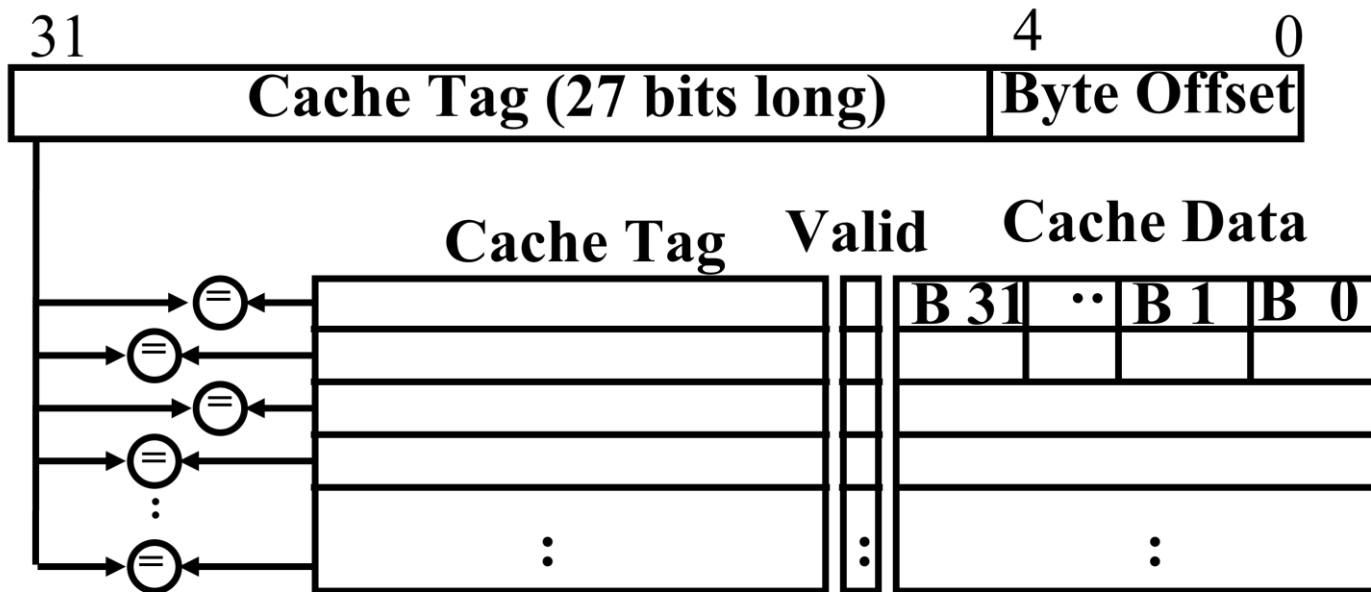
# Fully Associative Cache (1/5)

- Memory address fields
  - Tag: same as before
  - Offset: same as before
  - **Index: Non-exist**
- What does this mean?
  - No “rows”: any block can go anywhere in the cache
  - Must compare with all tags in entire cache to see if data is there



# Fully Associative Cache (2/5)

- Fully Associative Cache (e.g., 32 bytes block)
  - Compare tags in parallel

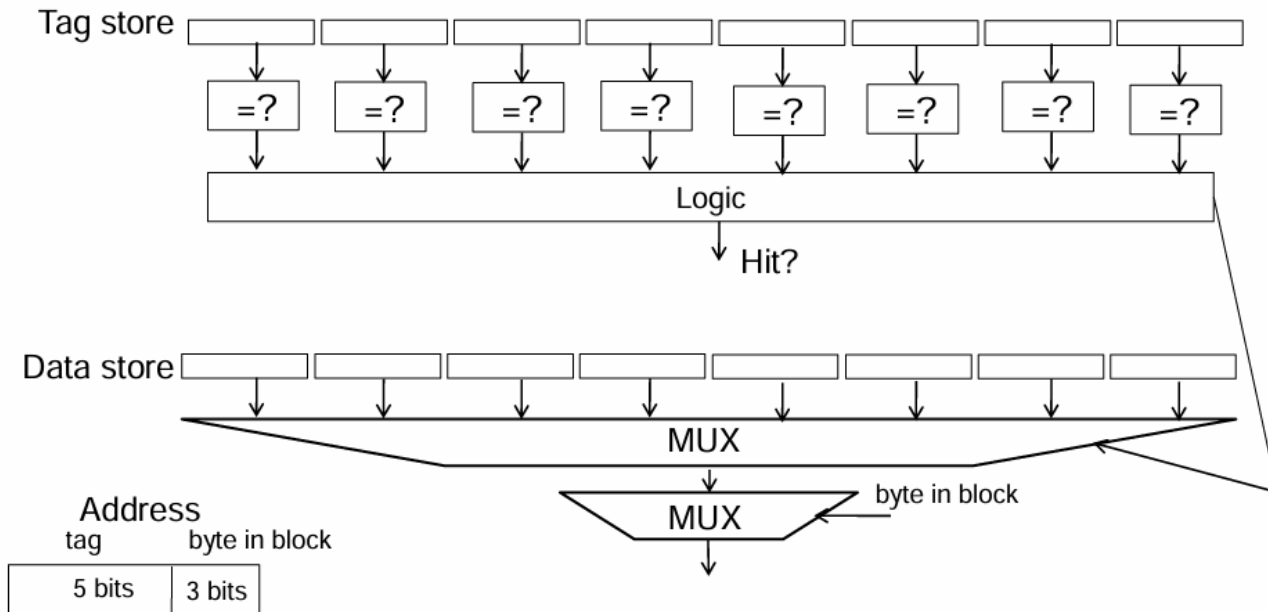




# Fully Associative Cache (3/5)

- Fully Associative Cache

- A block can be placed in any cache location





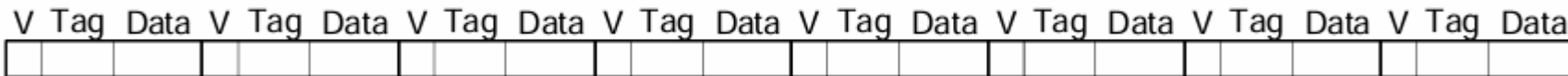
# Fully Associative Cache (4/5)

- **Benefit of Fully Assoc Cache**

- No conflict misses (since data can go anywhere)

- **Drawbacks of Fully Assoc Cache**

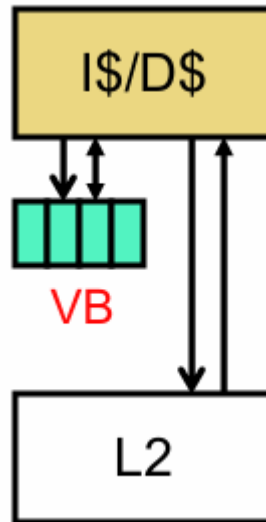
- Need hardware comparator for every single entry
- If we have a 64KB of data in cache with 4 bytes entries, we need 16K comparators
- Expensive to build





# Fully Associative Cache (5/5)

- **Victim buffer (VB):** small fully-associative cache
  - Reduce conflict miss
  - Sits on I\$/D\$ miss path
  - Small (e.g. 8 entries) so very fast
  - Blocks kicked out of I\$/D\$ placed in VB
  - On miss, check VB: hit? Place block back in I\$/D\$
  - 8 extra ways, shared among all sets
  - Very effective in practice





# Types of Cache Misses

- **“Three Cs” Model of Misses**
- **3<sup>rd</sup> C: Capacity Misses**
  - Miss that occurs because the cache has a limited size
  - Miss that would not occur if we increase the size of the cache
  - This is the primary type of miss for Fully Associative cache



# Outline

- Associativity
- Fully Associative Cache
- **N-way Set Associative Cache**
- Cache Write Issue
- Block Replacement Policy
- Multi-level Caches



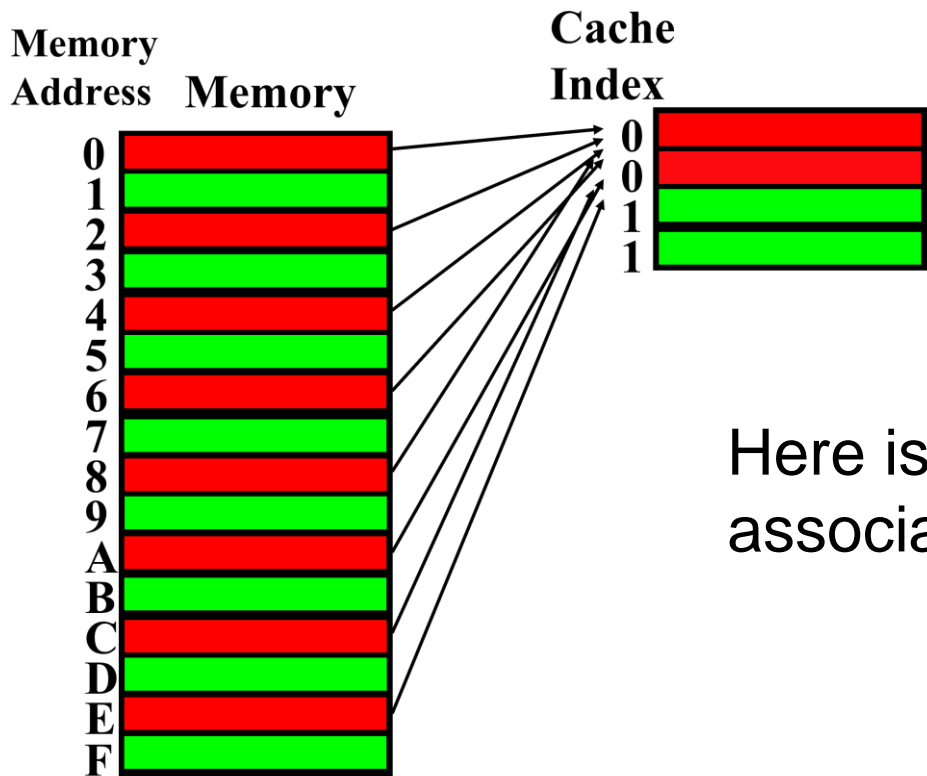


## N-way Set Associative Cache (1/8)

- Memory address fields:
  - **Tag**: same as before
  - **Offset**: same as before
  - **Index**: points us to the correct “row” (called a set in this case)
- What’s the difference?
  - Each set contains multiple blocks
  - Once we’ve found correct set, must compare with all tags in that set to find our data



# N-way Set Associative Cache (2/8)



Here is a simple 2-way set associative cache



# N-way Set Associative Cache (3/8)

- **Basic idea**

- Cache is directed-mapped w/respect to sets
- Each set is fully associative with N blocks in it

- **Given memory address**

- Find correct set using index value
- Compare Tag with all Tag values in the determined set
- If a match occurs, hit! Otherwise a miss
- Finally, use the offset field as usual to find the desired data within the block

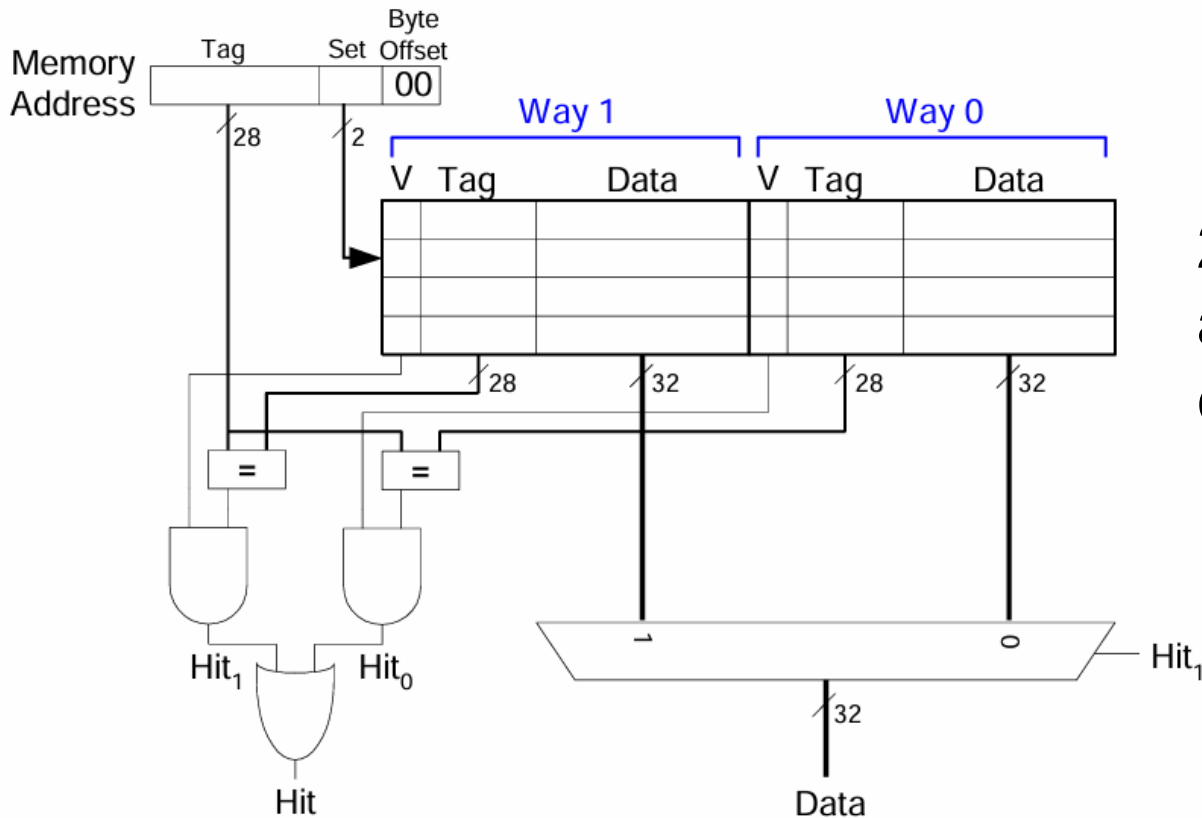


## N-way Set Associative Cache (4/8)

- What's so great about this?
  - Even a 2-way set associative cache avoids a lot of conflict misses
  - Hardware cost isn't that bad: only need  $N$  comparators
- In fact, for a cache with  $M$  blocks
  - It's **Direct-Mapped** if it's 1-way set associativity
  - It's **Fully Associative** if it's  $M$ -way set associativity
  - So these two are just special cases of the more general set associative design



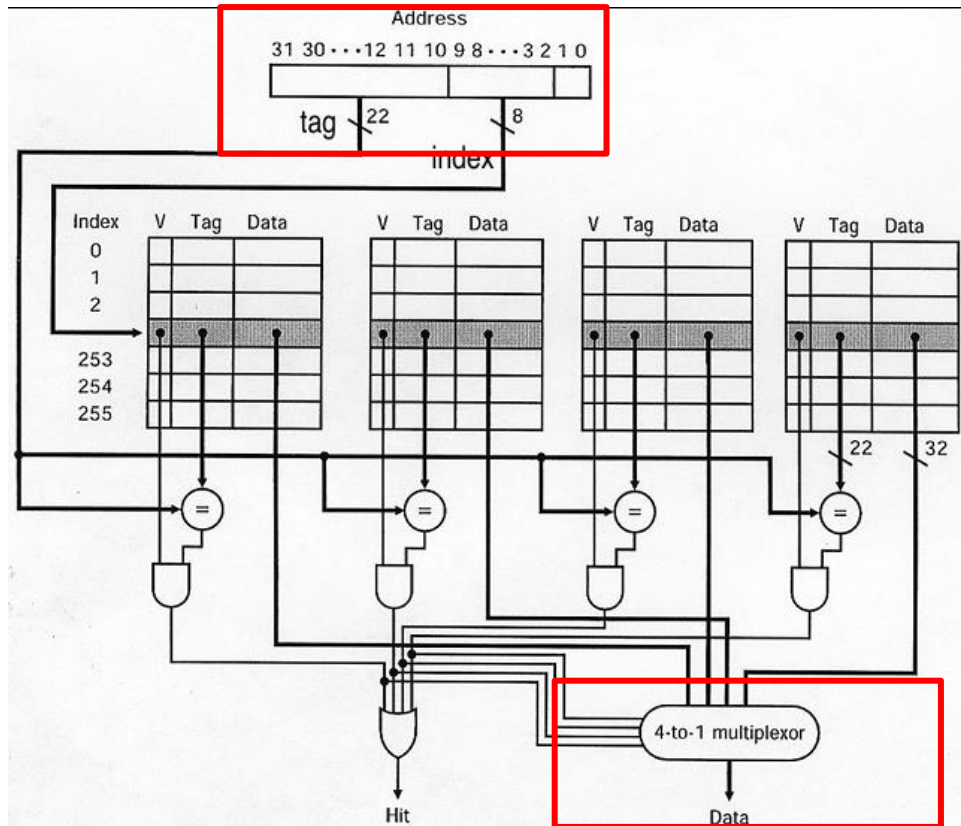
# N-way Set Associative Cache (5/8)



2-way set  
associative  
cache



# N-way Set Associative Cache (6/8)



4-way set  
associative  
cache

# ways = index  
length / offset  
length



# N-way Set Associative Performance (7/8)

```
# RISC-V assembly code
      addi $t0, $0, 5
loop:  beq  $t0, $0, done
      lw   $t1, 0x4($0)
      lw   $t2, 0x24($0)
      addi $t0, $t0, -1
      j    loop
done:
```

*Miss Rate =*

Way 1			Way 0			
V	Tag	Data	V	Tag	Data	
0			0			Set 3
0			0			Set 2
1	00...10	mem[0x00...24]	1	00...00	mem[0x00...04]	Set 1
0			0			Set 0



# N-way Set Associative Performance (8/8)

```
# RISC-V assembly code
```

```

    addi $t0, $0, 5
loop:  beq  $t0, $0, done
       lw   $t1, 0x4($0)
       lw   $t2, 0x24($0)
       addi $t0, $t0, -1
       j   loop
done:

```

*Miss Rate = 2/10*

*= 20%*

Associativity reduces  
conflict misses

Way 1			Way 0			
V	Tag	Data	V	Tag	Data	
0			0			Set 3
0			0			Set 2
1	00...10	mem[0x00...24]	1	00...00	mem[0x00...04]	Set 1
0			0			Set 0





# Outline

- Associativity
- Fully Associative Cache
- N-way Set Associative Cache
- **Cache Write Issue**
- Block Replacement Policy
- Multi-level Caches



## Cache Write Issue (1/9)

- So far we have looked at reading from cache
  - Instruction fetches, loads
- What about writing into cache
  - Stores, not an issue for instruction caches
- Several new issues
  - Tag/data access
  - Write-through vs. write-back
  - Write-allocate vs. write-not-allocate
  - Hiding write miss latency



## Cache Write Issue (2/9)

- Tag/Data access
  - Reads: read tag and data in parallel
    - Tag mis-match -> data is wrong (OK, just stall until good data arrives)
  - Writes: read tag, write data in parallel? No. Why?
    - Tag mis-match -> clobbered data (oops!)
    - For associative caches, which way was written into?



## Cache Write Issue (3/9)

- Tag/Data access
  - Writes are a pipelined two step (multi-cycle) process
    - Step 1: match tag
    - Step 2: write to matching way
    - Bypass (with address check) to avoid load stalls
    - May introduce structural hazards



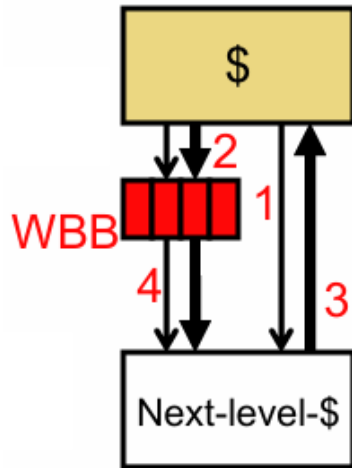
## Cache Write Issue (4/9)

- Write propagation: when to propagate new value to (lower level) memory?
  - **Option #1: Write-through:** immediately
    - On hit, update cache
    - Immediately send the write to the next level
  - **Option #2: Write-back:** when block is replaced
    - Requires additional “**dirty**” bit per block
    - Replace **clean** block: **no extra traffic**
    - Replace **dirty** block: **extra “writeback” of block**



## Cache Write Issue (5/9)

- **Option #2: Write-back:** when block is replaced
  - **Writeback-buffer (WBB)**
    - Hide latency of writeback (keep off critical path)
    - Step#1: Send “fill” request to next-level
    - Step#2: While waiting, write dirty block to buffer
    - Step#3: When new blocks arrives, put it into cache
    - Step#4: Write buffer contents to next-level





# Cache Write Issue (6/9)

- **Write-through**

- - Requires additional bus bandwidth
  - Consider repeated write hits
- - Next level must handle small writes (1, 2, 4, 8-bytes)
- + No need for dirty bits in cache
- + No need to handle “writeback” operations
  - Simplifies miss handling

- **Write-back**

- + Key advantages: uses less bandwidth
- Used in most CPU designs



# Cache Write Issue (7/9)

- **Write Miss Handling**

- **Write-allocate**: fill block from next level, then write it
  - + Decreases read misses (next read to block will hit)
  - - Requires additional bandwidth
  - Commonly used (especially with write-back caches)
- **Write-non-allocate**: just write to next level, no allocate
  - - Potentially more read misses
  - + Uses less bandwidth
  - Use with write-through





# Cache Write Issue (8/9)

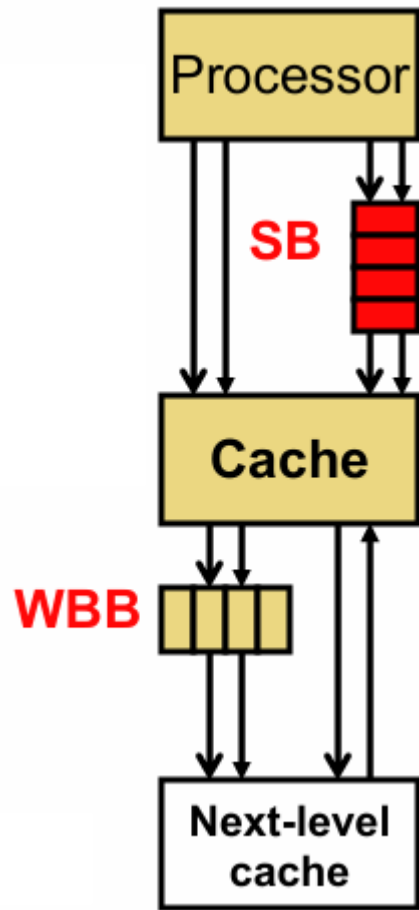
- **Write Miss and Store Buffers**
  - Read miss?
    - Load can't go on without the data
    - It must stall
  - Write miss?
    - No instruction is waiting for data
    - Why stall?



# Cache Write Issue (9/9)

- **Write Miss and Store Buffers**

- Stores put address/value to store buffer, keep going
- Store buffer writes stores to D\$ in the background
- Loads must search store buffer
- + Eliminates stalls on write misses (mostly)
- Store buffer vs. write-back buffer
  - **Store buffer**: in front of D\$, for hiding store misses
  - **Write-buffer**: behind D\$, for hiding writebacks





# Outline

- Associativity
- Fully Associative Cache
- N-way Set Associative Cache
- Cache Write Issue
- **Block Replacement Policy**
- Multi-level Caches



## Block Replacement Policy (1/3)

- If we have the choice, where should we write an incoming block?
  - If there are any locations with valid bit off (empty), then usually write the new block into the first one
  - If all possible locations already have a valid block, we must pick a **replacement policy**:
    - Rule by which we determine which block gets “cached out” on a miss



## Block Replacement Policy (2/3)

- On cache miss, which block in set to replace (kick out)?
  - If there are any locations with valid bit off (empty), then usually write the new block into the first one
  - If all possible locations already have a valid block, we must pick a **replacement policy**:
    - Rule by which we determine which block gets “cached out” on a miss



# Block Replacement Policy (3/3)

- **Block replacement options**

- Random
- FIFO (first-in first-out)
- LRU (least recently used)
  - Fit with temporal locality, LRU = least likely to be used in future
- NMRU (not most recently used)
  - Track which block in set is MRU
  - On replacement, pick a non-MRU block
  - One MRU pointer per set (vs. N LRU counters)



# Block Replacement Policy: LRU (1/8)

- **LRU (Least Recently Used)**

- **Idea:**

- cache out block which has been accessed (read or write)  
least recently

- **Pro:**

- temporal locality => recent past use implies likely future use

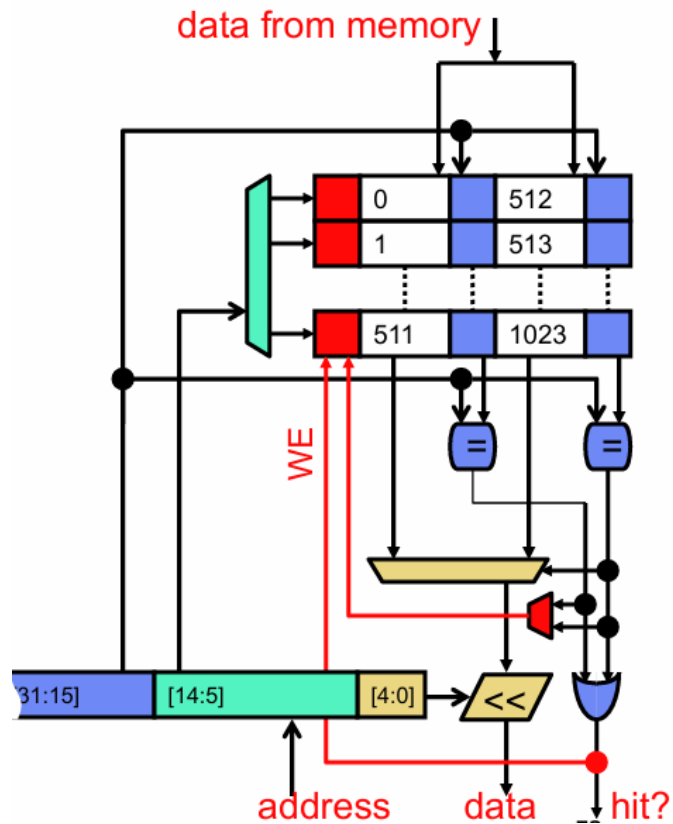
- **Con:**

- with 2-way set assoc, easy to keep track (on LRU bit)
- With 4-way or greater, requires complicated hardware and much time to keep track of this



# Block Replacement Policy: LRU (2/8)

- Add **LRU** field to each set
  - LRU data is encoded “way”
  - Hit? Update MRU
  - LRU bits updated on each access







## Block Replacement Policy: LRU (3/8)

- We have a 2-way set associative cache with a four word total capacity and one word blocks. We perform the following word access (ignore bytes for this problem)
  - 0, 2, 0, 1, 4, 0, 2, 3, 5, 4
- How many hits and how many misses will there be for the LRU block replacement policy?



# Block Replacement Policy: LRU (4/8)

**0: miss, bring into set 0 (loc 0)**

**2: miss, bring into set 0 (loc 1)**

**0: hit**

**1: miss, bring into set 1 (loc 0)**

Addresses 0, 2, 0, 1, 4, 0, ...

	loc 0	loc 1
set 0	0	<i>iru</i>
set 1		
set 0	<i>iru</i> 0	2
set 1		
set 0	0	<i>iru</i> 2
set 1		
set 0	0	<i>iru</i> 2
set 1	1	<i>iru</i>



# Block Replacement Policy: LRU (5/8)

1: miss, bring into set 1 (loc 0)

4: miss, bring into set 0 (loc 1, replace 2)

0: hit

	loc 0	loc 1
set 0	0	<i>lru</i> 2
set 1	1	<i>lru</i>

set 0	0	<i>lru</i> 4
set 1	1	<i>lru</i>

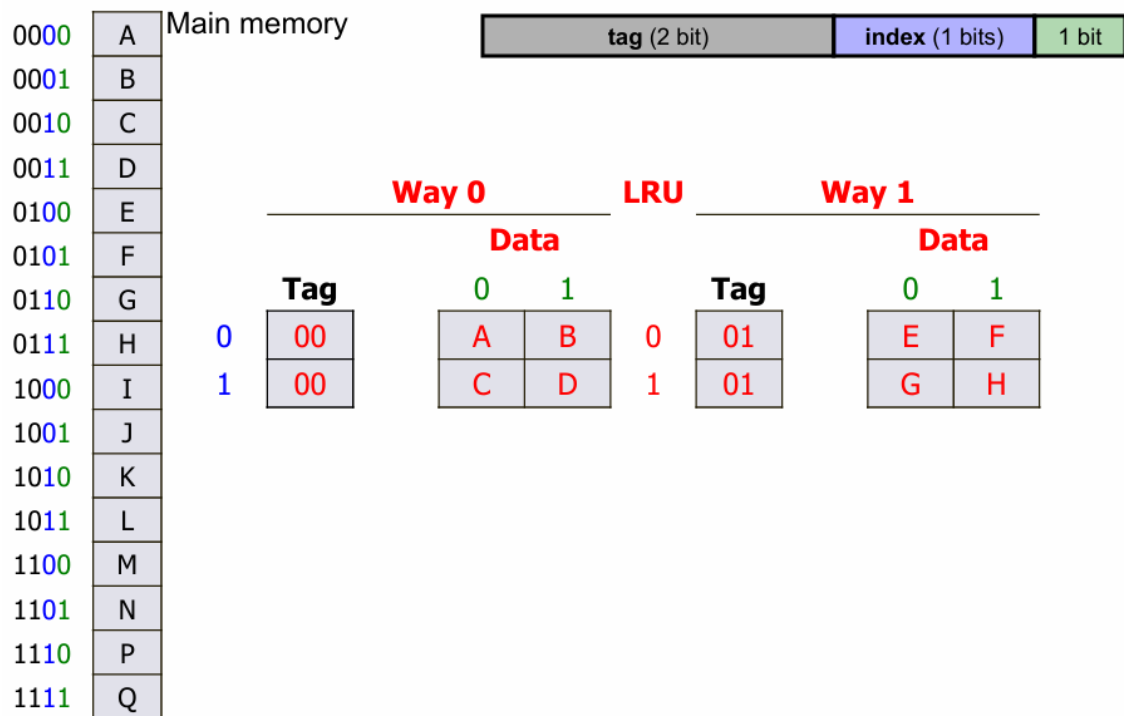
set 0	0	<i>lru</i> 4
set 1	1	<i>lru</i>

Addresses 0, 2, 0, 1, 4, 0, ...



# Block Replacement Policy: LRU (6/8)

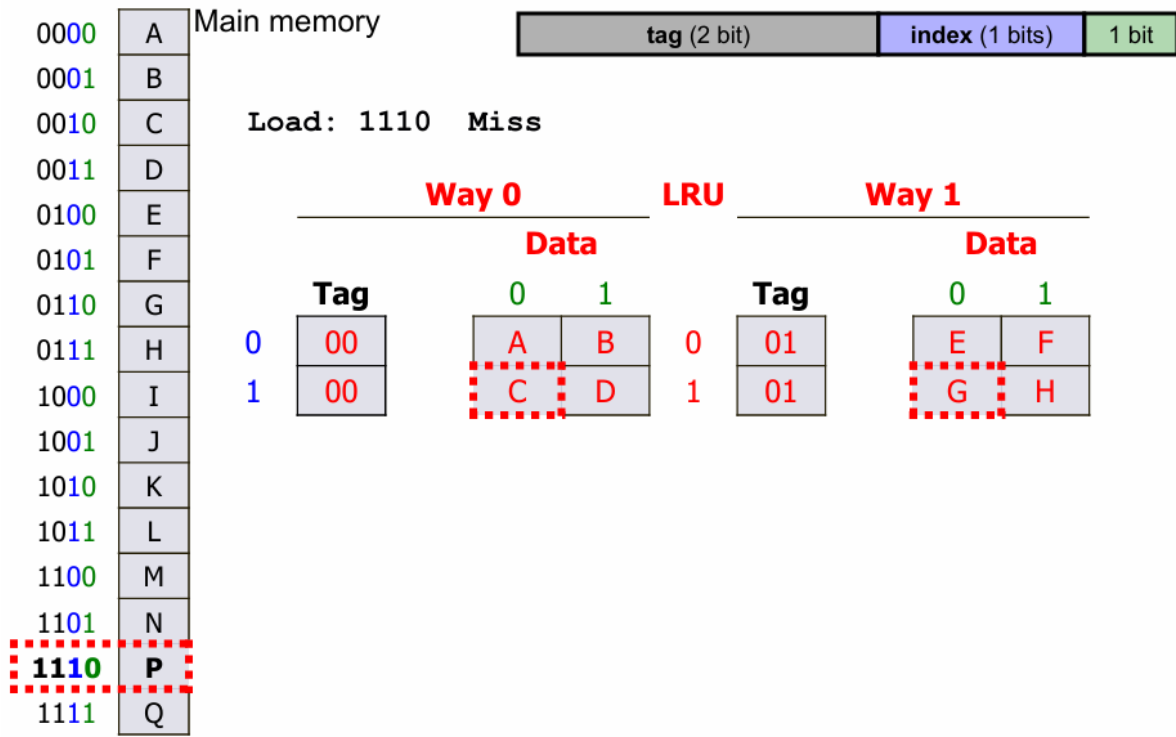
- 4-bit address, 8B Cache, 2B Blocks, 2-way





# Block Replacement Policy: LRU (7/8)

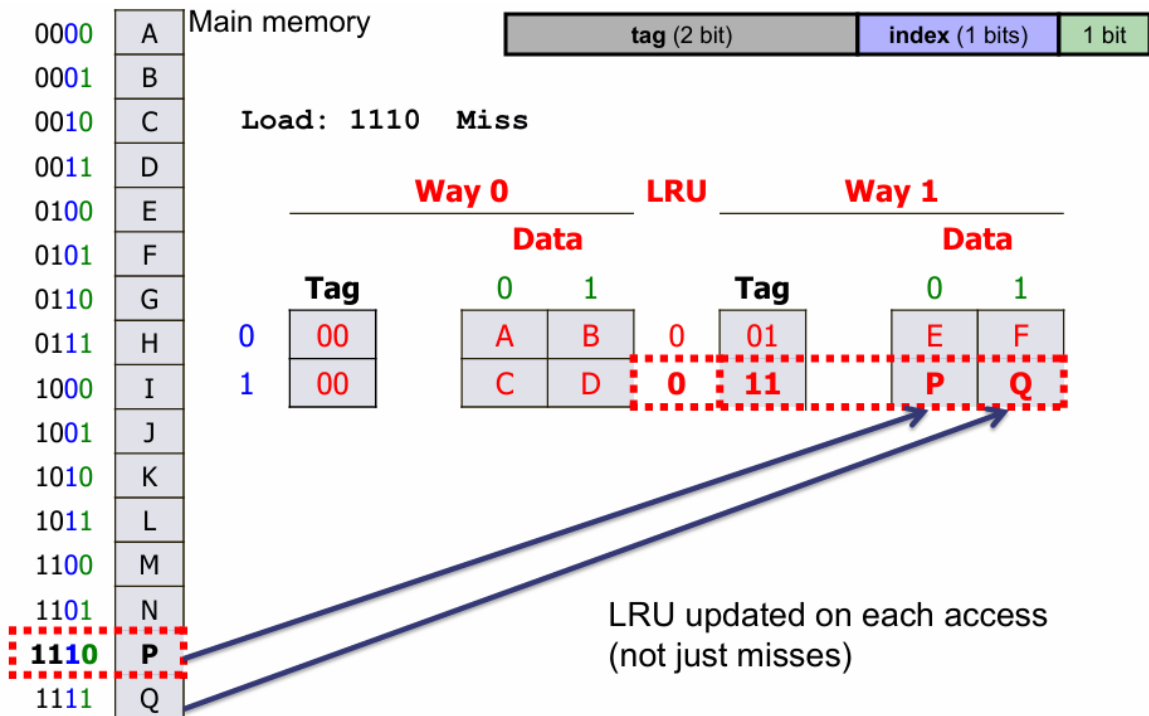
- 4-bit address, 8B Cache, 2B Blocks, 2-way





# Block Replacement Policy: LRU (8/8)

- 4-bit address, 8B Cache, 2B Blocks, 2-way





# Outline

- Associativity
- Fully Associative Cache
- N-way Set Associative Cache
- Cache Write Issue
- Block Replacement Policy
- **Multi-level Caches**



# Multi-level Caches (1/5)

- $T_{\text{access}}$  vs.  $\%_{\text{miss}}$  tradeoff
- Upper memory components (I\$, D\$) emphasize low  $t_{\text{access}}$ 
  - Frequent access ->  $t_{\text{access}}$  important
  - $T_{\text{miss}}$  is not bad ->  $\%_{\text{miss}}$  less important
  - Lower capacity and lower associativity (to reduce  $t_{\text{access}}$ )
  - Small-medium block-size (to reduce conflicts)
- Moving down (L2, L3) emphasis turns to  $\%_{\text{miss}}$ 
  - $T_{\text{miss}}$  is bad ->  $\%_{\text{miss}}$  important
  - High capacity, associativity, and block size (to reduce  $\%_{\text{miss}}$ )





## Multi-level Caches (2/5)

- Memory hierarchy parameters

Parameter	I\$/D\$	L2	L3	Main Memory
$t_{\text{access}}$	2ns	10ns	30ns	100ns
$t_{\text{miss}}$	<b>10ns</b>	<b>30ns</b>	<b>100ns</b>	<b>10ms (10M ns)</b>
Capacity	8KB–64KB	256KB–8MB	2–16MB	1-4GBs
Block size	16B–64B	32B–128B	32B-256B	NA
Associativity	2-8	4–16	4-16	NA



## Multi-level Caches (3/5)

- **Split vs. unified caches**
- **Split I\$/D\$:** instruction and data in different caches
  - To minimize structural hazards and  $t_{\text{access}}$
  - Larger unified I\$/D\$ would be slow, 2<sup>nd</sup> port even slower
  - Optimize I\$ and D\$ separately
    - Not write for I\$, smaller reads for D\$



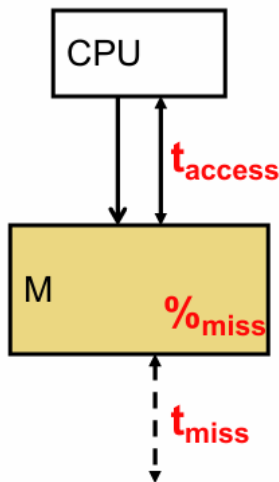
## Multi-level Caches (4/5)

- **Split vs. unified caches**
- **Unified L2, L3:** instruction and data together
  - To minimize  $\%_{\text{miss}}$
  - + Fewer capacity misses: unused instruction capacity can be used for data
  - - More conflict misses: instruction/data conflict
    - A much smaller effect in large caches
  - Instruction/data structural hazards are rare: simultaneous I\$/D\$ miss



# Multi-level Caches (5/5)

## • Memory performance equation



- For memory component M
  - **Access**: read or write to M
  - **Hit**: desired data found in M
  - **Miss**: desired data not found in M
    - Must get from another (slower) component
  - **Fill**: action of placing data in M
- $\%_{\text{miss}}$  (miss-rate):  $\# \text{misses} / \# \text{accesses}$
- $t_{\text{access}}$ : time to read data from (write data to) M
- $t_{\text{miss}}$ : time to read data into M

- Performance metric
  - $t_{\text{avg}}$ : average access time

$$t_{\text{avg}} = t_{\text{access}} + (\%_{\text{miss}} * t_{\text{miss}})$$



# Takeaway Question

- Parameters

- Baseline pipeline CPI = 1
- 30% of instructions are memory operations
- L1:  $t_{\text{access}} = 1$  cycle (included in CPI of 1),  $\%_{\text{miss}} = 5\%$  of accesses
- L2:  $t_{\text{access}} = 10$  cycle,  $\%_{\text{miss}} = 20\%$  of L2 accesses
- L2:  $t_{\text{access}} = 50$  cycle
- What is the new CPI?



# Takeaway Question

- Parameters

- 30% of instructions are memory operations
- L1:  $t_{\text{access}} = 1$  cycle (included in CPI of 1),  $\%_{\text{miss}} = 5\%$  of accesses
- L2:  $t_{\text{access}} = 10$  cycle,  $\%_{\text{miss}} = 20\%$  of L2 accesses
- L2:  $t_{\text{access}} = 50$  cycle
- What is the new CPI?
  - $\text{CPI} = 1 + 30\% * 5\% * t_{\text{missD\$}}$
  - $t_{\text{missD\$}} = t_{\text{avgL2}} = t_{\text{accL2}} + (\%_{\text{missL2}} * t_{\text{accMem}}) = 10 + (20\% * 50) = 20$  cycles
  - Thus,  $\text{CPI} = 1 + 30\% * 5\% * 20 = 1.3$  CPI



# Conclusion

- **Memory hierarchy**
  - Cache (SRAM) -> Memory (DRAM) -> swap (Disk)
  - Smaller, faster, more expensive
- **Cache ABCs (capacity, associativity, block size)**
  - 3C miss model: compulsory, capacity, conflict
- **Write issues**
  - Write-back vs. write-through/write-allocate vs. write-no-allocate