



# Lecture 10: Cache II

## **CS10014 Computer Organization**

Tsung Tai Yeh

Department of Computer Science  
National Yang Ming Chiao University



# 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
- Cache Replacement Policy
- Multi-level Caches





# Review: Cache Basics

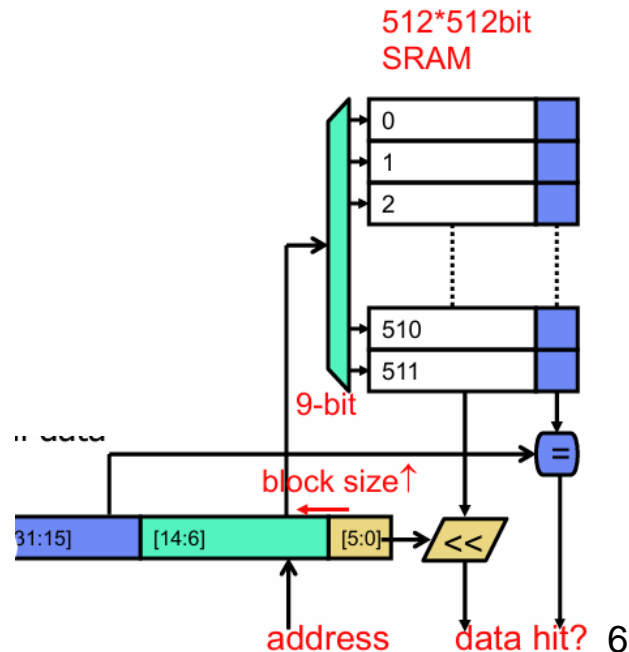
- 4-bit addresses -> 16 bytes memory
- 8 bytes cache, 2 bytes blocks
  - The number of sets: 4 (capacity / block size)
  - How address splits into offset/index/tag bits
    - **Offset:** least-significant  $\log_2(\text{block size}) = \log_2(2) = 1$  -> 000**0**
    - **Index:** next  $\log_2(\text{number-of-sets})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





# Review: Cache Basics

- 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





# Review: Cache Basics

- **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**

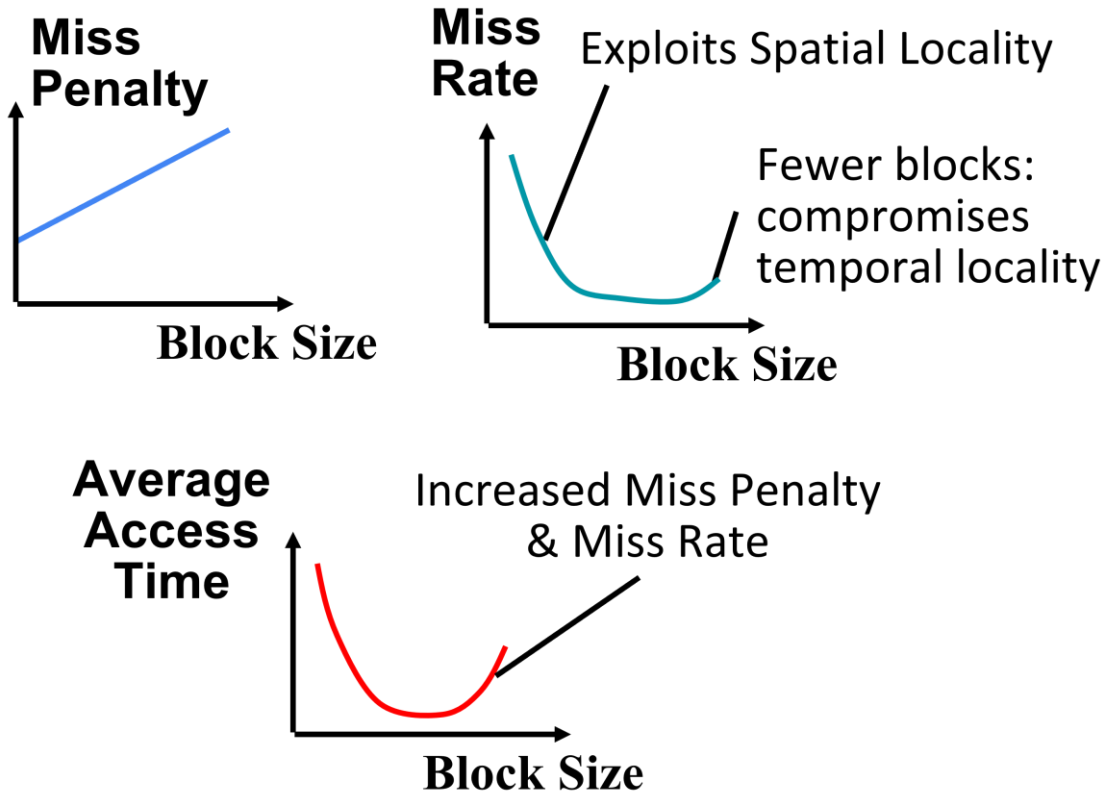
- % of requests that are found in current level cache

- **Miss Rate**

- $1 - \text{Hit Rate}$



# Review: Cache Basics





# Block Size Considerations

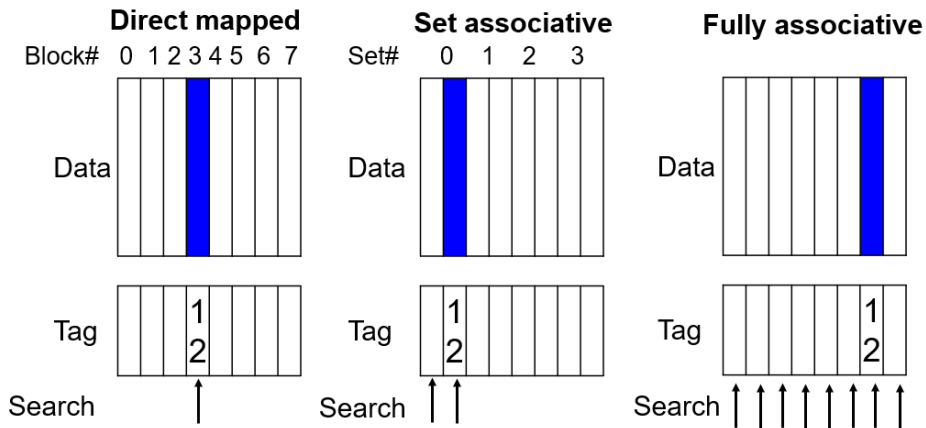
- Larger blocks should reduce miss rate
  - Due to spatial locality
- But in a fixed-sized cache
  - Larger blocks => fewer of them
    - More competition => increased miss rate
  - Larger blocks => pollution
- Larger miss penalty
  - Can override benefit of reduced miss rate
  - Early restart and critical-word-first can help



# Cache Basics

- **Cache associativity**

- In direct-mapped cache, there is only one cache block where memory block 12 can be found.  $(12 \text{ modulo } 8) = 4$
- In two-way set-associative cache, memory block 12 must be in set  $(12 \text{ mod } 4) = 0$
- In fully-associative, the block 12 can appear in any of 8 cache blocks





# Associativity

- **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 positions within the set on a miss

- **Fully associative**

- Block can be written into any positions



# Spectrum of Associativity

- For a cache with 8 entries

## 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 |  |  |  |  |

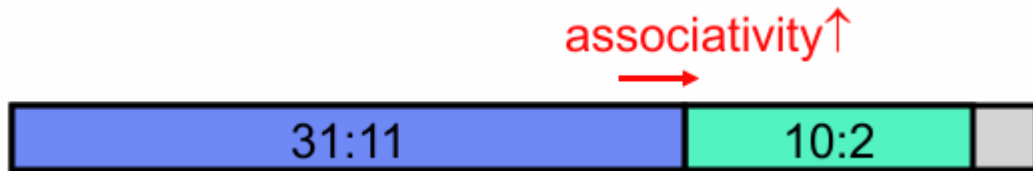
## Eight-way set associative (Fully Associative)

Tag Data Tag Data ....

|  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|
|  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|

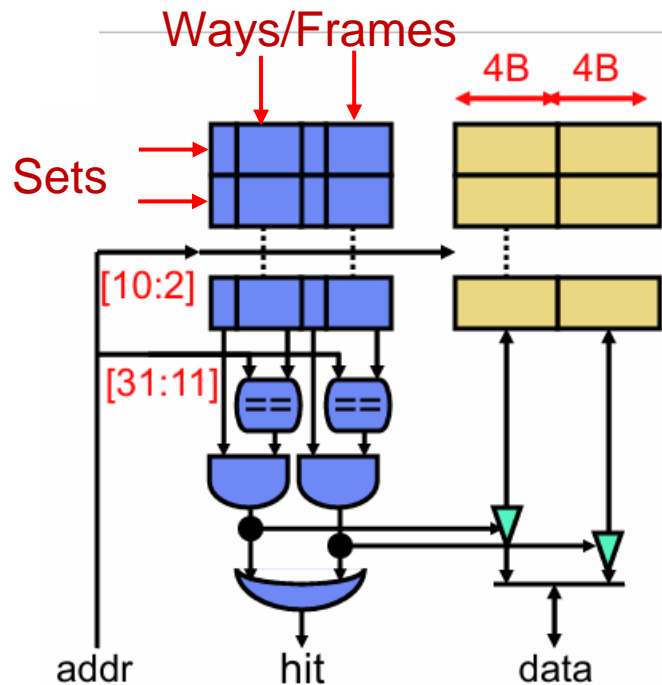


# Associativity



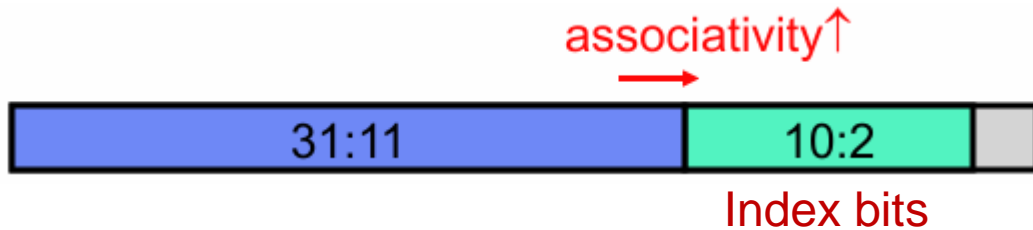
## • Set-associativity

- 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 -> **directed-mapped (DM)**
- 1-set-> **fully-associative (FA)**
- + Reduce conflicts
- - Increase latency<sub>hit</sub>
  - Additional tag match & muxing



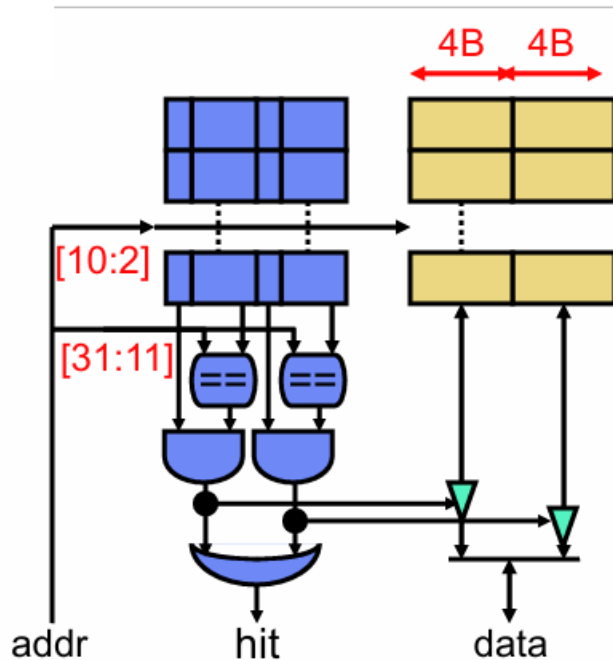


# Associativity



- Lookup algorithm

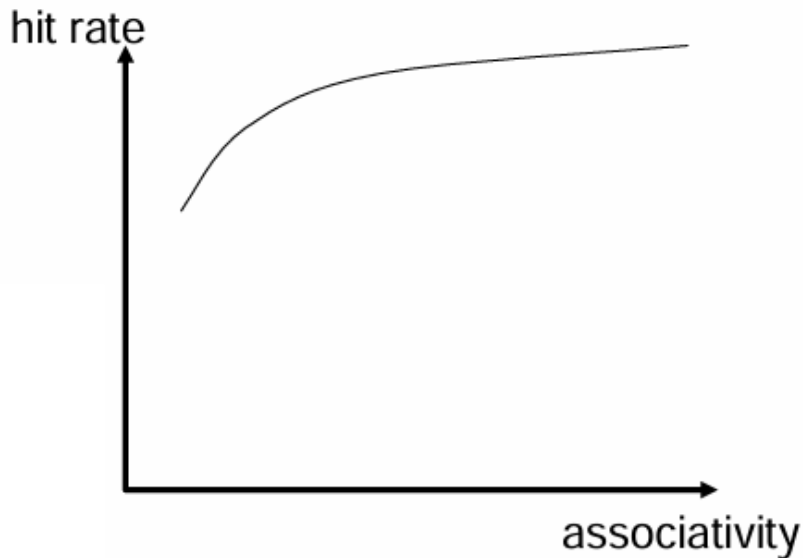
- Use index bits to find set
- Read data/tags in all frames in parallel
- **Any** (match and valid bit), Hit
- Notice tag/index/offset bits
  - Only 9-bit index (versus 10-bit for direct mapped)





# Associativity

- How many blocks can be presented 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





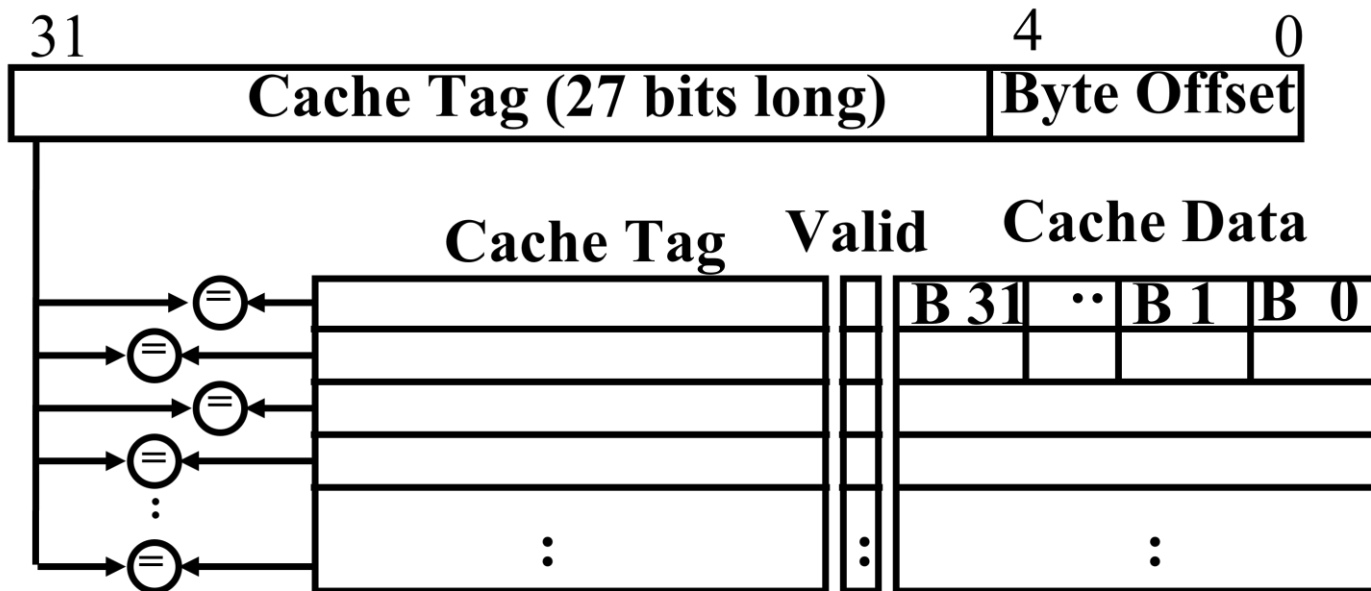
# Fully Associative Cache

- Memory address fields
  - Tag: same as before
  - Offset: same as before
  - **Index: Non-existent**
- 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

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

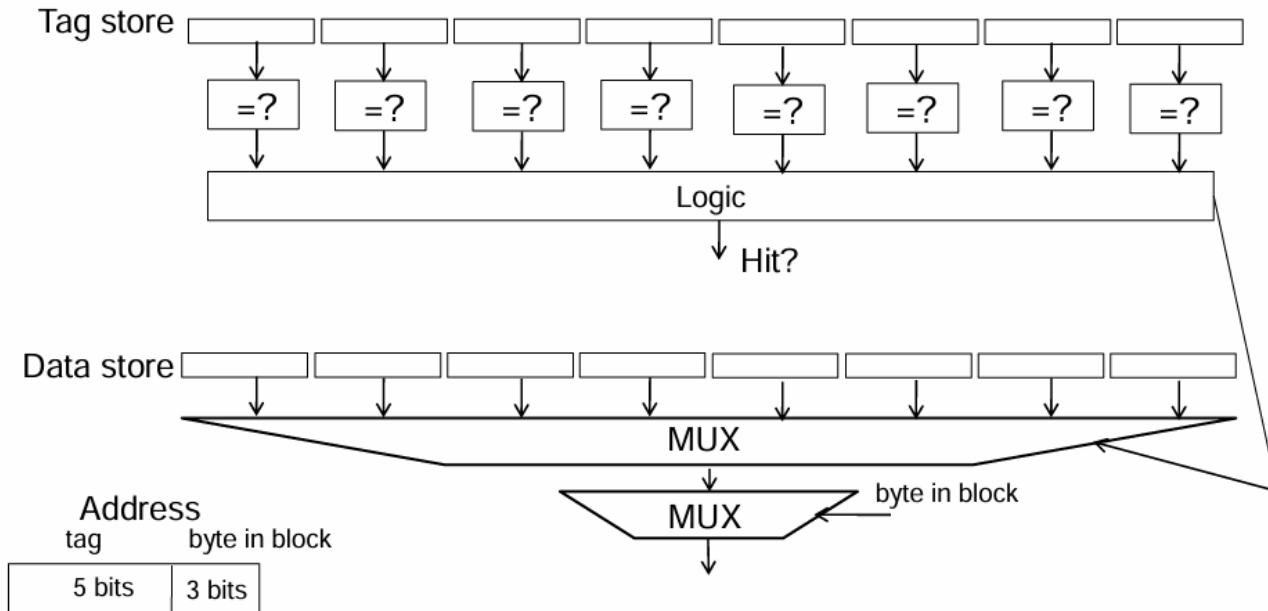




# Fully Associative Cache

- Fully Associative Cache

- A block can be placed in any cache location





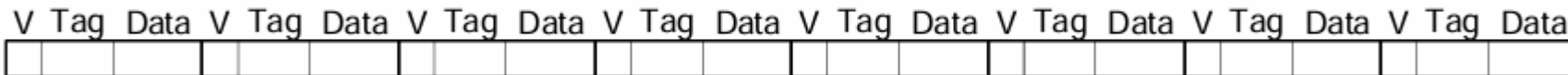
# Fully Associative Cache

- **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





# Associativity Example

- Compare 4-block caches
  - Block access sequence: 0, 8, 0, 6, 8
- Direct mapped

| Block address | Cache index | Hit/miss | Cache content after access |   |        |   |
|---------------|-------------|----------|----------------------------|---|--------|---|
|               |             |          | 0                          | 1 | 2      | 3 |
| 0             | 0           | miss     | Mem[0]                     |   |        |   |
| 8             | 0           | miss     | Mem[8]                     |   |        |   |
| 0             | 0           | miss     | Mem[0]                     |   |        |   |
| 6             | 2           | miss     | Mem[0]                     |   | Mem[6] |   |
| 8             | 0           | miss     | Mem[8]                     |   | Mem[6] |   |



# Associativity Example

- Compare 4-block caches
  - Block access sequence: 0, 8, 0, 6, 8
- 2-way set associative

| Block address | Cache index | Hit/miss | Cache content after access |               |       |  |
|---------------|-------------|----------|----------------------------|---------------|-------|--|
|               |             |          | Set 0                      |               | Set 1 |  |
| 0             | 0           | miss     | <b>Mem[0]</b>              |               |       |  |
| 8             | 0           | miss     | Mem[0]                     | <b>Mem[8]</b> |       |  |
| 0             | 0           | hit      | <b>Mem[0]</b>              | Mem[8]        |       |  |
| 6             | 0           | miss     | Mem[0]                     | <b>Mem[6]</b> |       |  |
| 8             | 0           | miss     | <b>Mem[8]</b>              | Mem[6]        |       |  |



# Associativity Example

- Compare 4-block caches
  - Block access sequence: 0, 8, 0, 6, 8
- Fully associative

| Block address |  | Hit/miss | Cache content after access |               |               |  |
|---------------|--|----------|----------------------------|---------------|---------------|--|
| 0             |  | miss     | <b>Mem[0]</b>              |               |               |  |
| 8             |  | miss     | Mem[0]                     | <b>Mem[8]</b> |               |  |
| 0             |  | hit      | <b>Mem[0]</b>              | Mem[8]        |               |  |
| 6             |  | miss     | Mem[0]                     | Mem[8]        | <b>Mem[6]</b> |  |
| 8             |  | hit      | Mem[0]                     | <b>Mem[8]</b> | Mem[6]        |  |

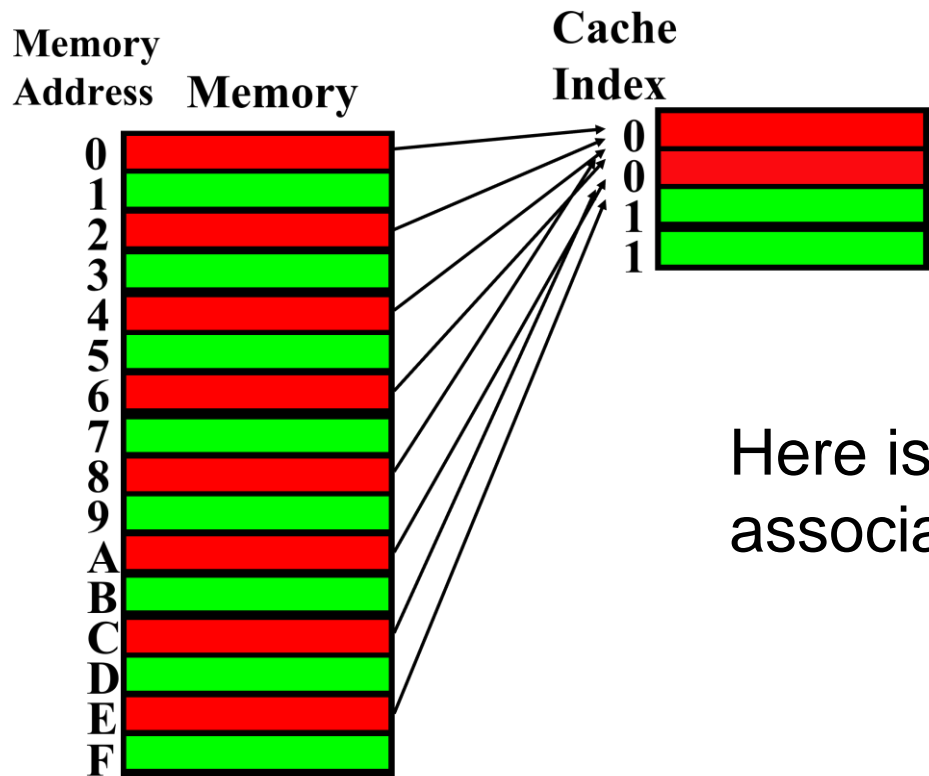


# N-way Set Associative Cache

- 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



Here is a simple 2-way set associative cache



# N-way Set Associative Cache

- **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



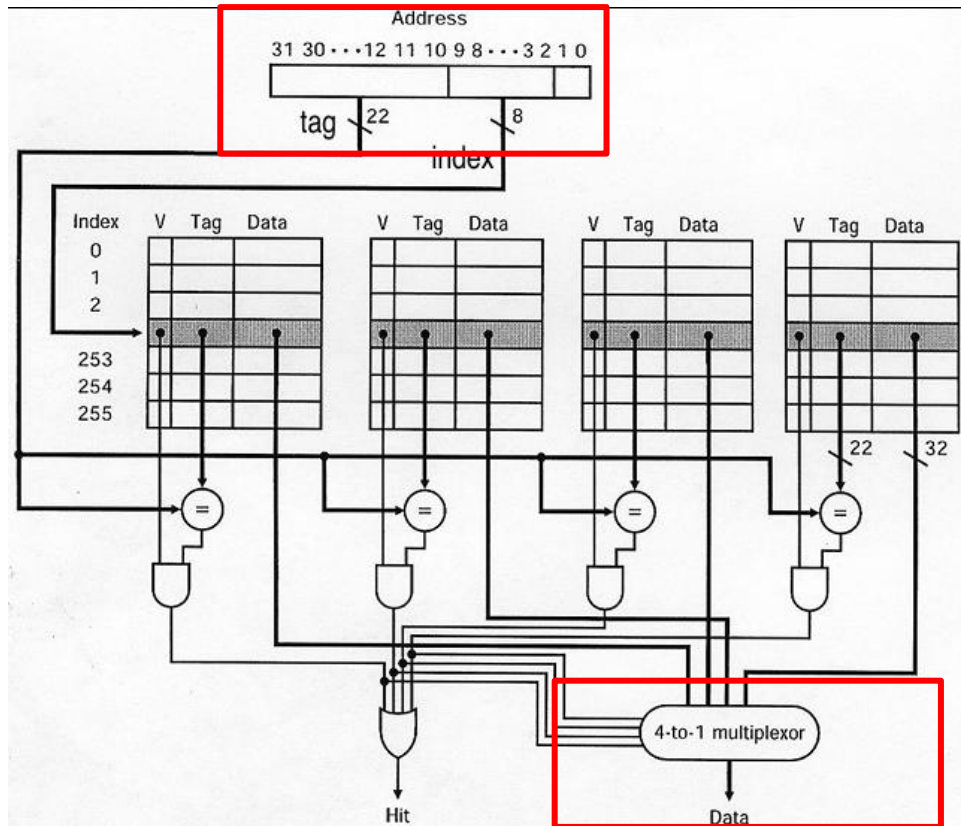
# N-way Set Associative Cache

- 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



4-way set  
associative  
cache

# ways = index  
length / offset  
length



# N-way Set Associative Performance

```
# 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

```
# 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 |



# Block Placement

- Replacement policy
  - Determine which block gets “cached out” on a miss
- Direct mapped (1-way associative)
  - One choice for placement
- N-way set associative
  - N choices within a set
- Fully associative
  - Any location

## Two-way set associative

| Set | Tag | Data | Tag | Data |
|-----|-----|------|-----|------|
| 0   |     |      |     |      |
| 1   |     |      |     |      |
| 2   |     |      |     |      |
| 3   |     |      |     |      |



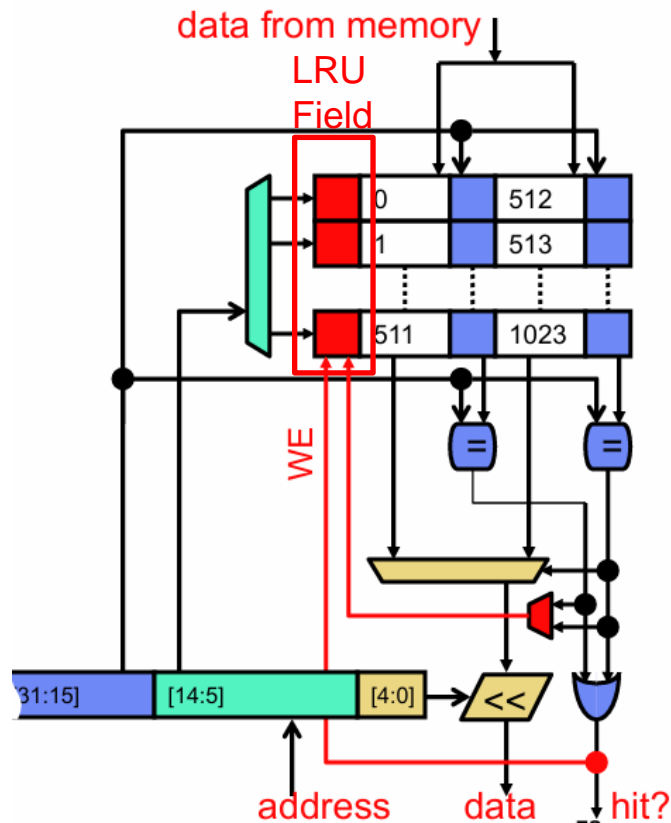
# Replacement Policy

- Choice of entry to replace on a miss
- Least-recently used (LRU)
  - Choose the one unused for the longest time
  - Temporal locality => recent past use implies likely future use
  - Simple for 2-way, manageable for 4-way, too hard beyond that
- Random
  - Gives approximately the same performance as LRU for high associativity



# Replacement Policy: LRU

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





# LRU Example

- A 2-way set associative cache
  - 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/misses will there be for the LRU?



# LRU Replacement Policy

**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>lru</i>   |
| set 1 |              |              |
| set 0 | <i>lru</i> 0 | 2            |
| set 1 |              |              |
| set 0 | 0            | <i>lru</i> 2 |
| set 1 |              |              |
| set 0 | 0            | <i>lru</i> 2 |
| set 1 | 1            | <i>lru</i>   |



# LRU Replacement Policy

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>   |

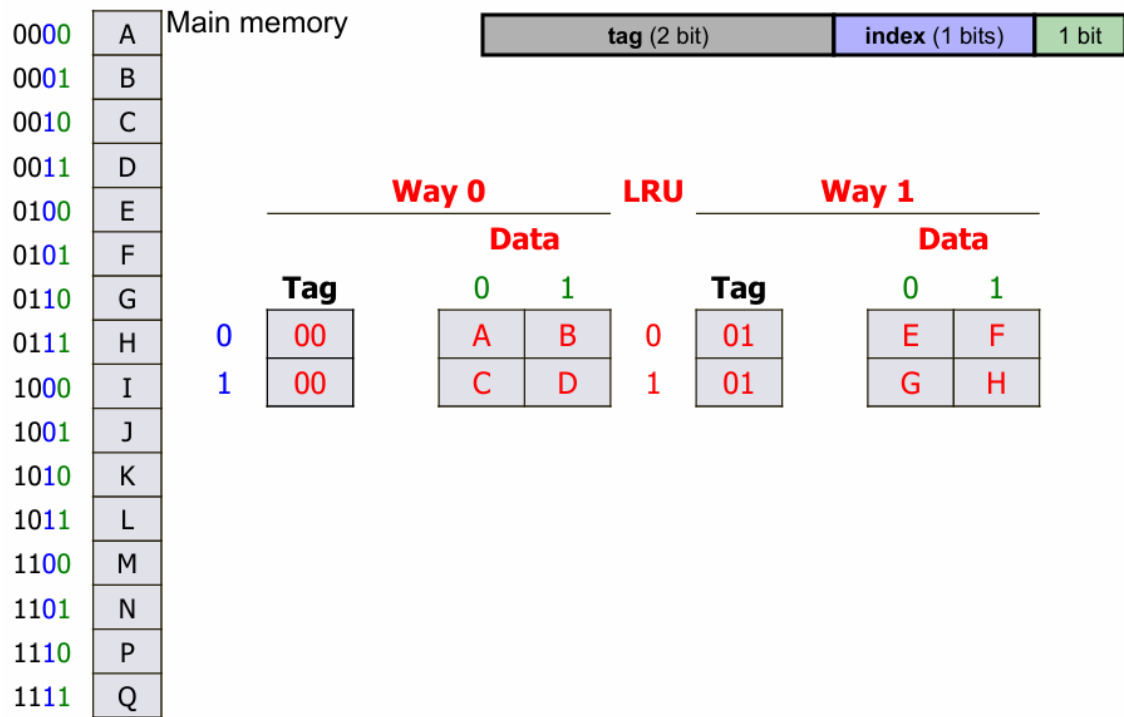
| set 0 | 0 | <i>lru</i> 4 |
|-------|---|--------------|
| set 1 | 1 | <i>lru</i>   |

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



# LRU Replacement Policy

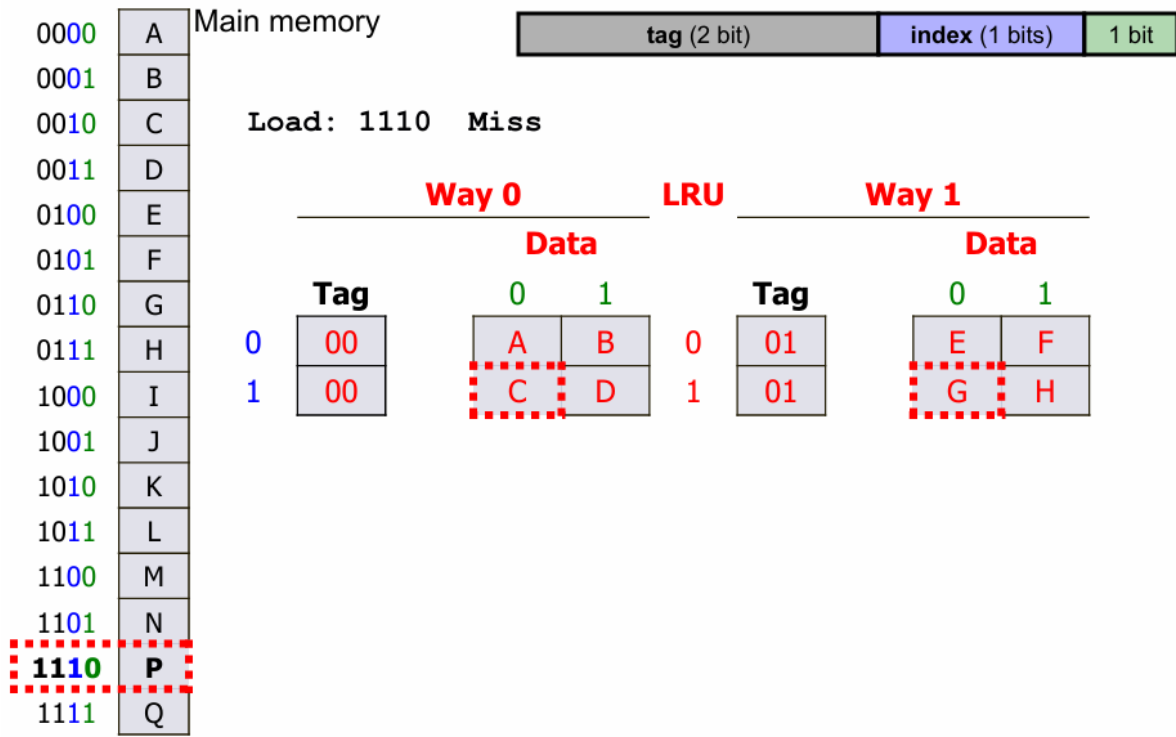
- 4-bit address, 8B Cache, 2B Blocks, 2-way





# LRU Replacement Policy

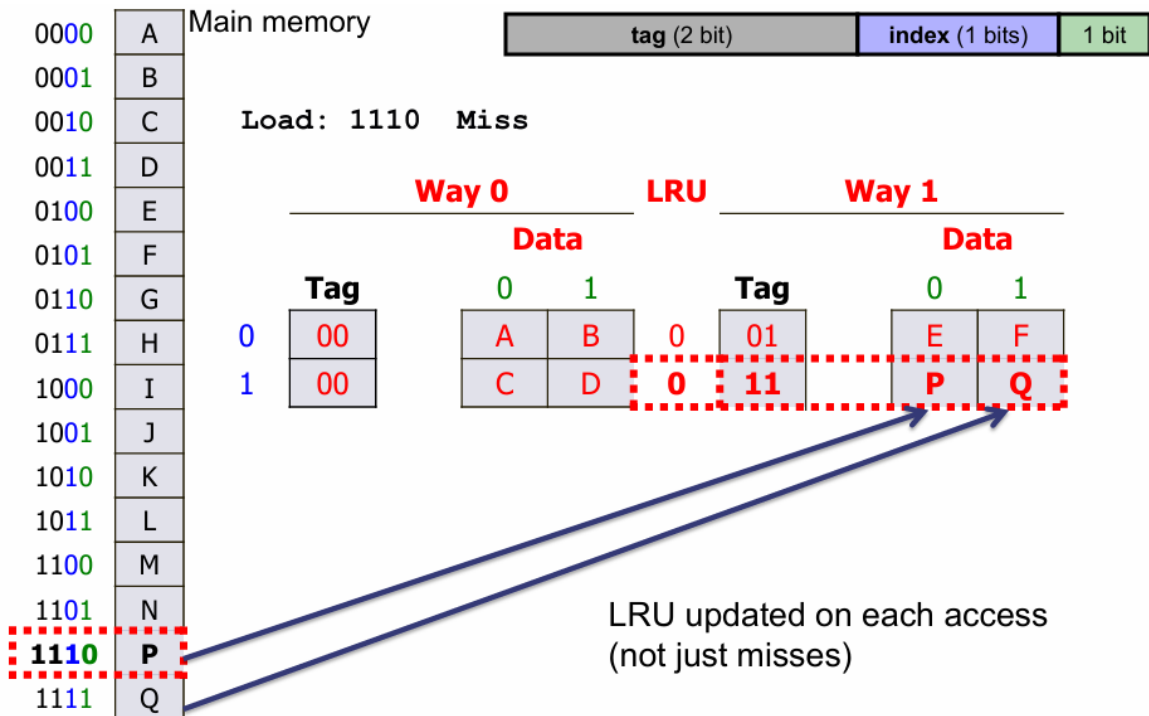
- 4-bit address, 8B Cache, 2B Blocks, 2-way





# LRU Replacement Policy

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





# Cache Write Issue

- 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

- 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

- 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



# Write Policy

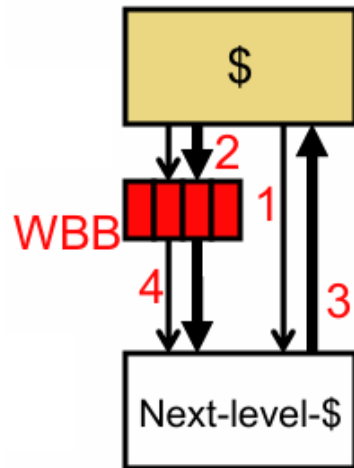
- **Write-through**

- Update both upper and lower levels
- Simplifies replacement, but may require write buffer
  - On hit, update cache
  - Immediately send the write to the next level
- But make writes take longer
  - If base CPI = 1, 10% of instructions are stores, write to memory takes 100 cycles
    - Effective CPI =  $1 + 0.1 \times 100 = 11$



# Write Policy

- Write buffer (WBB)
  - Hold data waiting to be written to memory
    - 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
  - CPU continues immediately
    - Only stalls on write if write buffer is already full





# Write Policy

- **Write-back**

- On data-write hit, just update the block in cache
  - Keep track of whether each block is dirty
- When a dirty block is replaced
  - Write it back to memory
  - Can use a write buffer to allow replacing block to be read first
- Need to keep more state
  - Requires additional “**dirty**” bit per block
  - Replace **clean** block: **no extra traffic**
  - Replace **dirty** block: **extra “writeback” of block**



# Write Policy

- **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



# Write Policy

- **Write Allocation**

- What should happen on a write miss
- Write miss
  - Occurs when a processor attempts to write data to a memory location
  - But that specific data block is not present in the cache
- Alternative for write-through
  - Allocate on miss: fetch the block
  - Write around: don't fetch the block
- For write-back
  - Usually fetch the block



# Cache Write Issue

- **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

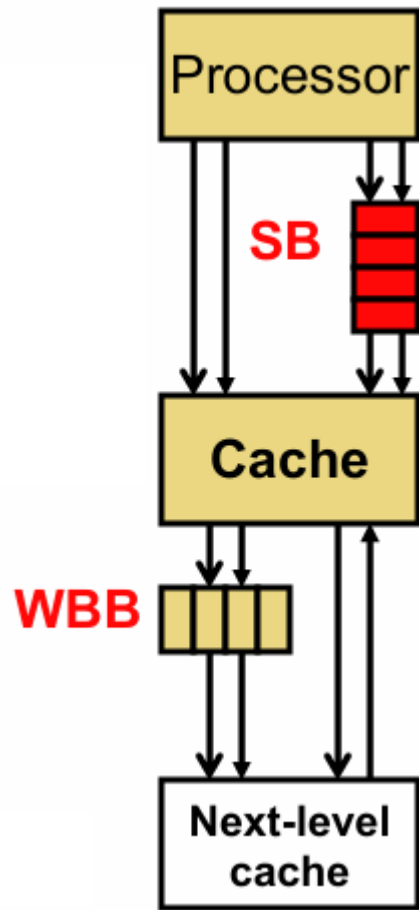
- **Write Miss and Store Buffers**
  - Read miss?
    - Load can't go on without the data
    - It must stall
  - Write miss?
    - No instruction is waiting for data
    - Why stall?



# Cache Write Issue

- **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





# Classifying Misses: 3C Model

- Divide cache misses into three categories
  - **Compulsory (cold)**: never seen this address before
    - **Would miss even in infinite cache**
  - **Capacity**: miss caused because cache is too small
    - **Would miss even in fully associative cache**
    - Identify? Consecutive accesses to block separated by access to at least  $N$  other distinct blocks ( $N$  is number of frames in cache)
  - **Conflict**: miss caused because cache associativity is too low
    - Identify? **All other misses**
  - **(Coherence)**: miss due to external invalidations
    - Only in shared memory multiprocessors (later)



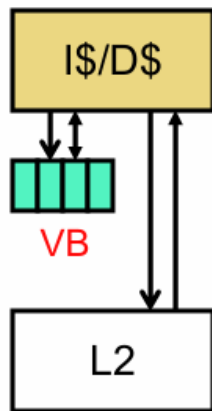
# Miss Rate: ABC

- Why do we care about 3C miss model?
  - So that we know what to do to eliminate misses
  - If you don't have conflict misses, increasing associativity won't help
- **More associativity** (assuming fixed capacity)
  - + Decreases conflict misses
  - Increases latency<sub>hit</sub>
- **Larger block size** (assuming fixed capacity)
  - Increases conflict/capacity misses (fewer frames)
  - + Decreases compulsory misses (spatial locality)
    - No significant effect on latency<sub>hit</sub>
- **More capacity**
  - + Decreases capacity misses
  - Increases latency<sub>hit</sub>



# Reducing Conflict Misses: Victim Buffer

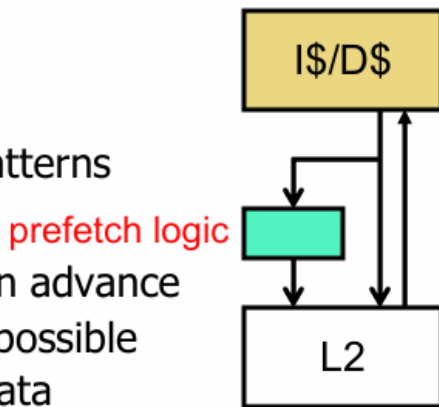
- Conflict misses: not enough associativity
  - High associativity is expensive, but also rarely needed
    - 3 blocks mapping to same 2-way set and accessed (XYZ)+
- **Victim buffer (VB)**: small fully-associative cache
  - Sits on I\$/D\$ miss path
  - Small (e.g., 8 entries) so very fast
  - Blocks 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
    - + Only a few sets will need it at any given time
  - + Very effective in practice





# Prefetching

- Bring data into cache proactively/**speculatively**
  - If successful, reduces number of caches misses
- Key: anticipate upcoming miss addresses accurately
  - Can do in software or hardware
- Simple hardware prefetching: **next block prefetching**
  - Miss on address **X** → anticipate miss on **X+block-size**
  - + Works for insns: sequential execution
  - + Works for data: arrays
- Table-driven hardware prefetching
  - Use **predictor** to detect strides, common patterns
- Effectiveness determined by:
  - **Timeliness**: initiate prefetches sufficiently in advance
  - **Coverage**: prefetch for as many misses as possible
  - **Accuracy**: don't pollute with unnecessary data





# Software Prefetching

- Use a special “prefetch” instruction
  - Tells the hardware to bring in data
  - Just a hint
- Inserted by programmer or compiler

- Example

```
int tree_add(tree_t* t) {  
    if (t == NULL) return 0;  
    __builtin_prefetch(t->left);  
    __builtin_prefetch(t->right);  
    return t->val + tree_add(t->right) + tree_add(t->left);  
}
```

- Multiple prefetches bring multiple blocks in parallel
  - More “Memory-level” parallelism (MLP)



# Software Restructuring: Data

- Capacity misses: poor spatial or temporal locality
  - Several code restructuring techniques to improve both
  - Compiler must know that restructuring preserves semantics
- **Loop interchange**: spatial locality
  - Example: row-major matrix:  $x[i][j]$  followed by  $x[i][j+1]$
  - Poor code:  $x[i][j]$  followed by  $x[i+1][j]$ 

```
for (j = 0; j<NCOLS; j++)
    for (i = 0; i<NROWS; i++)
        sum += X[i][j];
```
  - Better code

```
for (i = 0; i<NROWS; i++)
    for (j = 0; j<NCOLS; j++)
        sum += X[i][j];
```



# Software Restructuring: Data

- **Loop blocking**: temporal locality

- Poor code

```
for (k=0; k<NUM_ITERATIONS; k++)  
    for (i=0; i<NUM_ELEMS; i++)  
        X[i] = f(X[i]);    // say
```

- Better code

- Cut array into CACHE\_SIZE chunks
- Run all phases on one chunk, proceed to next chunk

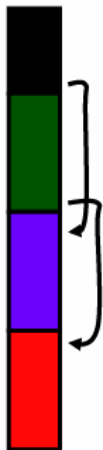
```
for (i=0; i<NUM_ELEMS; i+=CACHE_SIZE)  
    for (k=0; k<NUM_ITERATIONS; k++)  
        for (j=0; j<CACHE_SIZE; j++)  
            X[i+j] = f(X[i+j]);
```

– Assumes you know CACHE\_SIZE, do you?



# Software Restructuring: Code

- Compiler can lay out code for temporal and spatial locality
  - If (a) { **code1**; } else { **code2**; } **code3**;
  - But, **code2** case never happens (say, error condition)



- + Better locality
- + Fewer taken branches



- + Better locality for code after **code3**
- + Fewer taken branches



# Multi-level Caches

- Primary cache attached to CPU
  - Small, but fast
- Level-2 cache services misses from primary cache
  - Larger, slower, but still faster than main memory
- Main memory services L2 cache misses
- Some high-end systems include L3 cache



# Multi-level Cache Example

- Given

- CPU base CPI = 1, clock rate = 4GHz
- Miss rate/instruction = 2%
- Main memory access time = 100 ns

- With just primary cache

- Miss penalty =  $100\text{ns}/0.25\text{ns} = 400$  cycles
- Effective CPI =  $1 + 0.02 \times 400 = 9$



# Multi-level Cache Example

- Now add L2 cache
  - Access time = 5 ns
  - Global miss rate to main memory = 0.5%
- Primary miss with L2 hit
  - Penalty =  $5\text{ns}/0.25\text{ns} = 20$  cycles
- Primary miss with L2 miss
  - Extra penalty = 500 cycles
- $\text{CPI} = 1 + 0.02 \times 20 + 0.005 \times 400 = 3.4$
- Performance ratio =  $9/3.4 = 2.6$



# Multi-level Caches

- **Inclusion**
  - Bring block from memory into L2 then L1
    - A block in the L1 is always in the L2
  - If block evicted from L2, must also evict it from L1
    - Why? more on this when we talk about multicore
- **Exclusion**
  - Bring block from memory into L1 but not L2
    - Move block to L2 on L1 eviction
      - L2 becomes a large victim cache
    - Block is either in L1 or L2 (never both)
  - Good if L2 is small relative to L1
    - Example: AMD's Duron 64KB L1s, 64KB L2



# Multi-level Cache Considerations

- Primary cache
  - Focus on minimal hit time
- L2 cache
  - Focus on low miss rate to avoid main memory access
  - Hit time has less overall impact
- Results
  - L1 cache usually smaller than a single cache
  - L1 block size smaller than L2 block size



# 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
- DRAM:  $t_{\text{access}} = 50$  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
- DRAM:  $t_{\text{access}} = 50$  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



# Cache Performance Example

- Given
  - I-cache miss rate = 2%
  - D-cache miss rate = 4%
  - Miss penalty = 100 cycles
  - Base CPI (ideal cache) = 2
  - Load & Stores are 36% of instructions
- Miss cycles per instruction
  - I-cache:  $0.02 \times 100 = 2$
  - D-cache:  $0.36 \times 0.04 \times 100 = 1.44$
- Actual CPI =  $2 + 2 + 1.44 = 5.44$ 
  - Ideal CPU is  $5.44/2 = 2.72$  times faster



# Conclusion

- When CPU performance increased
  - Miss penalty becomes more significant
- Decreasing base CPI
  - Greater proportion of time spent on memory stalls
- Increasing clock rate
  - Memory stalls account for more CPU cycles
- Can't neglect cache behavior when evaluating system performance