

---

# **Chapter 5**

## **Large and Fast : Exploiting the Memory Hierarchy**

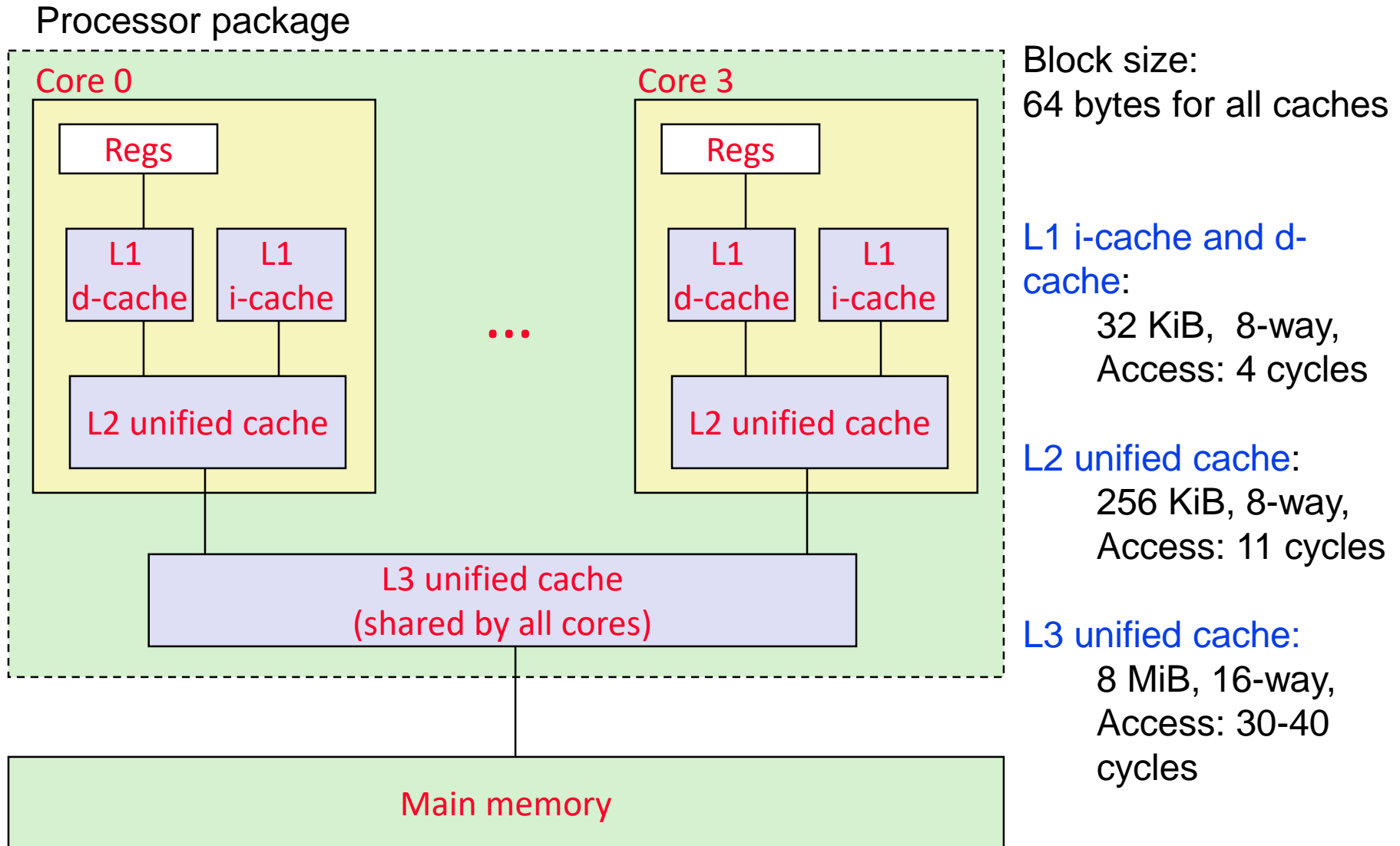
### **Cache Memory**

# Review : Memory Hierarchies

---

- ❑ Fundamental idea of a memory hierarchy:
  - | For each level  $k$ , the faster, smaller device at level  $k$  serves as a cache for the larger, slower device at level  $k+1$
  
- ❑ Why do memory hierarchies work?
  - | Because of *locality*, programs tend to access the data at level  $k$  more often than they access the data at level  $k+1$
  - | Thus, the storage at level  $k+1$  can be slower, and thus larger and cheaper per bit
  
- ❑ *Big Idea*: The memory hierarchy creates a large pool of storage that costs as much as the cheap storage near the bottom, but that serves data to programs at the rate of the fast storage near the top

# Example : Intel Core i7 Cache Hierarchy



# Cache Basics

---

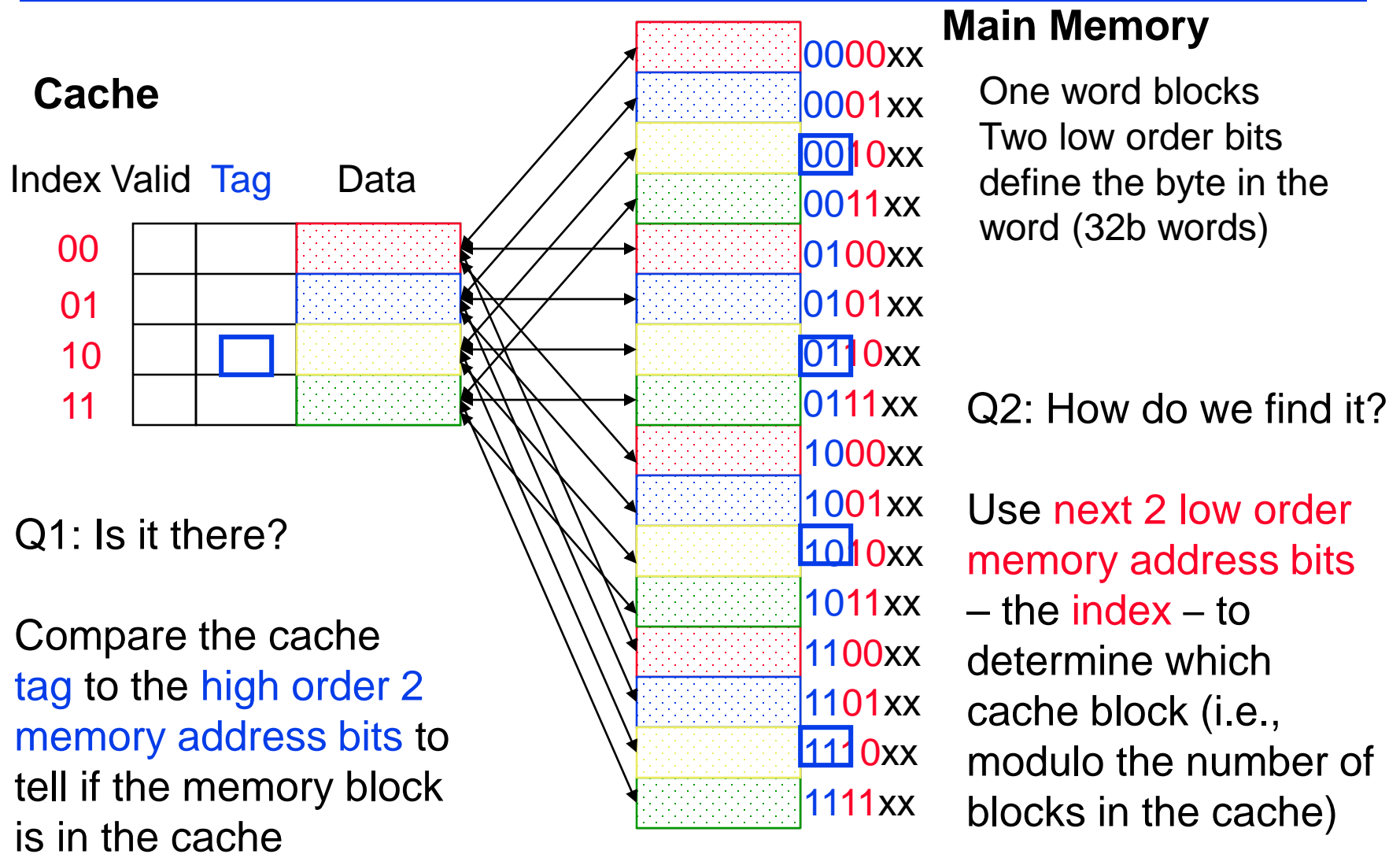
## ❑ Two questions to answer (in hardware):

- | Q1: How do we know if a data item is in the cache?
- | Q2: If it is, how do we find it?

## ❑ Direct mapped

- | Each memory block is mapped to exactly one block in the cache
  - lots of lower level blocks must **share** blocks in the cache
- | Address mapping (to answer Q2): **cache index** =  
**(block address) modulo (# of blocks in the cache)**
- | Have a **tag** associated with each cache block that contains the address information (the upper portion of the address) required to identify the block (to answer Q1)

# Caching: A Simple First Example



---

❑ Place Data in Cache **by Hashing Address!**

❑ **Tags** Differentiate Blocks in Same Index

# Direct Mapped Cache

❑ Consider the main memory word reference string

Start with an empty cache - all blocks initially marked as not valid

0 1 2 3 4 3 4 15

**0 miss**

00	Mem(0)

**1 miss**

00	Mem(0)
00	Mem(1)

**2 miss**

00	Mem(0)
00	Mem(1)
00	Mem(2)

**3 miss**

00	Mem(0)
00	Mem(1)
00	Mem(2)
00	Mem(3)

**4 miss**

01

<del>00</del>	<del>Mem(0)</del>
00	Mem(1)
00	Mem(2)
00	Mem(3)

4

**3 hit**

01	Mem(4)
00	Mem(1)
00	Mem(2)
00	Mem(3)

**4 hit**

01	Mem(4)
00	Mem(1)
00	Mem(2)
00	Mem(3)

**15 miss**

11

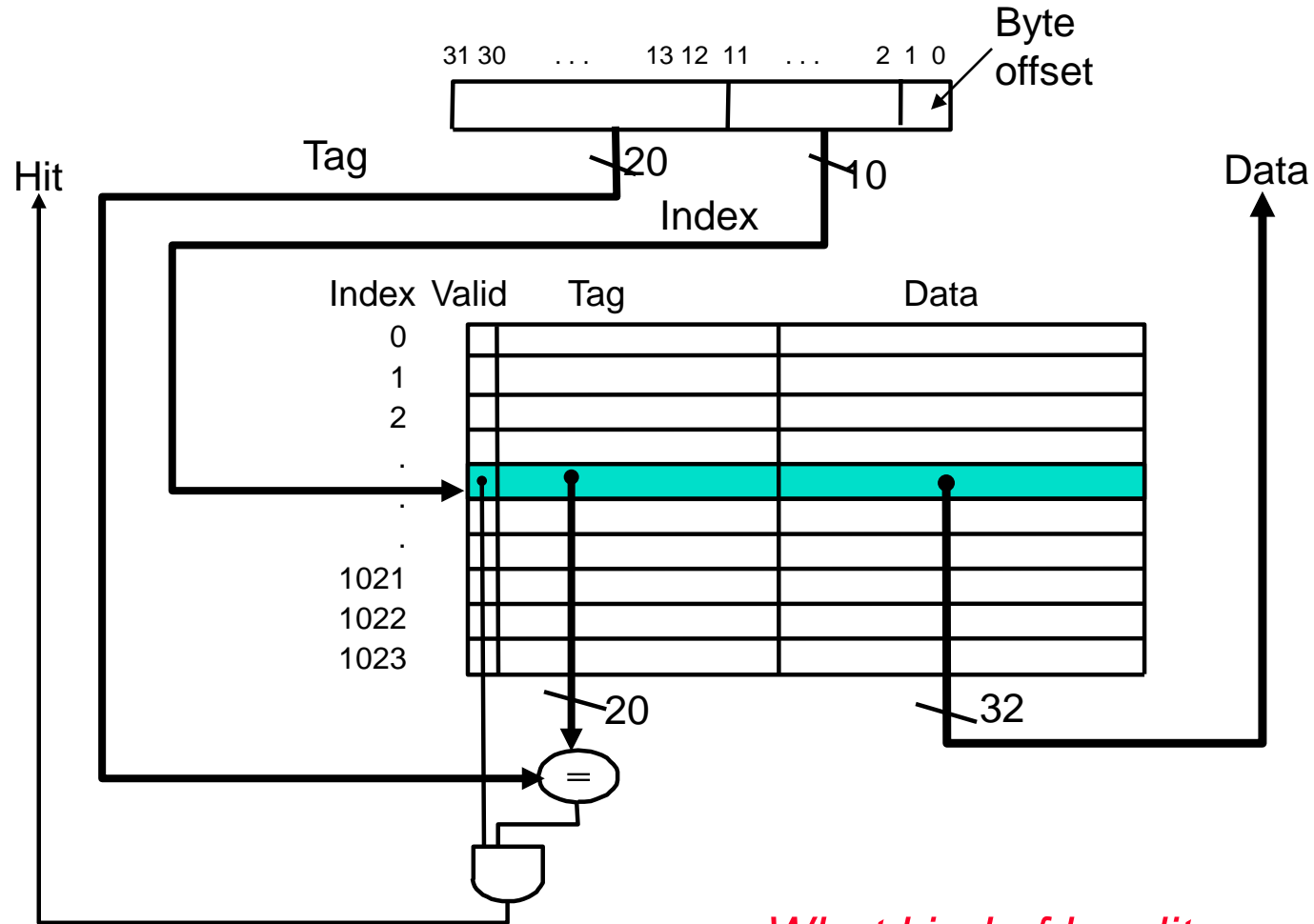
01	Mem(4)
00	Mem(1)
00	Mem(2)
<del>00</del>	<del>Mem(3)</del>

15

| 8 requests, 6 misses

# Direct Mapped Cache Example (Fig. 5.10)

- One word blocks, cache size = 1K words (or 4KB)



When is there a cache hit?  
( Valid[index] = TRUE ) **AND**  
( Tag[ index ] = Tag[ memory address ] )

*What kind of locality are we taking advantage of?*



# Cache Field Sizes

---

- ❑ The number of bits in a cache includes both the storage for data and for the tags
  - | 32-bit byte address
  - | For a direct mapped cache with  $2^n$  blocks,  $n$  bits are used for the index
  - | For a block size of  $2^m$  words ( $2^{m+2}$  bytes),  $m$  bits are used to address the word within the block and 2 bits are used to address the byte within the word
  
- ❑ What is the size of the tag field?
  
- ❑ The total number of bits in a direct-mapped cache is then  
$$2^n \times (\text{block size} + \text{tag field size} + \text{valid field size})$$

## Cache Field Sizes Examples (pp. 390)

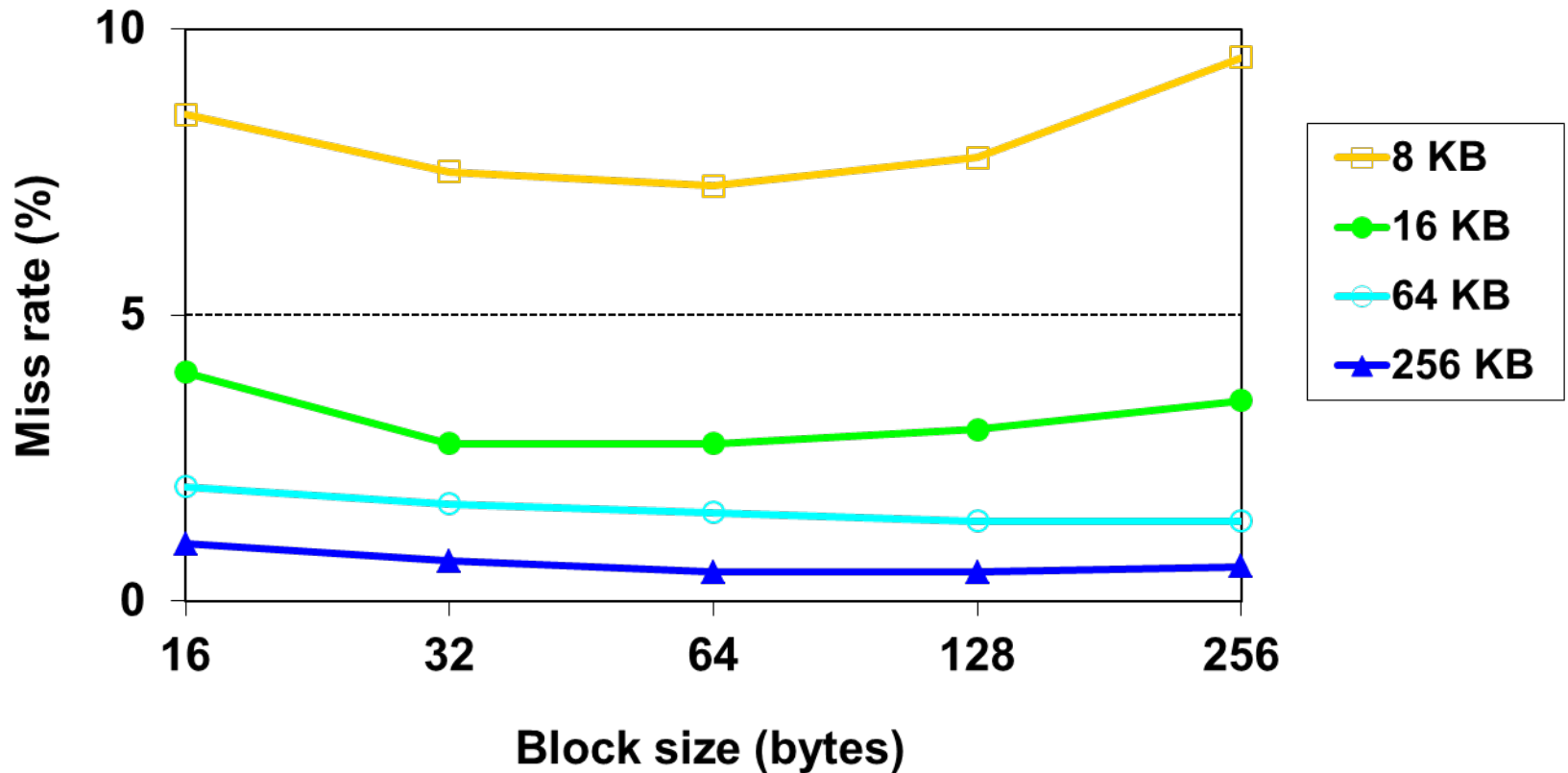
---

- ❑ How many total bits are required for a direct mapped cache with 16KiB of data and 4-word blocks, assuming a 32-bit address?
- ❑ Mapping an Address to a Multiword Cache Block

# Summary

---

# Miss Rate vs Block Size vs Cache Size



- ❑ Miss rate goes up if the block size becomes a significant fraction of the cache size because the number of blocks that can be held in the same size cache is smaller (increasing **capacity** misses)

# Handling Cache Misses (pp. 392 ~ 393)

❑ Instruction Cache Miss

❑ Data Cache Miss

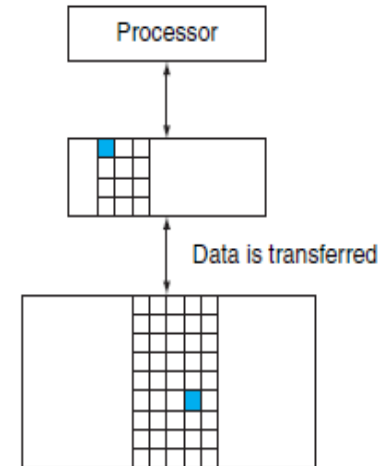
# Handling Writes

❑ Writes?

❑ Write Through

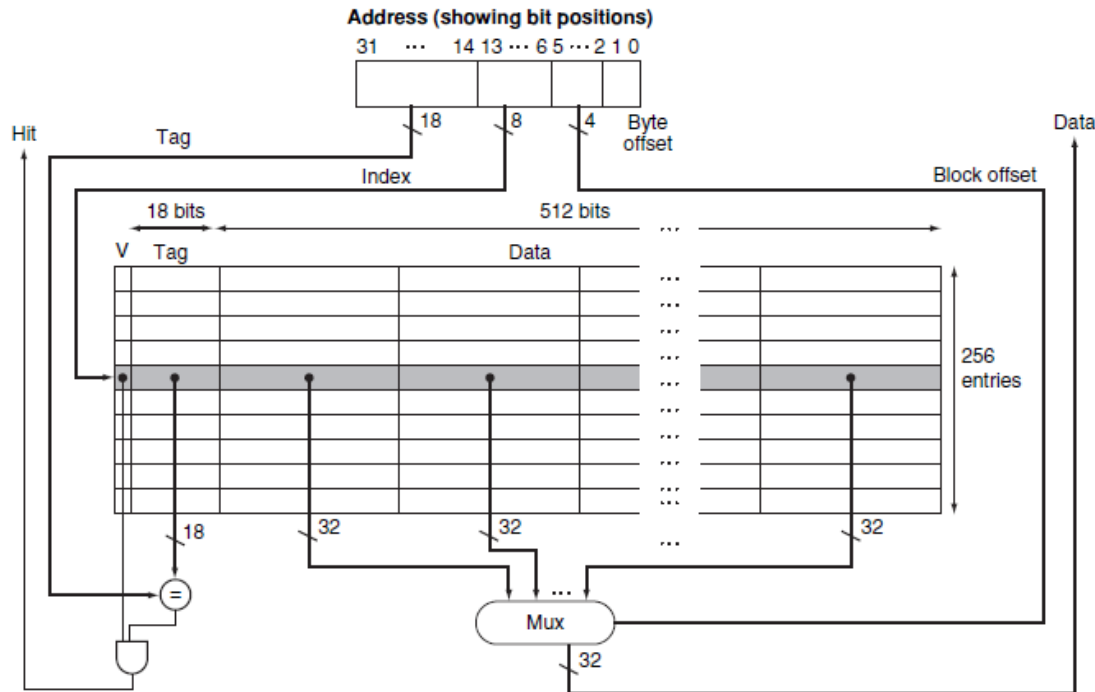
| Write buffer

❑ Write Back (copy back)



# Example Caches

## ❑ 16KiB caches in the Intrinsity FastMATH Fig. 5.12



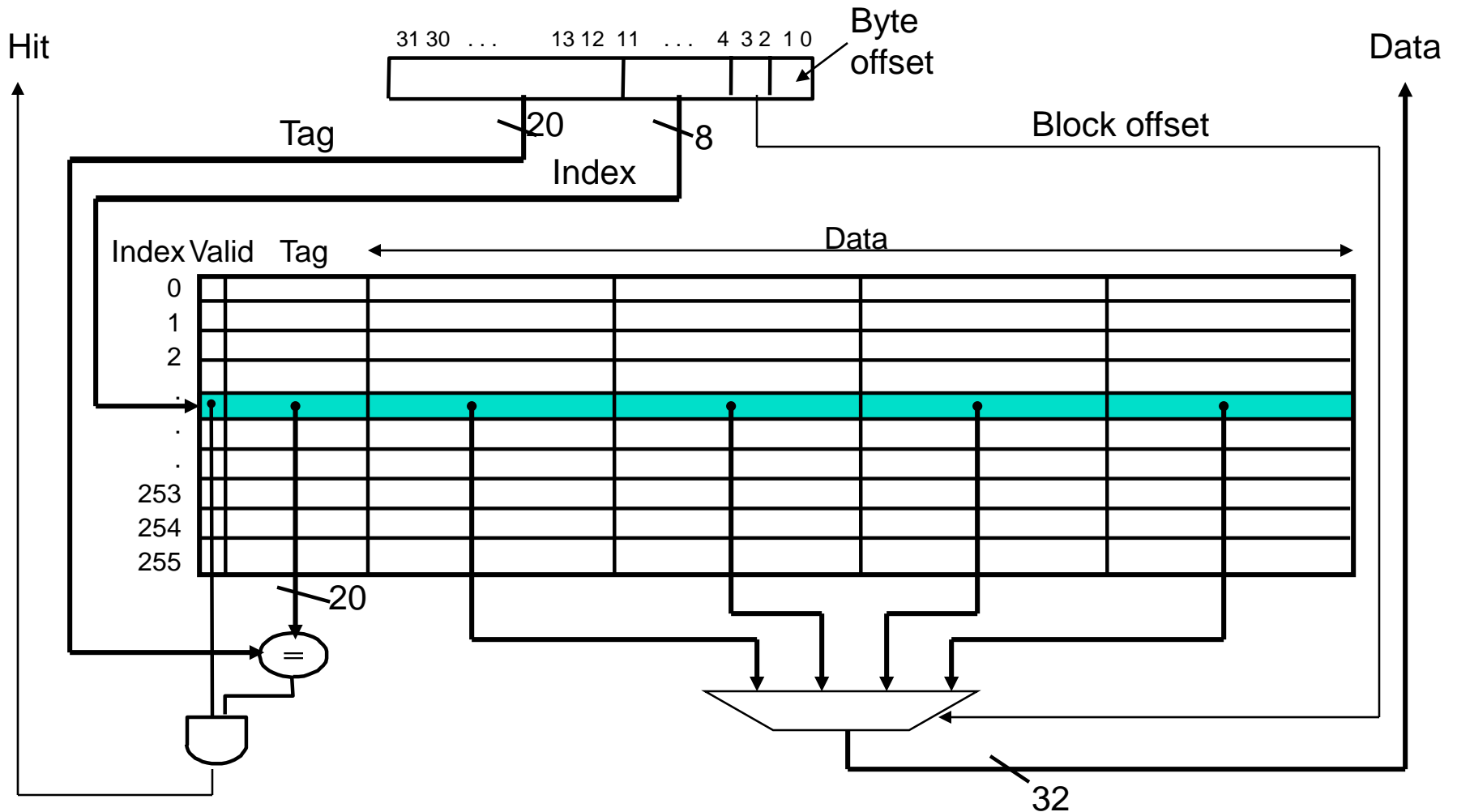
## ❑ 16KiB cache with 4-word(16-byte) blocks ?

## ❑ Split cache pp. 397

## ❑ Check yourself. pp.398

# Multiword Block Direct Mapped Cache

- Four words/block, cache size = 1K words



*What kind of locality are we taking advantage of?*



# Taking Advantage of Spatial Locality

- ❑ Let cache block hold more than one word

Start with an empty cache - all  
blocks initially marked as not valid

0 1 2 3 4 3 4 15

**0**


**1**


**2**


**3**


**4**


**3**


**4**


**15**


# Taking Advantage of Spatial Locality

## ❑ Let cache block hold more than one word

Start with an empty cache - all blocks initially marked as not valid

0 1 2 3 4 3 4 15

**0 miss**

00	Mem(1)	Mem(0)

**1 hit**

00	Mem(1)	Mem(0)

**2 miss**

00	Mem(1)	Mem(0)
00	Mem(3)	Mem(2)

**3 hit**

00	Mem(1)	Mem(0)
00	Mem(3)	Mem(2)

**4 miss**

<del>01</del>	<del>Mem(1)</del>	<del>Mem(0)</del>
00	Mem(3)	Mem(2)

**3 hit**

01	Mem(5)	Mem(4)
00	Mem(3)	Mem(2)

**4 hit**

01	Mem(5)	Mem(4)
00	Mem(3)	Mem(2)

**15 miss**

<del>11</del>	<del>Mem(5)</del>	<del>Mem(4)</del>
00	<del>Mem(3)</del>	<del>Mem(2)</del>

| 8 requests, 4 misses

# Sources of Cache Misses (pp. 459)

---

❑ **Compulsory** (cold start or process migration, first reference):

- | First access to a block, “cold” fact of life, not a whole lot you can do about it. If you are going to run “millions” of instruction, compulsory misses are insignificant
- | Solution: increase block size (increases miss penalty; very large blocks could increase miss rate)

❑ **Capacity**:

- | Cache cannot contain all blocks accessed by the program
- | Solution: increase cache size (may increase access time)

❑ **Conflict** (collision):

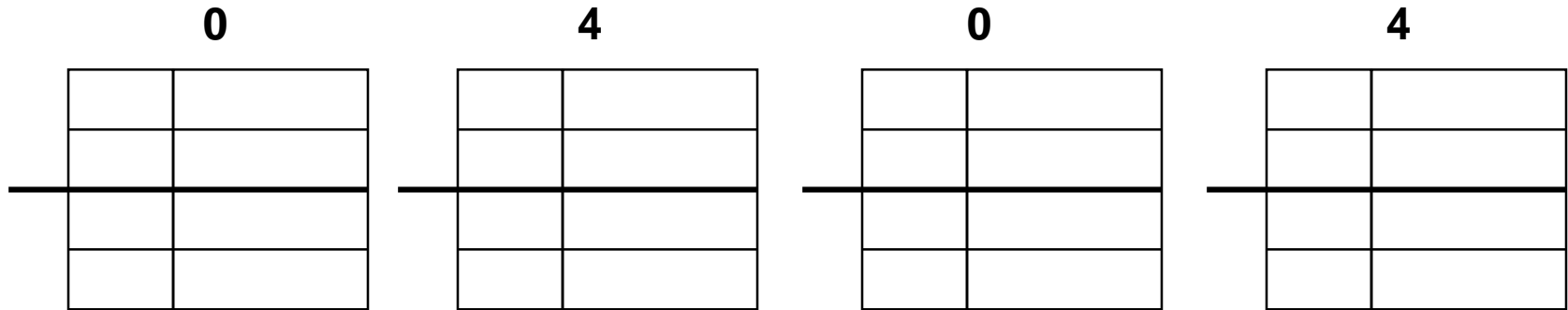
- | Multiple memory locations mapped to the same cache location
- | Solution 1: increase cache size
- | Solution 2: increase associativity (stay tuned) (may increase access time)

# Another Reference String Mapping

❑ Consider the main memory word reference string

Start with an empty cache - all  
blocks initially marked as not valid

0 4 0 4 0 4 0 4



# Another Reference String Mapping

❑ Consider the main memory word reference string

Start with an empty cache - all  
blocks initially marked as not valid

0 4 0 4 0 4 0 4

0


4


0


4


0


4


0

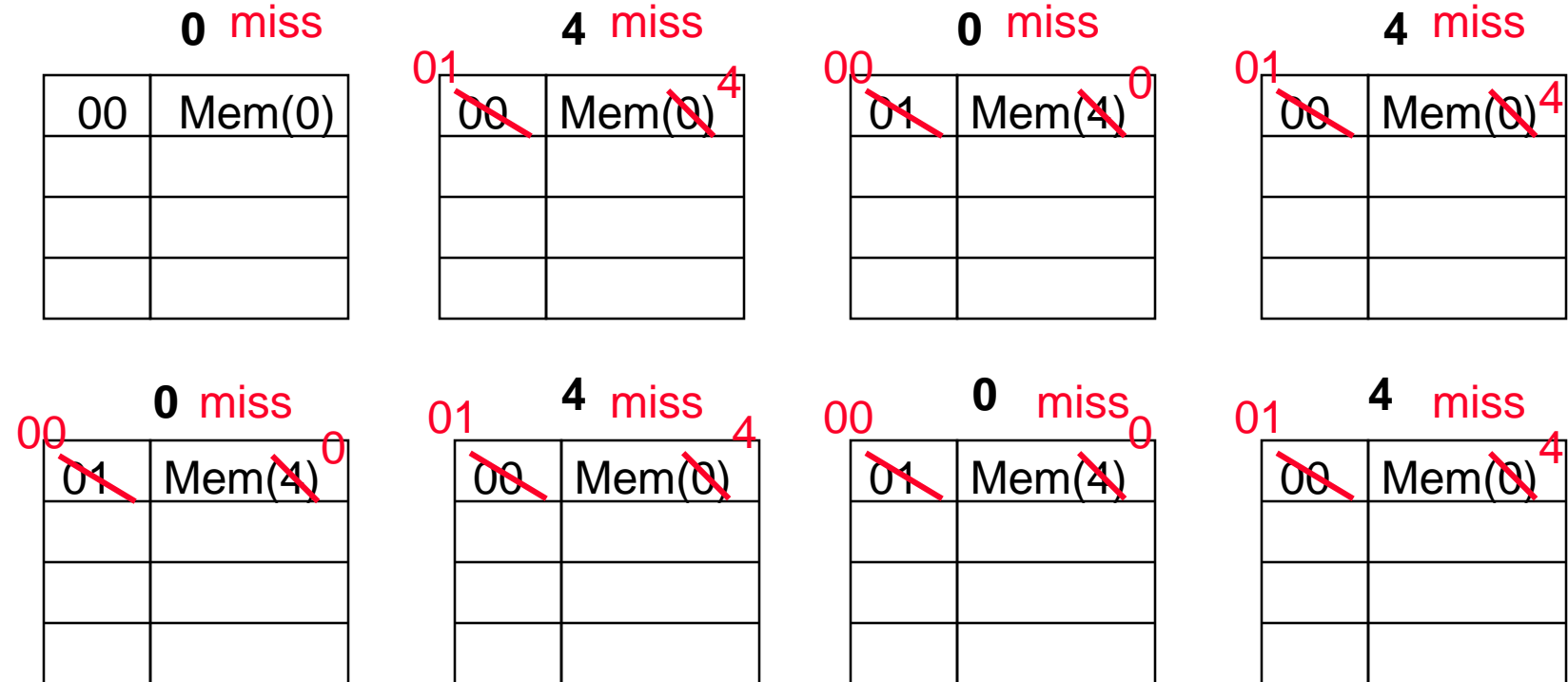

4


# Another Reference String Mapping

❑ Consider the main memory word reference string

Start with an empty cache - all blocks initially marked as not valid

0 4 0 4 0 4 0 4



| 8 requests, 8 misses

❑ Ping pong effect due to **conflict** misses - two memory locations that map into the same cache block

# Summary

---

# Measuring and Improving Cache Performance

- ❑ Example pp. 402
- ❑ Suppose our on-chip SRAM (cache) has 0.8 ns access time, but the fastest DRAM (main memory) we can get has an access time of 10ns. How high a hit rate do we need to sustain an average access time of 1ns?

Let  $h$  be the desired hit rate.

$$1 = 0.8h + (1 - h) \times (10 + 0.8) = 0.8h + 10.8 - 10.8h$$

$$10h = 9.8 \rightarrow h = 0.98$$

Hence we need a hit rate of **98%**.

- ❑ Average Memory Access Time (AMAT)

- | The average to access memory considering both hits and misses

$$\text{AMAT} = \text{Time for a hit} + \text{Miss rate} \times \text{Miss penalty}$$



# Reducing Cache Miss Rates

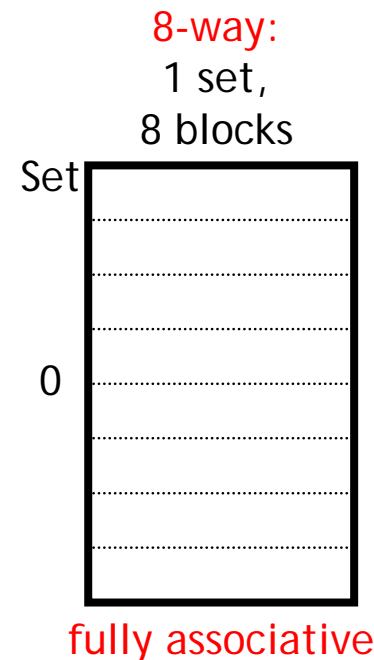
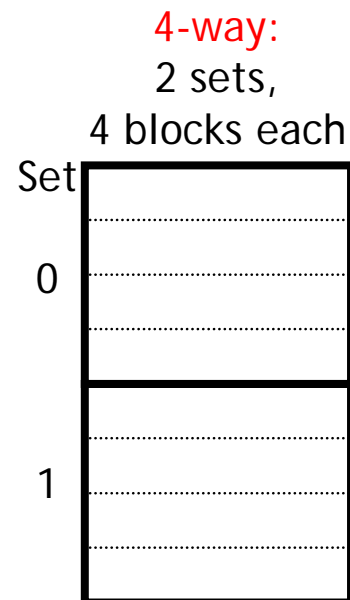
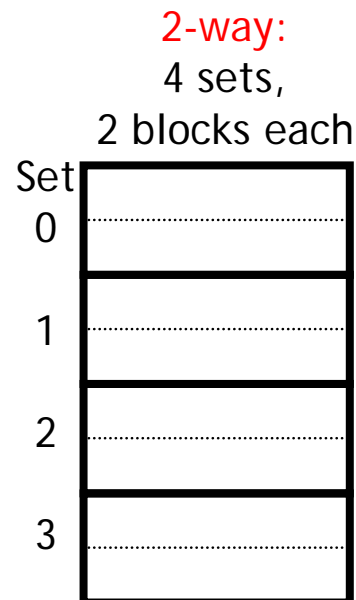
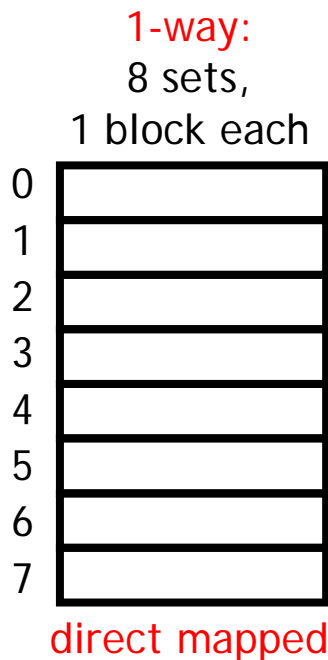
---

## 1. Allow more flexible block placement

- ❑ Figure 5.14, Figure 5.15
- ❑ In a **direct mapped cache** a memory block maps to **exactly one cache block**
- ❑ At the other extreme, could allow a memory block to be mapped to *any* cache block – **fully associative cache**
- ❑ A compromise is to divide the cache into **sets** each of which consists of *n* “ways” (***n*-way set associative**).
  - | A memory block maps to a unique set (specified by the index field) and can be placed in any way of that set (so there are *n* choices)  
 **$(\text{block address}) \bmod (\# \text{ sets in the cache})$**

# Associativity

- ❑ What if we could store data in any place in the cache?
  - | More complicated hardware = more power consumed, slower
- ❑ So we *combine* the two ideas:
  - | Each address maps to exactly one **set**
  - | Each set can store block in more than one **way**



# Set Associative Cache Example

**Cache**

Way	Set	V	Tag	Data
0	0			
	1			
1	0			
	1			

Q1: Is it there?

Compare *all* the cache **tags** in the set to the **high order 3 memory address bits** to tell if the memory block is in the cache

**Main Memory**

	000	0xx
	000	1xx
	001	0xx
	001	1xx
	010	0xx
	010	1xx
	011	0xx
	011	1xx
	100	0xx
	100	1xx
	101	0xx
	101	1xx
	110	0xx
	110	1xx
	111	0xx
	111	1xx

One word blocks  
Two low order bits define the byte in the word (32b words)

Q2: How do we find it?

Use **next 1 low order memory address bit** to determine which cache set (i.e., modulo the number of sets in the cache)

# Another Reference String Mapping

❑ Consider the main memory word reference string

Start with an empty cache - all blocks initially marked as not valid

0 4 0 4 0 4 0 4

0 miss

4 miss

0 hit

4 hit

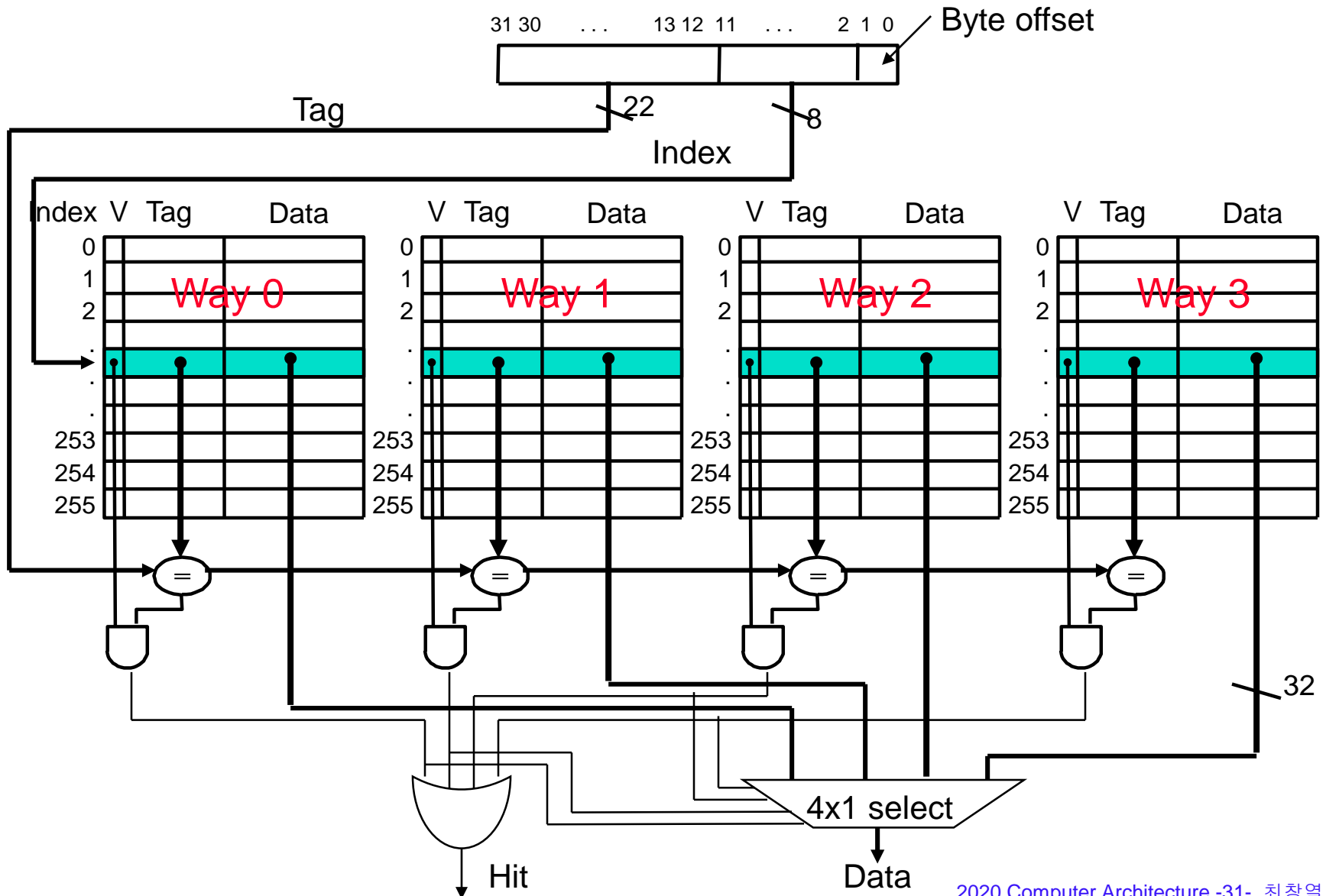
000	Mem(0)	000	Mem(0)	000	Mem(0)	000	Mem(0)
		010	Mem(4)	010	Mem(4)	010	Mem(4)

| 8 requests, 2 misses

❑ Solves the ping pong effect in a direct mapped cache due to **conflict** misses since now two memory locations that map into the same cache set can co-exist!

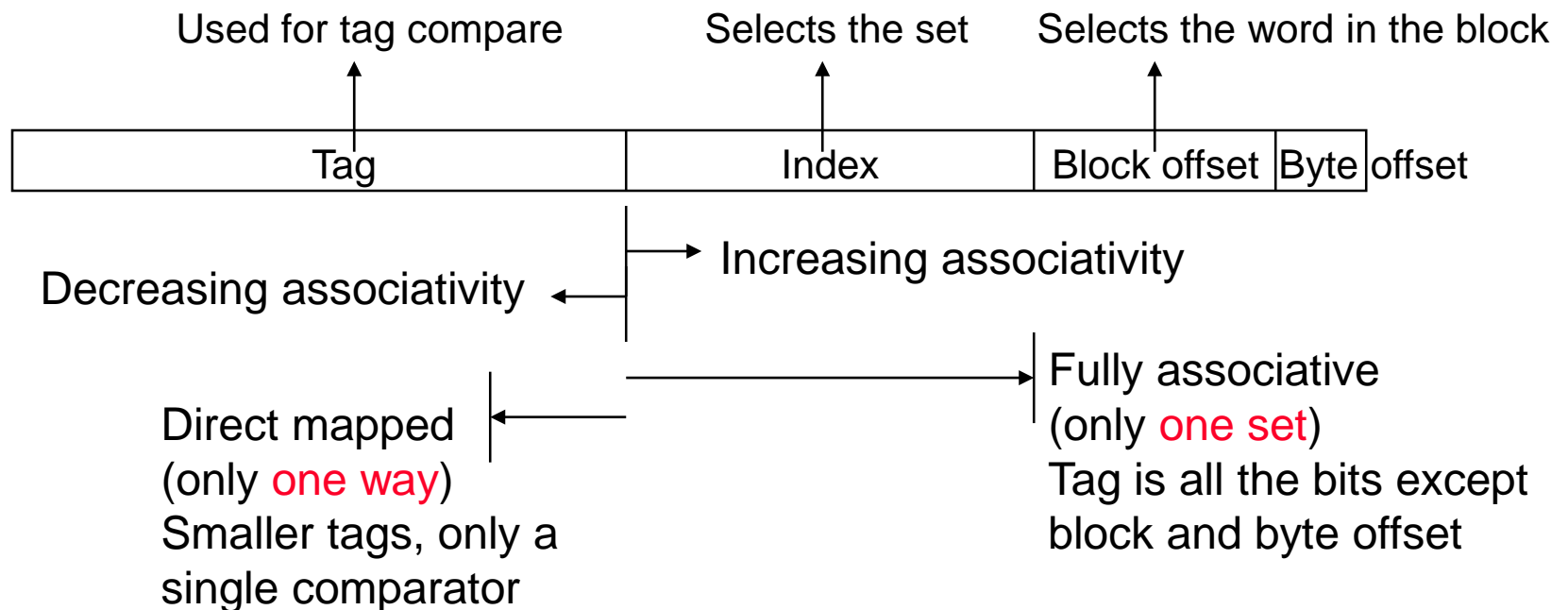
# Four-Way Set Associative Cache (Fig. 5.18)

❑  $2^8 = 256$  sets each with four ways (each with one block)



# Range of Set Associative Caches

- For a fixed size cache, each increase by a factor of two in associativity doubles the number of blocks per set (i.e., the number of ways) and halves the number of sets – decreases the size of the index by 1 bit and increases the size of the tag by 1 bit



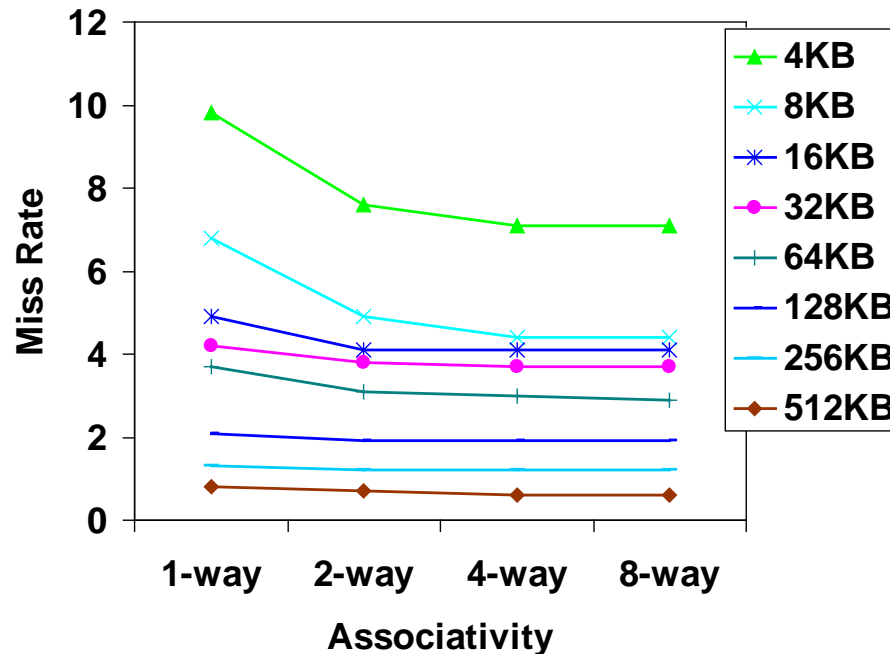
# Costs of Set Associative Caches

---

- ❑ When a miss occurs, which way's block do we pick for replacement?
  - | **Least Recently Used (LRU)**: the block replaced is the one that has been unused for the longest time
    - Must have hardware to keep track of when each way's block was used relative to the other blocks in the set
    - For 2-way set associative, takes **one bit per set** → set the bit when a block is referenced (and reset the other way's bit)
- ❑ N-way set associative cache costs
  - | N comparators (delay and area)
  - | MUX delay (set selection) before data is available
  - | Data available **after** set selection (and Hit/Miss decision). In a direct mapped cache, the cache block is available **before** the Hit/Miss decision
- ❑ Size of Tags versus Set Associativity Example : pp. 409

# Benefits of Set Associative Caches

- ❑ The choice of direct mapped or set associative depends on the cost of a miss versus the cost of implementation



- ❑ Largest gains are in going from direct mapped to 2-way (20%+ reduction in miss rate)



# Reducing the Miss Penalty (pp. 410)

## 2. Use multiple levels of caches

- ❑ With advancing technology have more than enough room on the die for bigger L1 caches *or* for a second level of caches – normally a **unified** L2 cache (i.e., it holds both instructions and data) and in some cases even a unified L3 cache
- ❑ **Example** Performance of Multilevel Caches

# Multilevel Cache Design Considerations

---

- ❑ Design considerations for L1 and L2 caches are very different
  - | Primary cache should focus on **minimizing hit time** in support of a shorter clock cycle
    - Smaller with smaller block sizes
  - | Secondary cache(s) should focus on **reducing miss rate** to reduce the penalty of long main memory access times
    - Larger with larger block sizes
    - Higher levels of associativity
- ❑ The miss penalty of the L1 cache is significantly reduced by the presence of an L2 cache – so it can be smaller (i.e., faster) but have a higher miss rate
- ❑ For the L2 cache, hit time is less important than miss rate
- ❑ Check yourself.

# Summary

## ❑ Common Framework for Memory Hierarchy

## ❑ Caches in the ARM Cortex-A8 and Intel Core i7- 920

Characteristic	ARM Cortex-A8	Intel Nehalem
L1 cache organization	Split instruction and data caches	Split instruction and data caches
L1 cache size	32 KiB each for instructions/data	32 KiB each for instructions/data per core
L1 cache associativity	4-way (I), 4-way (D) set associative	4-way (I), 8-way (D) set associative
L1 replacement	Random	Approximated LRU
L1 block size	64 bytes	64 bytes
L1 write policy	Write-back, Write-allocate(?)	Write-back, No-write-allocate
L1 hit time (load-use)	1 clock cycle	4 clock cycles, pipelined
L2 cache organization	Unified (instruction and data)	Unified (instruction and data) per core
L2 cache size	128 KiB to 1 MiB	256 KiB (0.25 MiB)
L2 cache associativity	8-way set associative	8-way set associative
L2 replacement	Random(?)	Approximated LRU
L2 block size	64 bytes	64 bytes
L2 write policy	Write-back, Write-allocate (?)	Write-back, Write-allocate
L2 hit time	11 clock cycles	10 clock cycles
L3 cache organization	–	Unified (instruction and data)
L3 cache size	–	8 MiB, shared
L3 cache associativity	–	16-way set associative
L3 replacement	–	Approximated LRU
L3 block size	–	64 bytes
L3 write policy	–	Write-back, Write-allocate
L3 hit time	–	35 clock cycles