

# Lecture 12

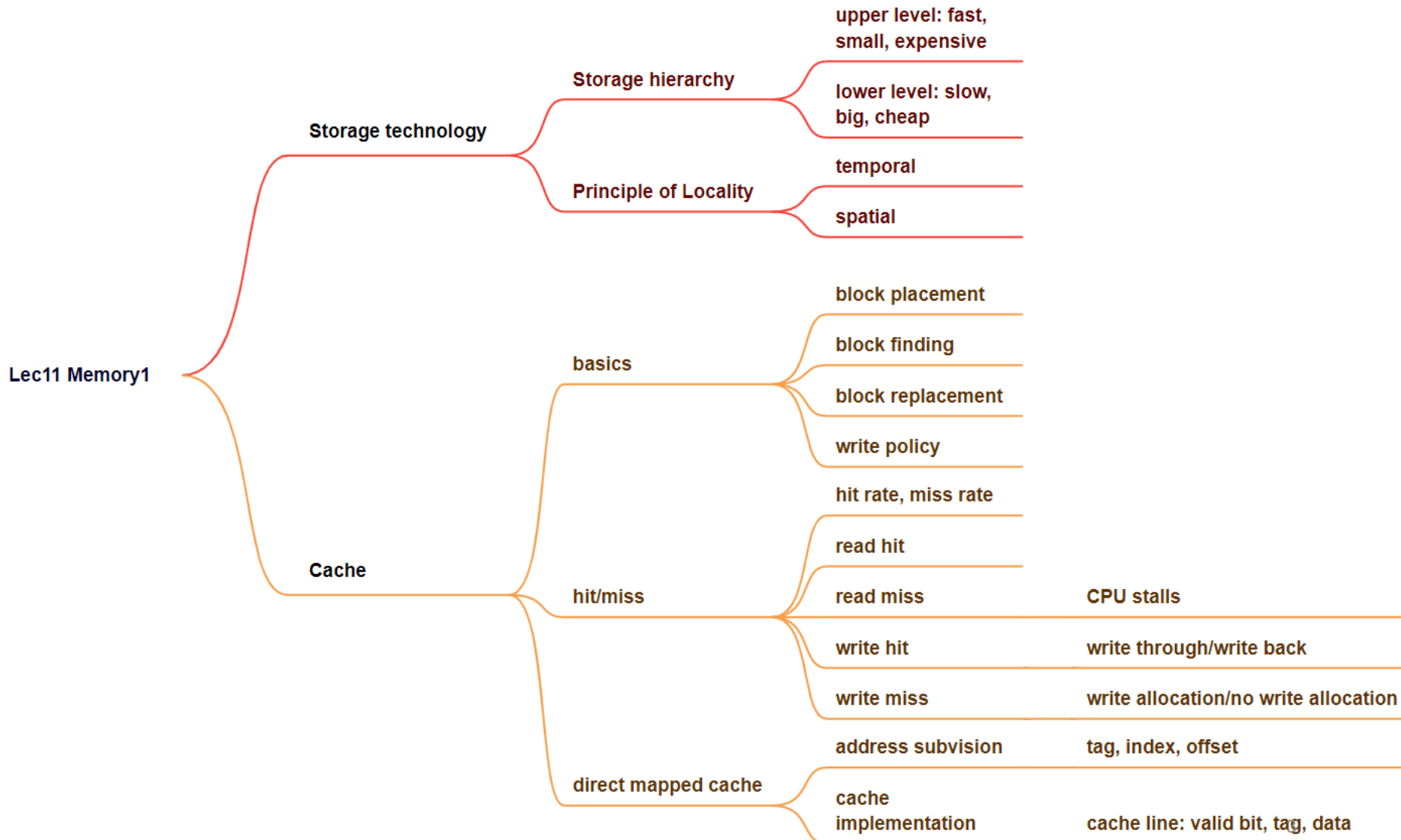
## Memory Hierarchy (2)

CS202 2023 Spring

# Today's Agenda

- Recap
  - Memory technology
  - Memory hierarchy
  - The basics of caches
- Context
  - Measuring and improving cache performance
  - Virtual memory
- Reading: Textbook 5.4, 5.7

# Recap



# Outline

- **Measuring cache performance**
- Improving performance – Associative cache
  - Fully associative
  - n-ways Set associative
- Improving performance – Multilevel Caches
- Virtual Memory

# Measuring Cache Performance

- Components of CPU time
  - Program execution cycles
    - Includes cache hit time
  - Memory stall cycles
    - Mainly from cache misses
- With simplifying assumptions:

Memory stall cycles

$$= \frac{\text{Memory accesses}}{\text{Program}} \times \text{Miss rate} \times \text{Miss penalty}$$

$$= \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Misses}}{\text{Instruction}} \times \text{Miss penalty}$$

# 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:  $2\% \times 100 = 2$
  - D-cache:  $36\% \times 4\% \times 100 = 1.44$
- Effective CPI =  $2 + 2 + 1.44 = 5.44$ 
  - Ideal CPU is  $5.44/2 = 2.72$  times faster

# Average Access Time

- Hit time is also important for performance
- Average memory access time (**AMAT**)

$$\text{AMAT} = \text{Hit time} + \text{Miss rate} \times \text{Miss penalty}$$

- Example
  - CPU with 1ns clock, hit time = 1 cycle, miss penalty = 20 cycles, l-cache miss rate = 5%
  - $\text{AMAT} = 1 + 0.05 \times 20 = 2\text{ns}$ 
    - 2 cycles per instruction

# Performance Summary

- 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

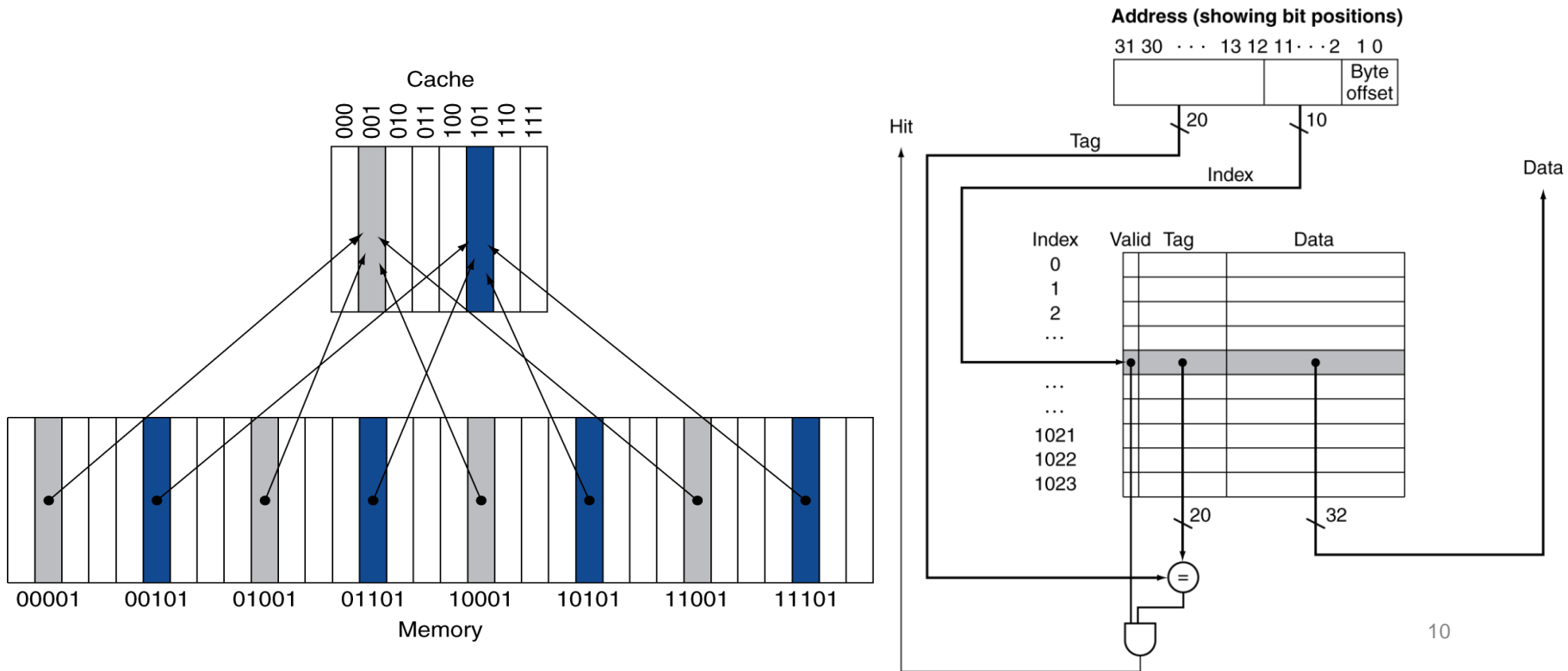


# Outline

- Measuring cache performance
- **Improving performance – Associative cache**
  - Fully associative
  - n-ways Set associative
- Improving performance – Multilevel Caches
- Virtual Memory

# Recall: Direct Mapped Cache

- Direct mapped cache:
  - Location determined by address
  - One data in memory is mapped to only one location in cache
  - Capacity of cache is not fully exploited
  - Miss rate is high



# Recall: Direct Mapped Cache

16	10 000	Miss	000
3	00 011	Miss	011
16	10 000	Hit	000

Index	V	Tag	Data
<b>000</b>	<b>Y</b>	<b>10</b>	<b>Mem[10000]</b>
001	N		
010	Y	11	Mem[11010]
<b>011</b>	<b>Y</b>	<b>00</b>	<b>Mem[00011]</b>
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		



18	10 010	Miss	010
16	10 000	Hit	000

Index	V	Tag	Data
000	Y	10	Mem[10000]
001	N		
<b>010</b>	<b>Y</b>	<b>10</b>	<b>Mem[10010]</b>
011	Y	00	Mem[00011]
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

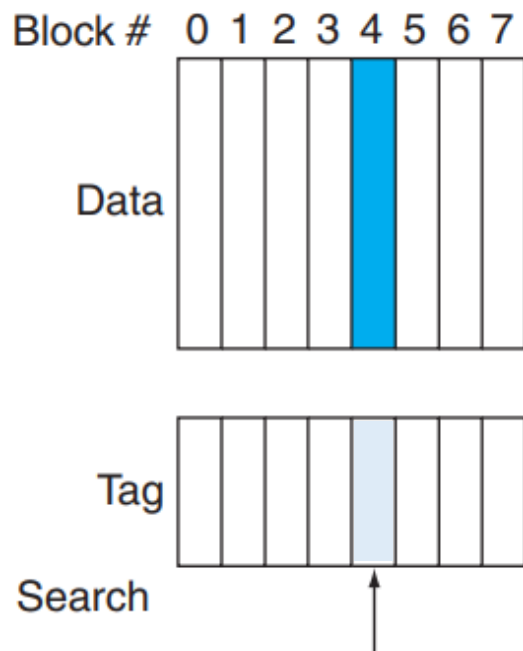
# Associative Caches

- Fully associative
  - Allow a given block to go in any cache entry
  - Requires all entries to be searched at once
  - Comparator per entry (expensive)
- $n$ -way set associative
  - Each set contains  $n$  entries
  - Block number determines which set
    - (Block number) modulo (#Sets in cache)
  - Search all entries in a given set at once
  - $n$  comparators (less expensive)

# Associative Cache Example

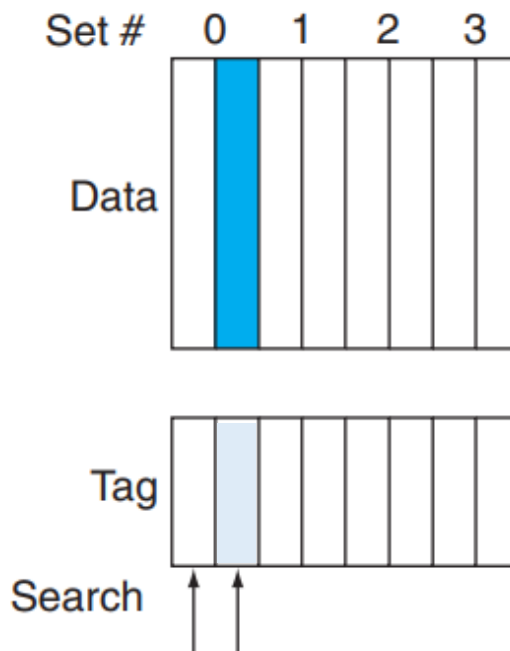
- Example: Placement of a block whose address is 12:

$$12 \bmod 8 = \text{block 4}$$



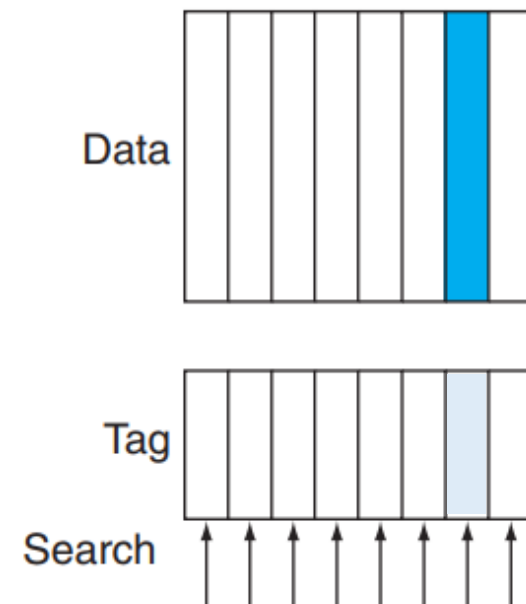
**Direct mapped**

$$12 \bmod 4 = \text{set 0, can be in either way of the set}$$



**Set associative**

can appear in any of the eight cache blocks



**Fully associative**

# Spectrum of Associativity

- An eight-block cache configured as direct mapped, two-way set associative, four-way set associative, and fully associative.

**One-way set associative**

(direct mapped)

Block	way,0	
	Tag	Data
0		
1		
2		
3		
4		
5		
6		
7		

**Two-way set associative**

Set	way,0		way,1	
	Tag	Data	Tag	Data
0				
1				
2				
3				

**Four-way set associative**

Set	way,0		way,1		way,2		way,3	
	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	Tag	Data

# Associativity Example

- 4-block caches, 1byte/block
  - Direct mapped, 2-way set associative, fully associative
  - Block access sequence: 0, 8, 0, 6, 8
- Direct mapped
  - $\text{index} = (\text{Block address}) \bmod (\text{\#Blocks})$  needs 2 bit for index

Block address	Cache index	Hit/miss	Cache content after access			
			way 0 block 0	way 0 block 1	way 0 block 3	way 0 block 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 (cont.)

- 2-way set associative

Block access sequence: 0, 8, 0, 6, 8

- index = (Block address) mod (**#Sets**)

needs 1 bit for index

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

- Fully associative (no index)

needs 0 bit for index

Block address		Hit/miss	Cache content after access			
			way <sub>0</sub>	way <sub>1</sub>	way <sub>2</sub>	way <sub>3</sub>
0		miss	Mem[0]			
8		miss	Mem[0]	Mem[8]		
0		hit	Mem[0]	Mem[8]		
6		miss	Mem[0]	Mem[8]	Mem[6]	
8		hit	Mem[0]	Mem[8]	Mem[6]	



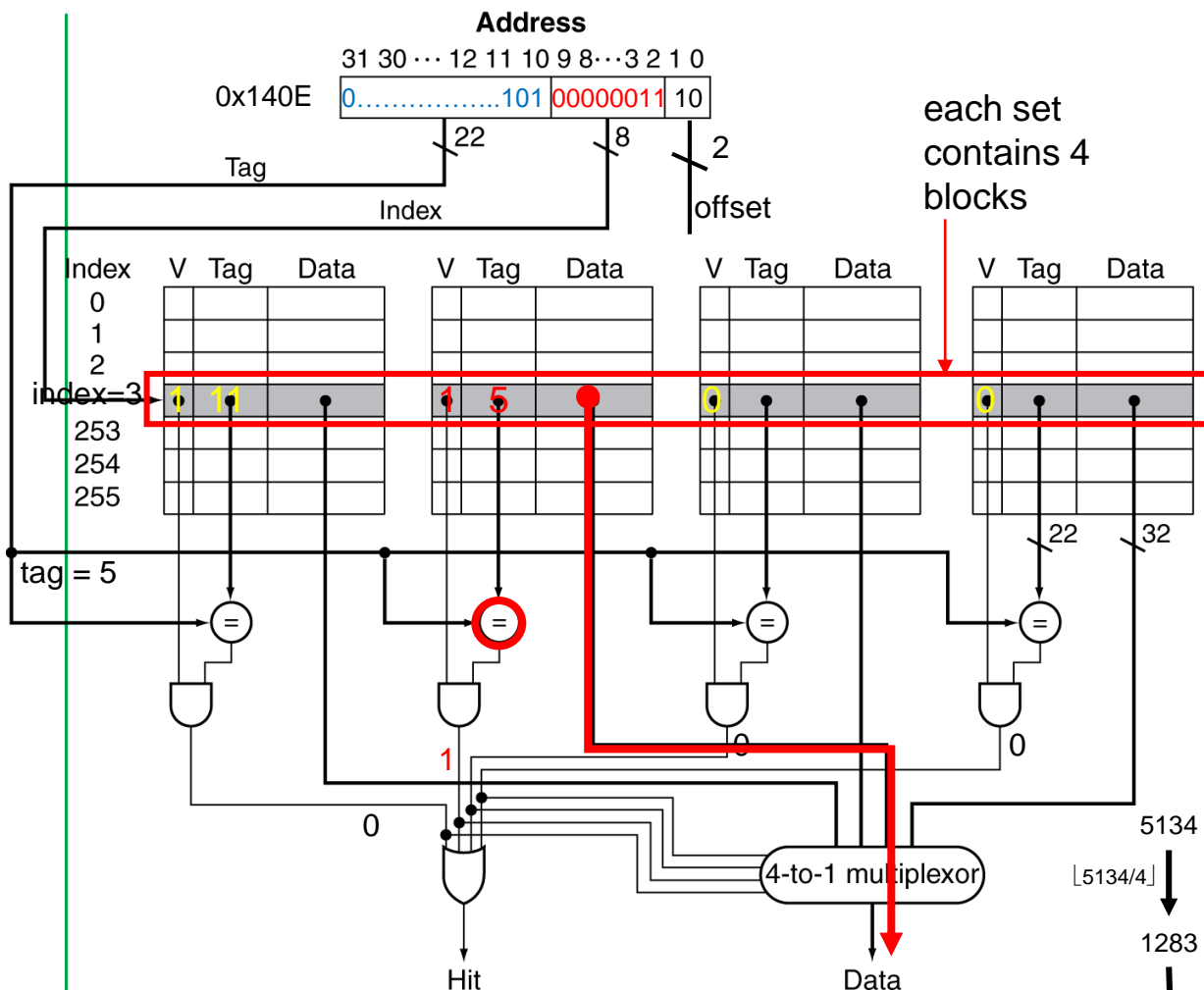
# How Much Associativity

- Increased associativity decreases miss rate
  - But with diminishing returns
- Miss rate simulation of a system with 64KB D-cache, 16-word blocks, SPEC2000
  - 1-way: 10.3%
  - 2-way: 8.6%
  - 4-way: 8.3%
  - 8-way: 8.1%

# Address Subdivision

- 4K blocks cache, 4-word/block, 32-bit address
- $\text{offset} = \log_2 16 = 4\text{bits} \rightarrow 28\text{bits for index} + \text{tag for all types}$
- direct mapped (1-way set associative)
  - same number of sets as blocks  $\rightarrow 4096$  blocks(sets)
  - $\text{index} = \log_2 4096 = 12$  bits
  - $\text{tag} = 32 - 4 - 12 = 16$  bits
- 2-way set associative
  - 2 blocks/set  $\rightarrow 2048$  sets
  - $\text{index} = \log_2 2048 = 11$  bits
  - $\text{tag} = 32 - 4 - 11 = 17$  bits
- 4-way set associative
  - 4 blocks/set  $\rightarrow 1024$  sets
  - $\text{index} = \log_2 1024 = 10$  bits
  - $\text{tag} = 32 - 4 - 10 = 18$  bits
- fully associative
  - just 1 set  $\rightarrow$  no index
  - $\text{tag} = 32 - 4 = 28$

# Set Associative Cache Organization



4KB cache, 1word/block  
To what **set** number does address **0x140E** map?

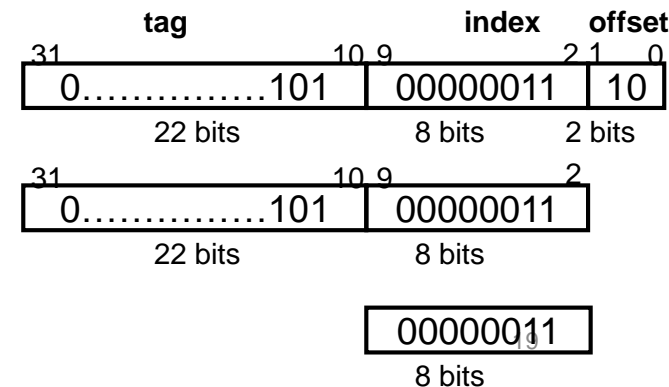
**offset:**  $\log_2(4\text{byte}) = 2\text{bits}$

**#blocks:**  $4\text{KB}/4\text{B} = 1\text{K blocks}$

**#sets:**  $1\text{K blocks}/4 \text{ way} = 256 \text{ sets}$

**index:**  $\log_2(256\text{sets}) = 8\text{bits}$

$0x140E_{\text{hex}} = \underbrace{101}_{\text{tag}} \underbrace{00000011}_{\text{index}} \underbrace{10}_{\text{offset}}$



Increasing associativity  
shrinks index, expands tag

# Replacement Policy

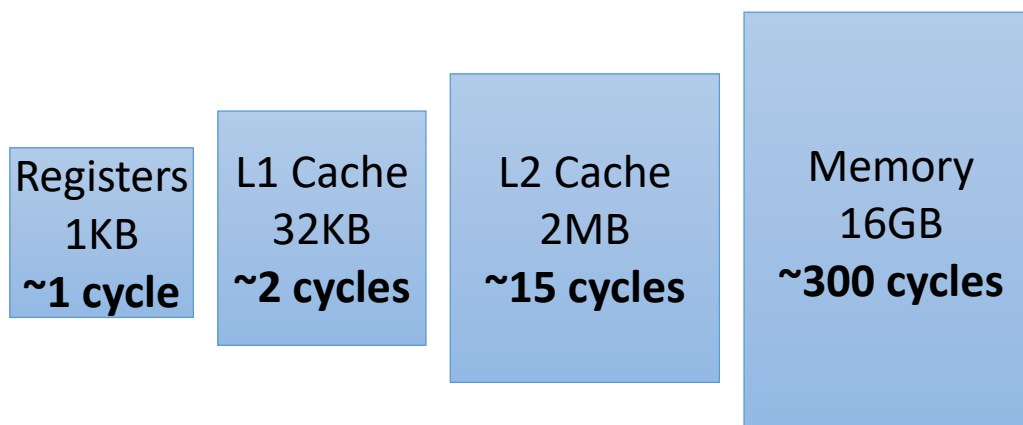
- Direct mapped: no choice
- Set associative
  - Prefer non-valid entry, if there is one
  - Otherwise, choose among entries in the set
- Least-recently used (LRU)
  - Choose the one unused for the longest time
    - Simple for 2-way, manageable for 4-way, too hard beyond that
- Random
  - Gives approximately the same performance as LRU for high associativity

# Outline

- Measuring cache performance
- Improving performance – Associative cache
  - Fully associative
  - n-ways Set associative
- **Improving performance – Multilevel Caches**
- Virtual Memory

# Multilevel 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 L-2 cache misses
- Some high-end systems include L-3 cache



# Multilevel Cache Example

- Given
  - CPU base CPI = 1, clock rate = 4GHz
  - Primary cache Miss rate/instruction = 2%
  - Main memory access time = 100ns
- Solution: With just primary cache
  - Miss penalty =  $100\text{ns} / 0.25\text{ns} = 400$  cycles
  - Effective CPI = Base CPI + miss penalty/instruction  
 $= 1 + 2\% \times 400 = 9$

Registers

L1 Cache  
**1 cycle**

**miss 2%**

Memory  
**400 cycles**

# Example (cont.)

- Given: After adding L-2 cache
  - L2 cache Access time = 5ns
  - **L2 global miss rate**/instruction = 0.5%
- Solution:
  - L-1 miss (2%) first need to access the L2 cache
    - Penalty = 5ns/0.25ns = 20 cycles
  - L-1 miss with L-2 also miss (0.5%)
    - Extra penalty = 400 cycles
  - **Effective CPI =  $1 + 2\% \times 20 + 0.5\% \times 400 = 3.4$**
  - Speedup =  $9/3.4 = 2.6$
- Global miss rate
  - The fraction of references that miss in all levels
- Local miss rate
  - The fraction of references to one level of a cache that miss
  - **L2 local miss rate**/instruction =  $0.5\%/2\% = 25\%$
  - **Effect CPI =  $1 + 2\% \times (20 + 25\% \times 400) = 3.4$**

Registers

L1 Cache  
**1 cycle**

**global miss 2%**

L2 Cache  
**20 cycles**

**global miss 0.5%**

Memory  
**400 cycles**



# Multilevel Cache Considerations

- Primary cache
  - Focus on minimal hit time
- L-2 cache
  - Focus on low miss rate to avoid main memory access
  - Hit time has less overall impact
- Results
  - L-1 cache usually smaller than a single cache
  - L-1 block size smaller than L-2 block size

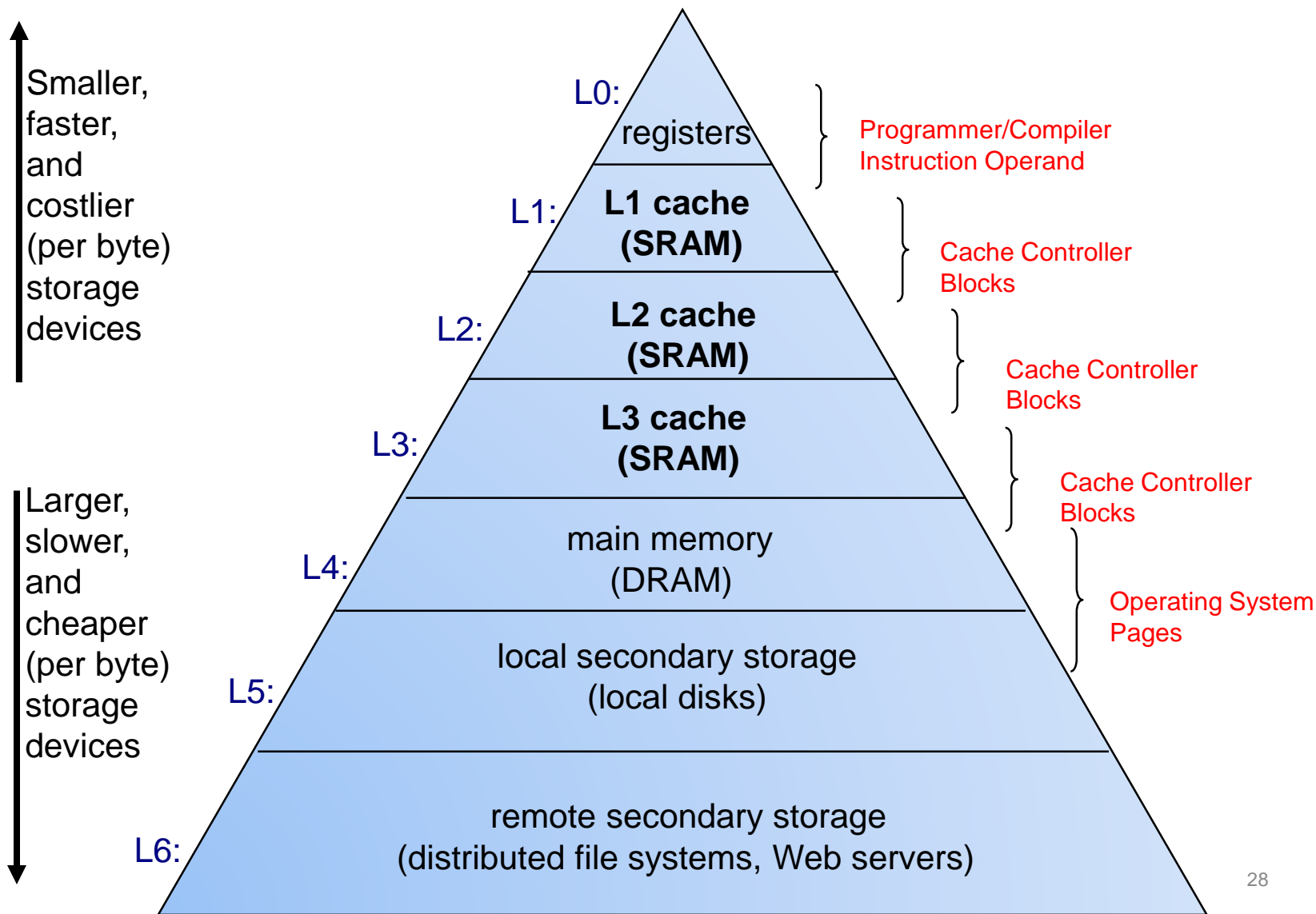
# Outline

- Measuring cache performance
- Improving performance – Associative cache
  - Fully associative
  - n-ways Set associative
- Improving performance – Multilevel Caches
- **Virtual Memory**

# Virtual Memory

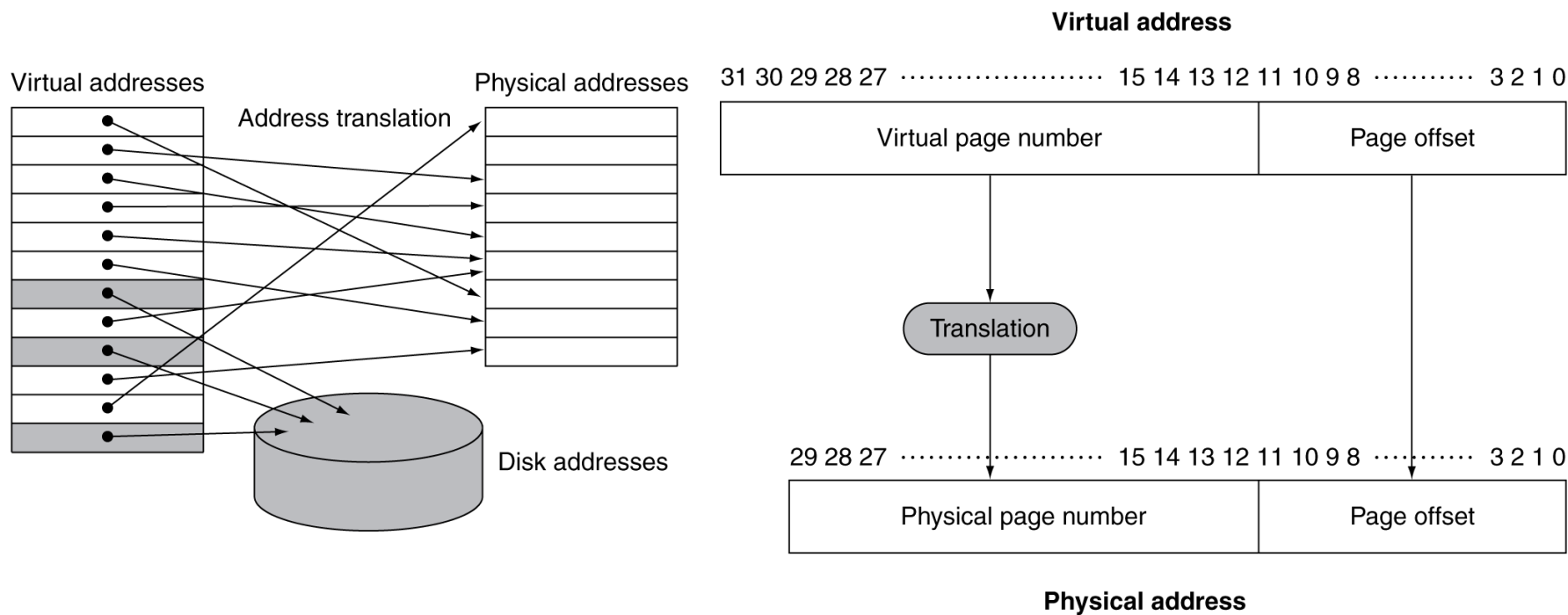
- Use main memory as a “cache” for secondary (disk) storage
  - Managed jointly by CPU hardware and the operating system (OS)
- Programs share main memory
  - Each gets a private virtual address space holding its frequently used code and data
  - Protected from other programs
- CPU and OS translate virtual addresses to physical addresses
  - VM “block” is called a **page**
  - VM translation “miss” is called a **page fault**

# Recall: Memory Hierarchy



# Address Translation

- Fixed-size pages (e.g., 4K)



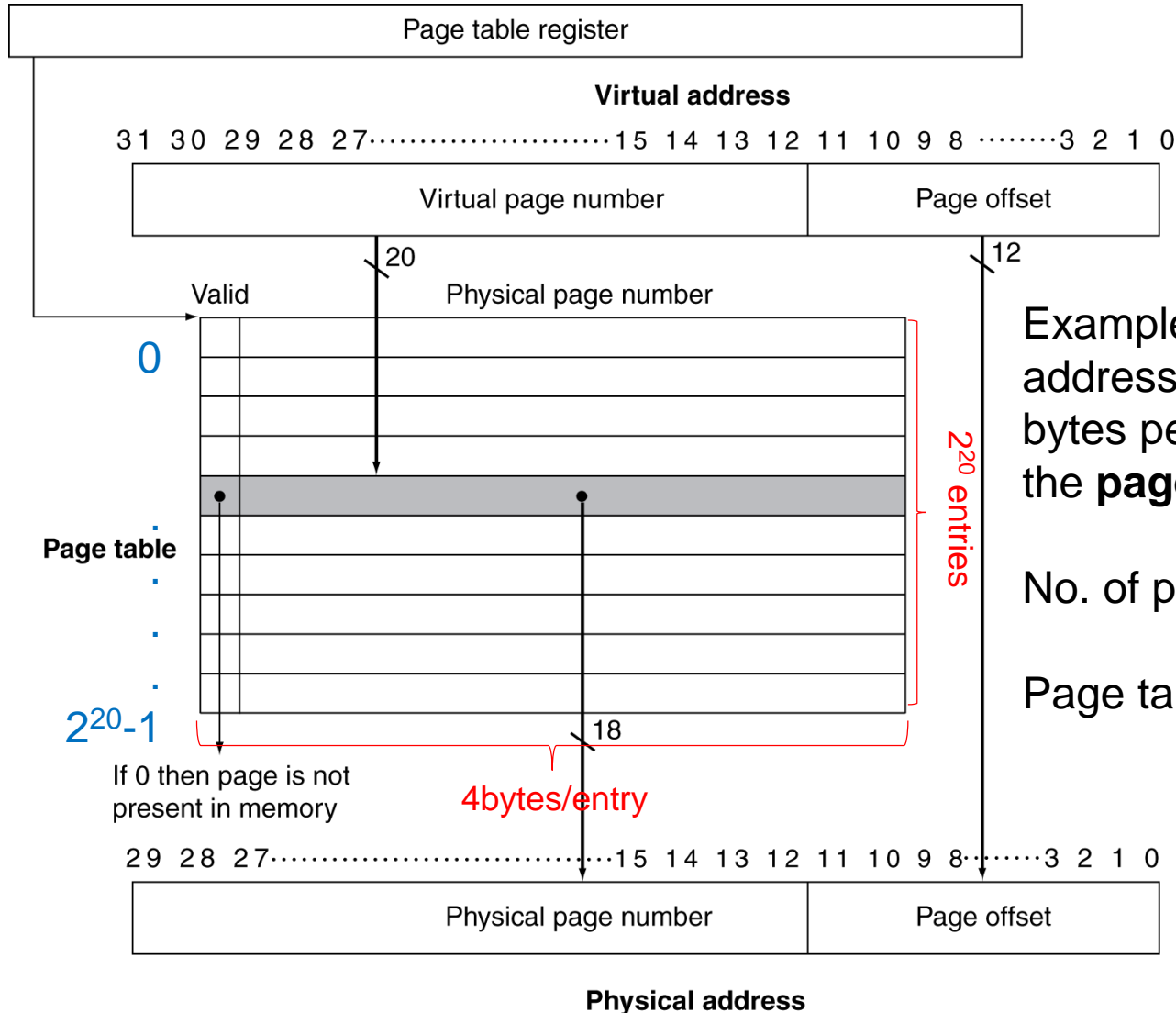
# Page Tables

- Stores placement information
  - Array of page table entries, indexed by virtual page number
  - Page table register in CPU points to page table in physical memory
- If page is present in memory
  - PTE stores the physical page number
  - Plus other status bits (referenced bit, dirty bit, ...)
- If page is not present
  - PTE can refer to location in swap space on disk

# Page Fault Penalty

- On page fault, the page must be fetched from disk
  - Takes millions of clock cycles
  - Handled by OS code
- Try to minimize page fault rate
  - Fully associative placement
  - Smart replacement algorithms

# Translation Using a Page Table



Example: With a 32-bit virtual address, 4 KB page size, and 4 bytes per page table entry, what's the **page table size**?

No. of page table entries  
 $= 2^{32} / 2^{12} = 2^{20}$

Page table size  
 $= 2^{20} \times 4 = 4\text{MB}$



# Mapping Pages to Storage

VA: 0x6C8<sub>hex</sub>

Virtual page  
number offset

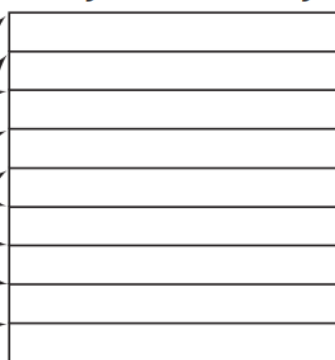
6	
---	--

Page table

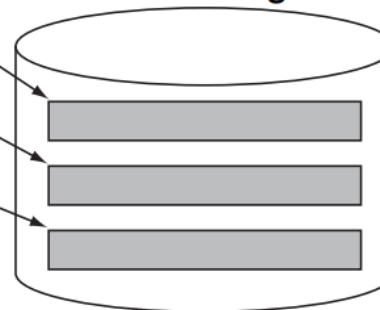
Physical page num  
or disk address

Valid		
0		-
1	1	5
2	1	2
3	1	6
4	1	7
5	0	-
6	1	3
7	1	8
...	0	-
...	1	0
...	1	4
...	0	-
15	1	1

Physical memory



Disk storage



Example: **16 virtual pages**,  
**256B/page**, virtual address:  
12 bits.

To what physical address  
does virtual address 0x6C8  
map to?

page offset:  $\log_2(256B) = 8$   
bits

0x6C8<sub>hex</sub> = 011011001000  
                                vpn      offset

Virtual page number = 1  
Physical page number = 5  
Physical address =  
1111001000 = 0x3C8<sub>hex</sub>

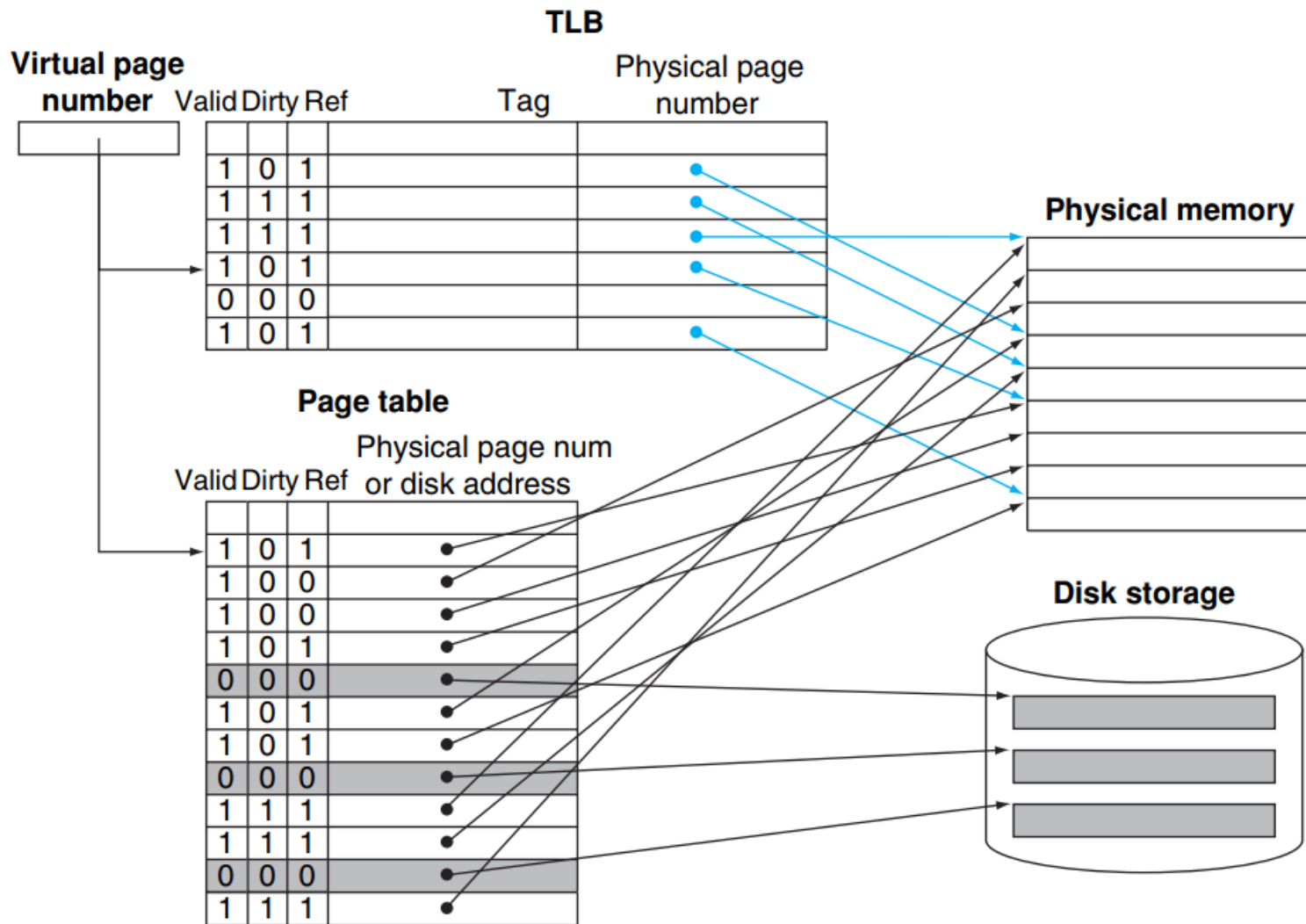
# Replacement and Writes

- To reduce page fault rate, prefer least-recently used (LRU) replacement
  - **Reference bit** (aka **use bit**) in PTE set to 1 on access to page
  - Periodically cleared to 0 by OS
  - A page with reference bit = 0 has not been used recently
- Disk writes take millions of cycles
  - Block at once, not individual locations
  - Write through is impractical
  - Use write-back
  - **Dirty bit** in PTE set when page is written

# Fast Translation Using a TLB

- Problems of Page Table
  - Access to page table is too slow
    - First access the PTE (one main memory access)
    - Then the actual memory access for data
  - Page table is too big
- But access to page tables has good locality
  - So use a fast cache of PTEs within the CPU
  - Called a **Translation Look-aside Buffer (TLB)**
  - Typical: 16–512 PTEs, 0.5–1 cycle for hit, 10–100 cycles for miss, 0.01%–1% miss rate
  - Misses could be handled by hardware or software

# Fast Translation Using a TLB

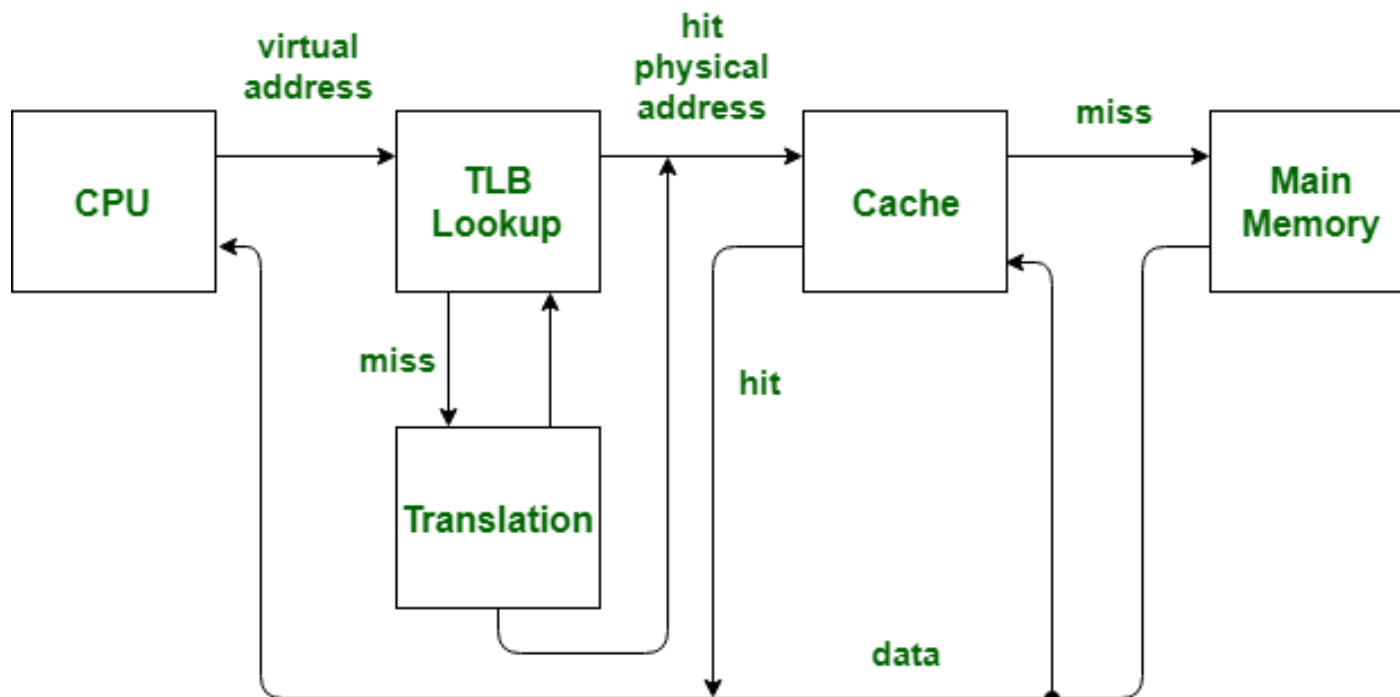


# TLB Misses

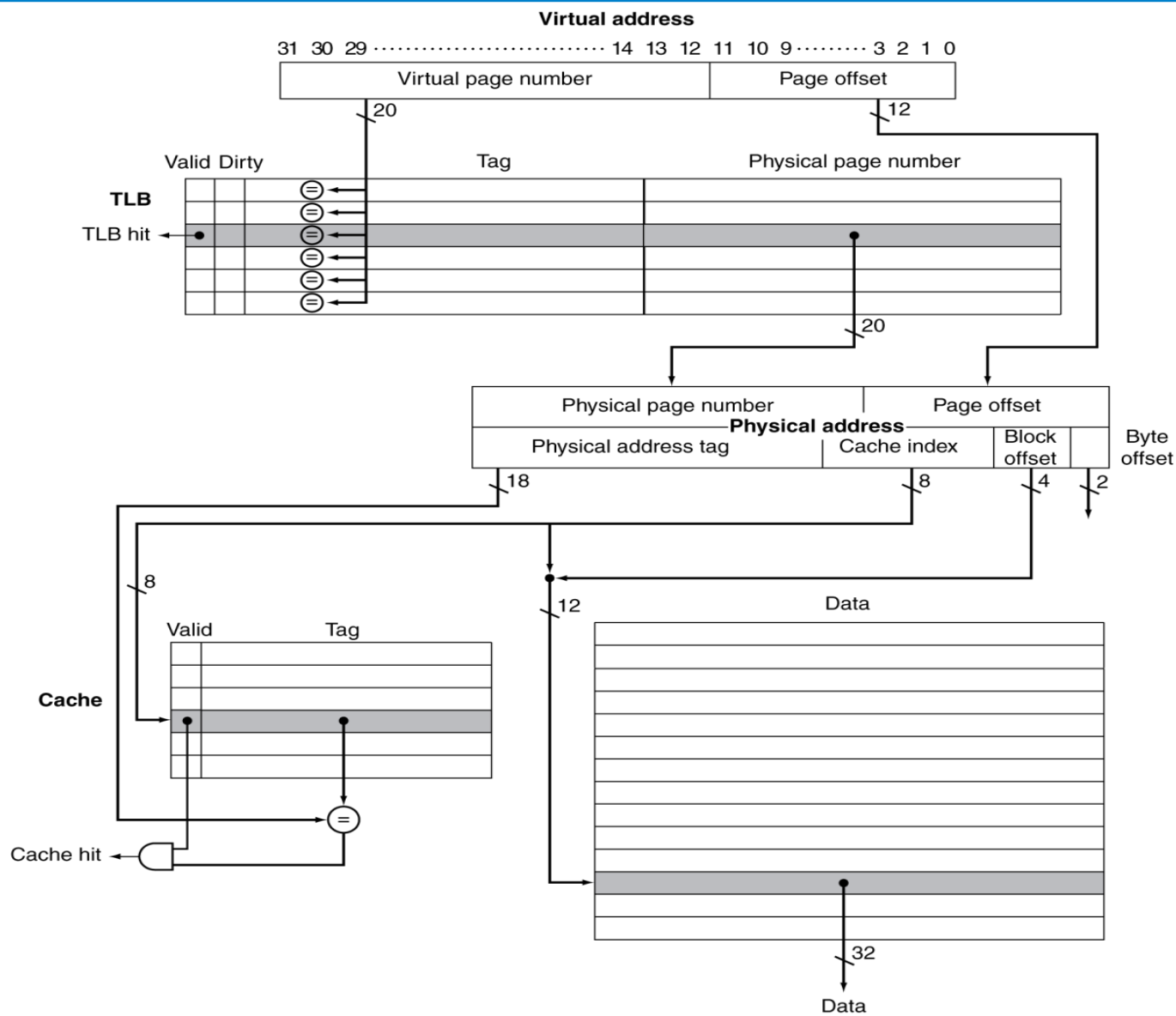
- If page is in memory
  - Load the PTE from memory and retry
  - Could be handled in hardware
    - Can get complex for more complicated page table structures
  - Or in software
    - Raise a special exception, with optimized handler
- If page is not in memory (page fault)
  - OS handles fetching the page and updating the page table (software)
  - Then restart the faulting instruction

# Making Address Translation Practical

- In VM, memory acts like a cache for disk
  - Page table maps virtual page numbers to physical frames
  - Use a page table cache for recent translation
    - => Translation Lookaside Buffer (TLB)



# TLB and Cache Interaction



# Memory Protection

- Different tasks can share parts of their virtual address spaces
  - But need to protect against errant access
  - Requires OS assistance
- Hardware support for OS protection
  - Privileged supervisor mode (aka kernel mode)
  - Privileged instructions
  - Page tables and other state information only accessible in supervisor mode
  - System call exception (e.g., syscall in MIPS)



# Check Yourself

- Match the definitions between left and right

L1 cache — A cache for a cache

L2 cache ~~—~~ A cache for disks

Main memory ~~—~~ A cache for a main memory

TLB — A cache for page table entries