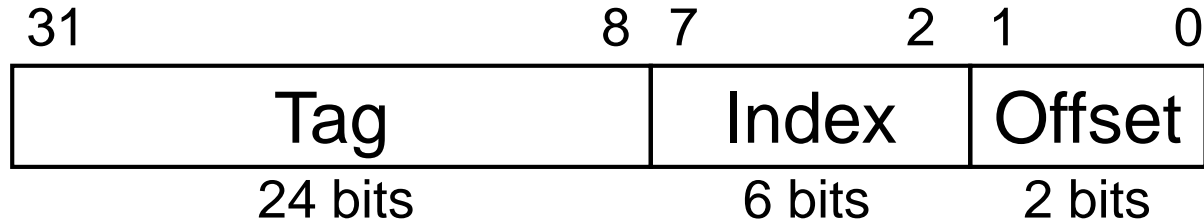


Example Cache Field Sizes

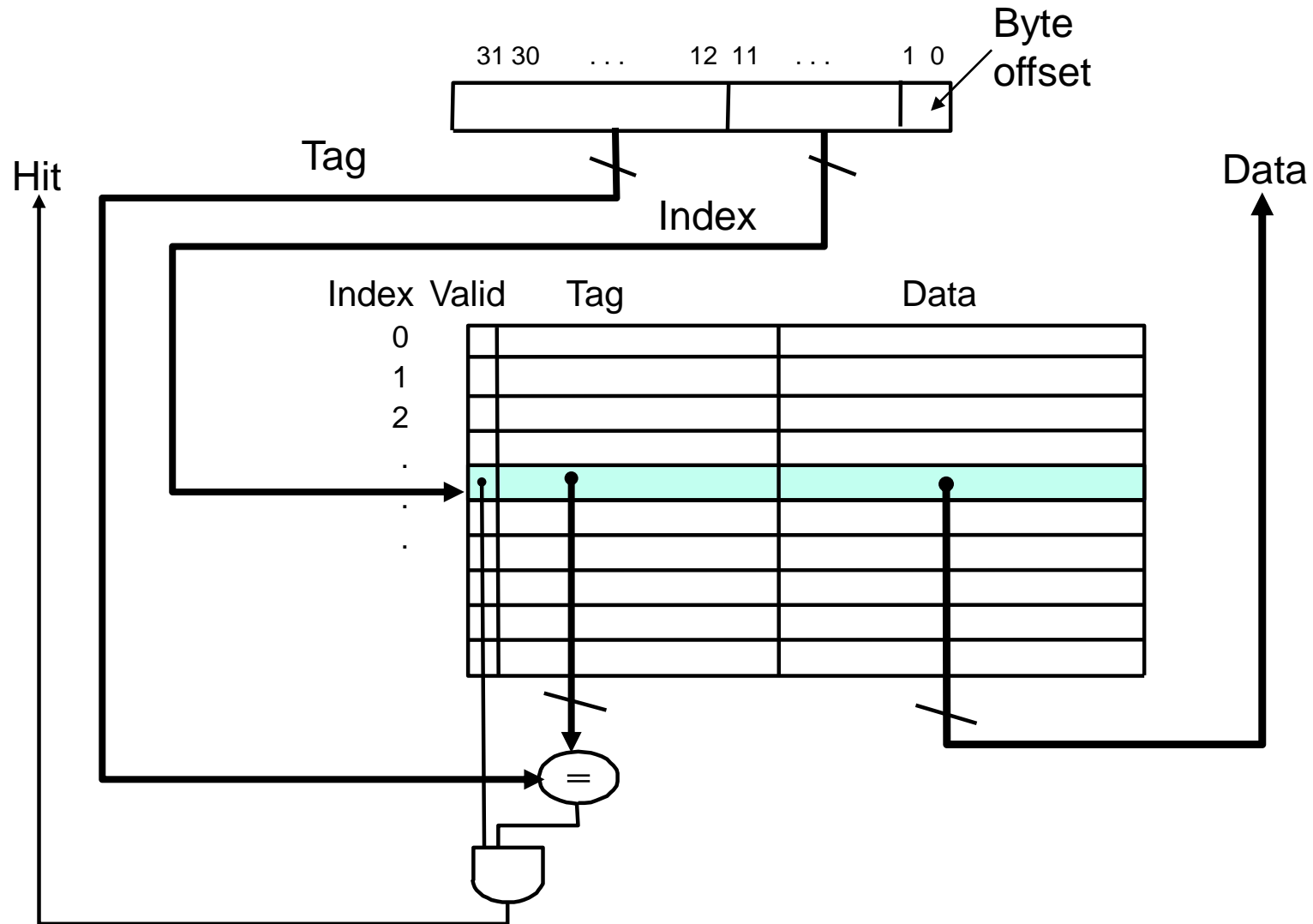
- Usually 32-bit byte addresses
- 64 blocks, 1 Word/block Direct Mapped Cache
 - 2 bits to access each byte (**Byte Offset**)
 - How many bits to index all blocks? i.e., $(2^6 = 64) \Rightarrow$ 6 bits for **index**
 - **Tag** bits: $32 - (2 + 6) = 24$



- How many bits will be in cache?
 - Blocks times bits needed for Valid, Tag and Data $\Rightarrow 64 \times (32+24+1)$

MIPS Direct Mapped Cache Example

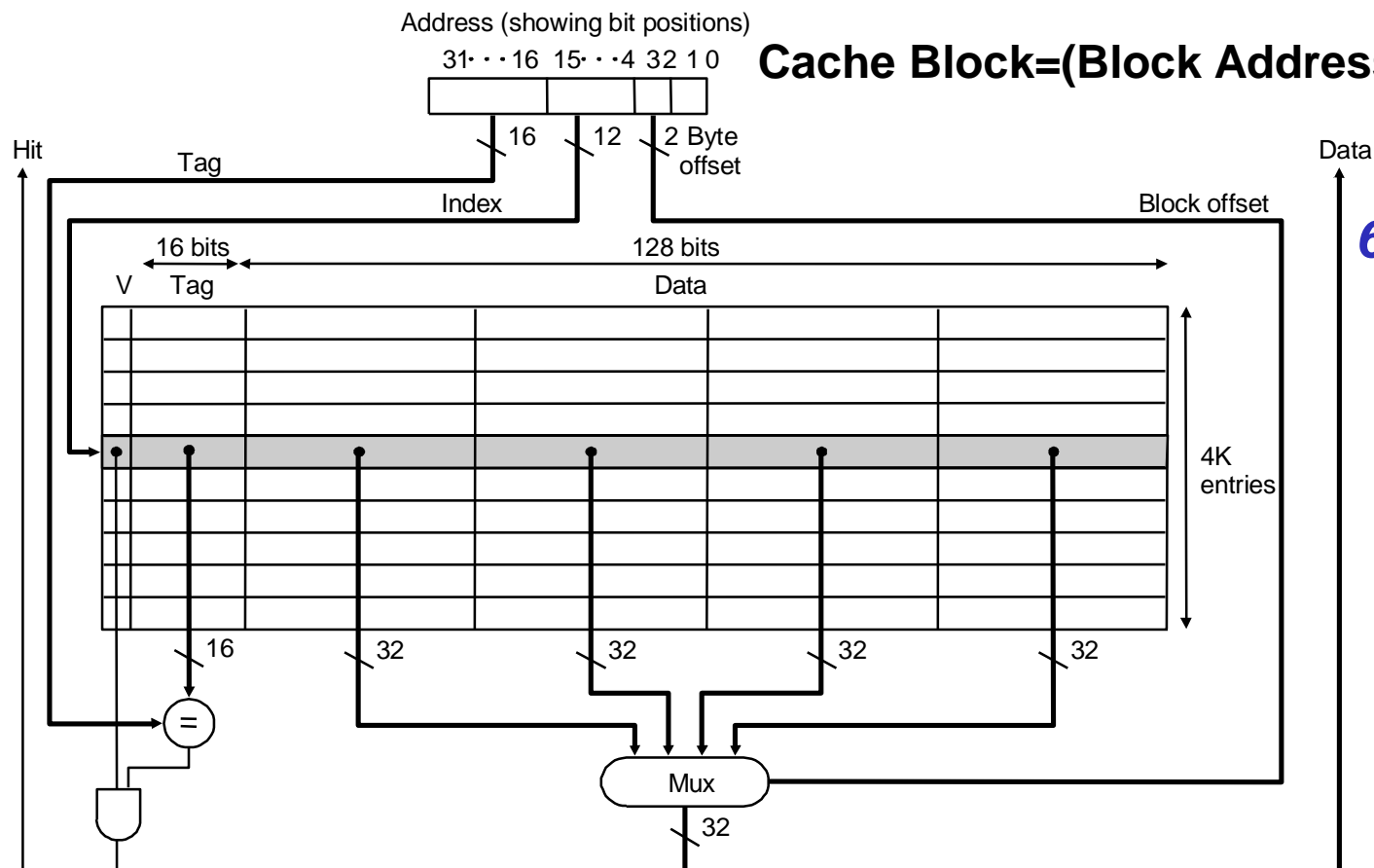
- **1- word** blocks, cache size = 4KB



What kind of locality are we taking advantage of?

Multiword Block Direct Mapped Cache

- Cache Block (sometimes called a cache *line*)
 - A cache entry that has in its own cache tag
 - Previous example used 4-byte blocks
- Larger cache blocks take advantage of *spatial* locality



**64KB Direct-Mapped Cache
with 16B Blocks**

Figure 5.9 from text.

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

00	Mem(1)	Mem(0)
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

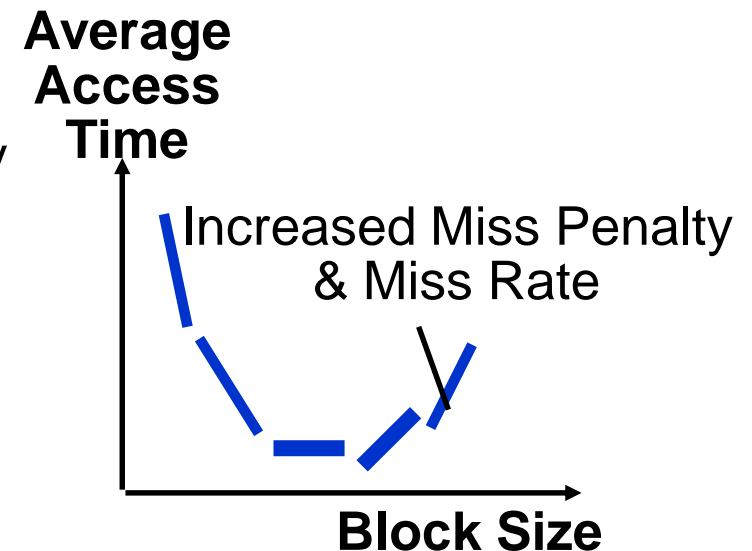
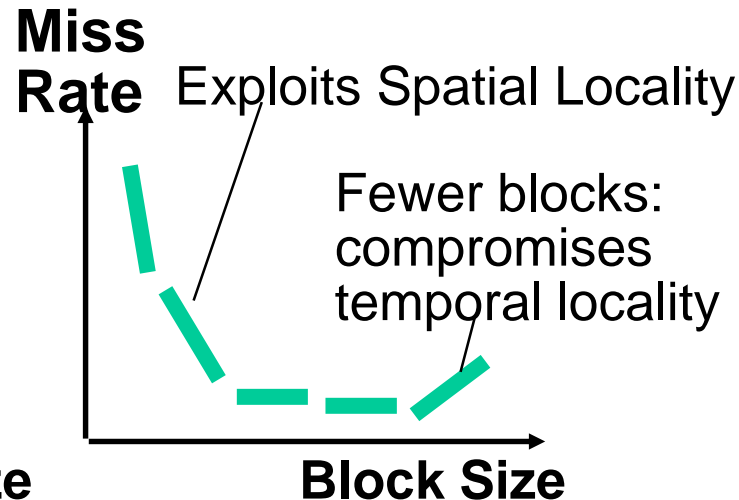
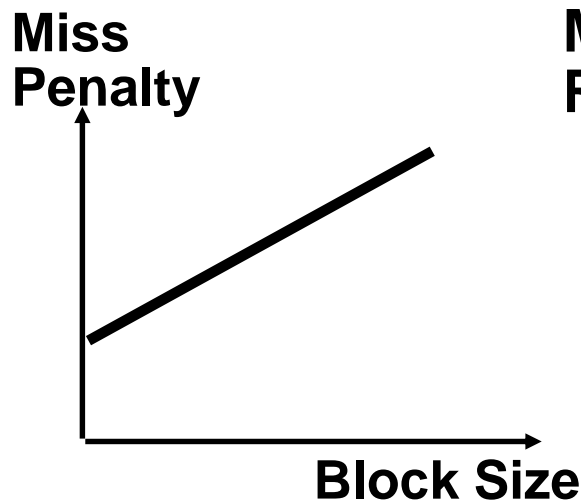
01	Mem(5)	Mem(4)
00	Mem(3)	Mem(2)

| **8 requests, 4 misses**

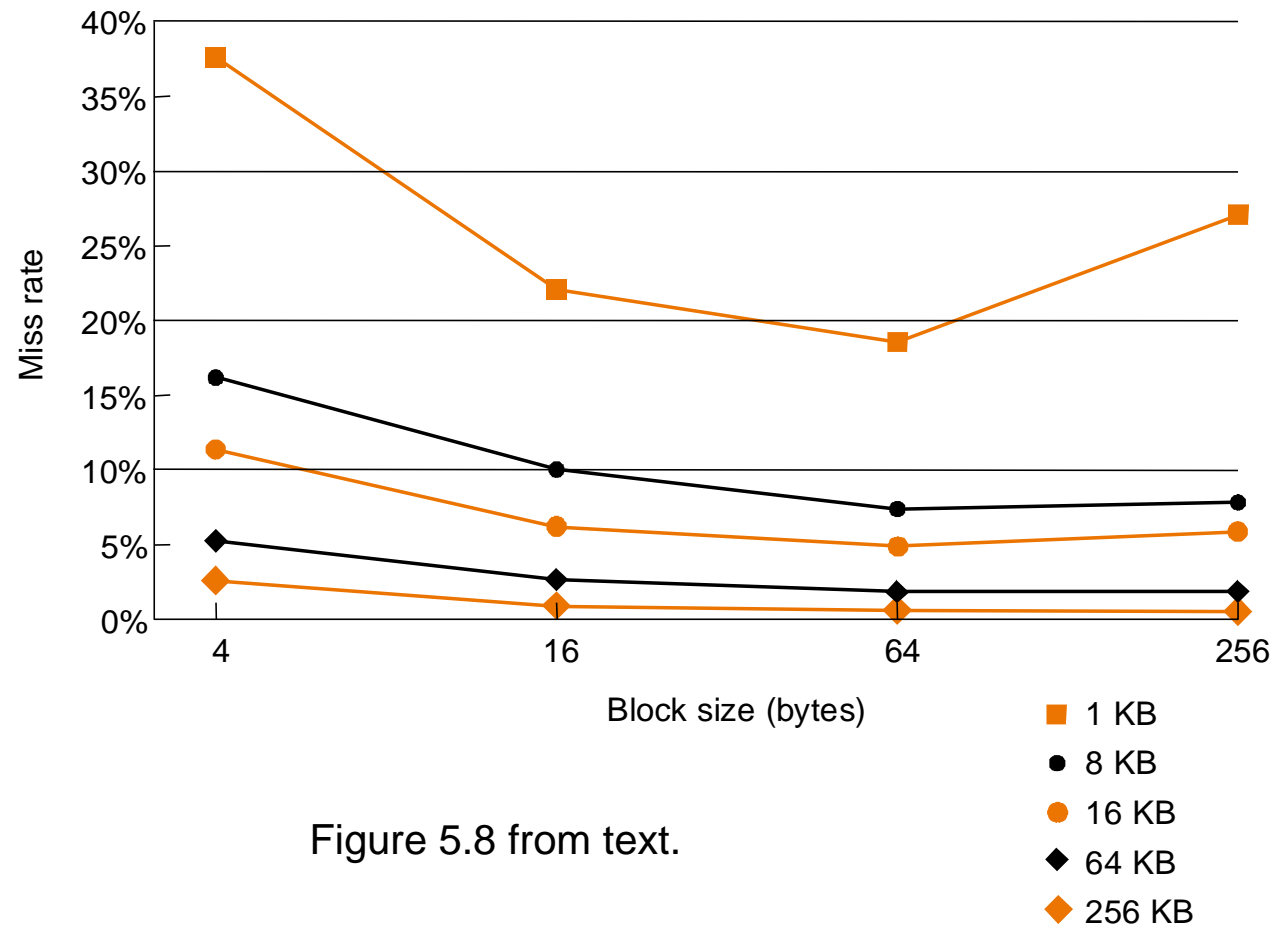
Block Size Tradeoff (1/2)

- Larger block size better exploits *spatial* locality but
 - Larger block size means larger *miss penalty*
 - Takes longer time to transfer the block
 - If block size is too big
 - Average access time goes up
 - Overall miss rate can increase b/c temporal locality is reduced when the replaced data would have been reused (too few lines in cache)
 - consider an extreme example for the previous cache: a single 64KB block

$$\text{Average Access Time} = \text{Hit Time} (1 - \text{MR}) + \text{Miss Penalty} \times \text{MR}$$



Block Size Tradeoff (2/2)



- Data from simulating a direct-mapped cache for VAX system traces
- Note miss rate trends
 - Larger block size causes higher miss rate

Block Size Considerations

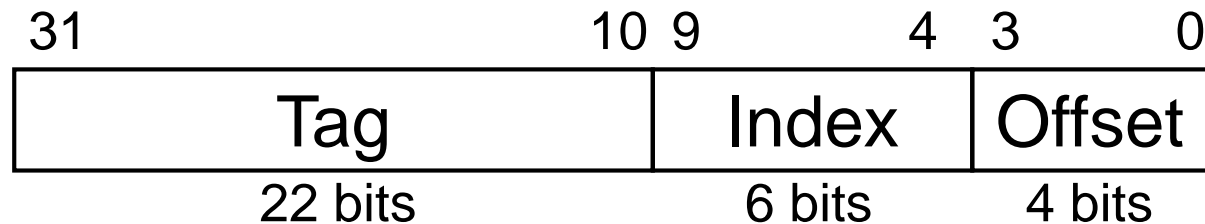
- Larger blocks should reduce miss rate
 - Due to spatial locality
- But in a fixed-sized cache
 - Larger blocks \Rightarrow fewer of them
 - More competition \Rightarrow increased miss rate
 - Larger blocks \Rightarrow pollution
- Larger miss penalty
 - Can override benefit of reduced miss rate

Example: Mapping

- 64 blocks, 16 bytes/block
 - To what block number does address 1200 map?

Block address = $\lfloor 1200/16 \rfloor = 75$

Block number = 75 modulo 64 = 11



Cache Misses

- On cache hit, CPU proceeds normally
- On cache miss
 - Stall the CPU pipeline
 - Fetch block from next level of hierarchy
 - Instruction cache miss
 - Restart instruction fetch
 - Data cache miss
 - Complete data access

Sources of Cache Misses

- 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)

Write-Through

- On data-write hit, could just update the block in cache
 - But then cache and memory would be inconsistent
- Write through: also update memory
- But makes writes take longer
 - e.g., if base CPI = 1, 10% of instructions are stores, write to memory takes 100 cycles
 - Effective CPI = $1 + 0.1 \times 100 = 11$
- Solution: **write buffer**
 - Holds data waiting to be written to memory
 - CPU continues immediately
 - Only stalls on write if write buffer is already full

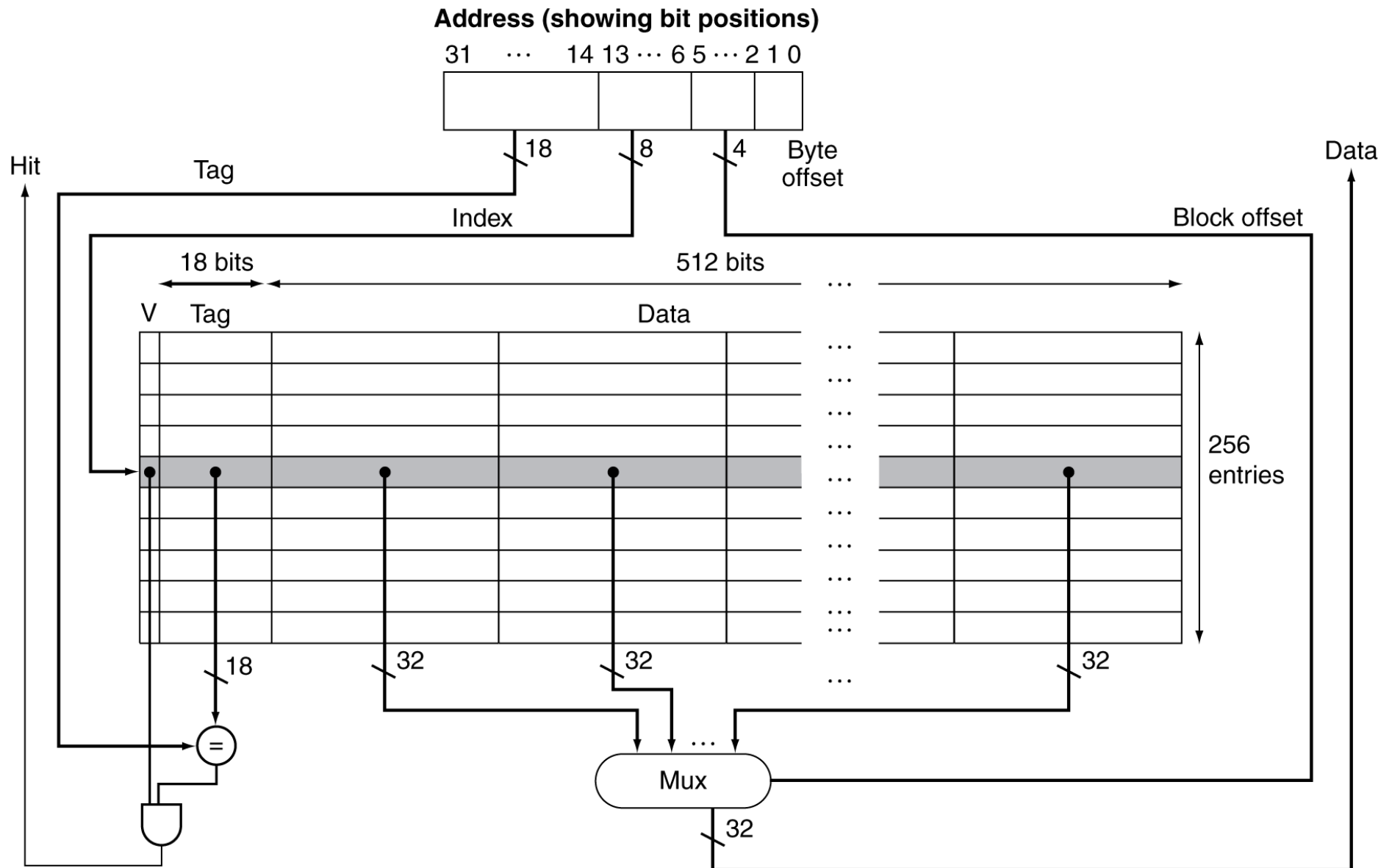
Write-Back

- Alternative: 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

Example: Intrinsity FastMATH

- Embedded processor- uses MIPS architecture and simple cache implementation
 - 12-stage pipeline
 - Instruction and data access on each cycle
- Split cache: separate I-cache and D-cache
 - Each 16KB: $256 \text{ blocks} \times 16 \text{ words/block}$
 - D-cache: write-through or write-back
- SPEC2000 miss rates
 - I-cache: 0.4%
 - D-cache: 11.4%
 - Weighted average: 3.2%

Example: Intrinsity FastMATH



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: $0.02 \times 100 = 2$
 - 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

Average Access Time

- Hit time is also important for performance
- Average memory access time (AMAT)
 - $AMAT = \text{Hit time} + \text{Miss rate} \times \text{Miss penalty}$
- Example
 - CPU with 1ns clock, hit time = 1 cycle, miss penalty = 20 cycles, I-cache miss rate = 5%
 - $AMAT = 1 + 0.05 \times 20 = 2\text{ns}$
 - 2 cycles per instruction

Performance Summary

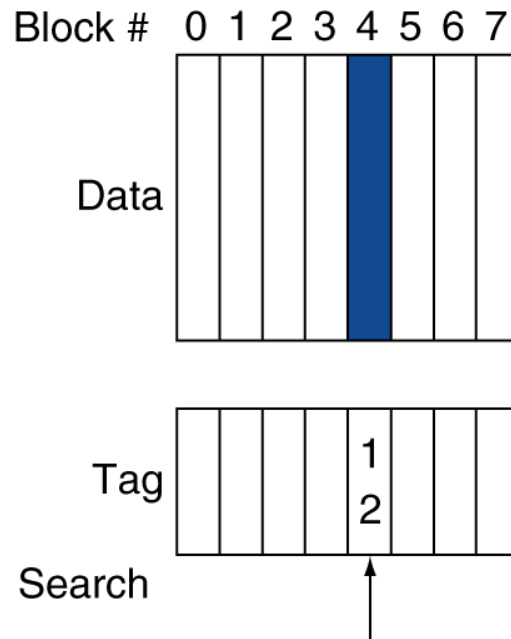
- When CPU performance increased
 - Miss penalty becomes more significant
 - stall will take increasing fraction of execution time)
- Decreasing base CPI
 - Greater proportion of time spent on memory stalls
 - Ex. from 2 to 1, $\text{CPI stall} / \text{CPI perf} = 4.44$ times faster
- Increasing clock rate
 - Memory stalls account for more CPU cycles
- Can't neglect cache behavior when evaluating system performance

Associative Caches

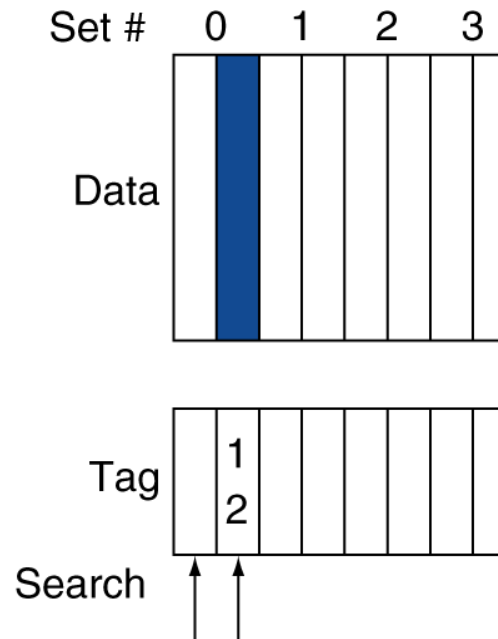
- Fully associative
 - Allow a given block to go in any cache entry
 - Requires all entries to be searched at once
 - Comparator per entry (expensive)
- n -way set associative
 - Each set contains n entries
 - Block number determines which set
 - (Block number) modulo (#Sets in cache)
 - Search all entries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

Associative Cache Example

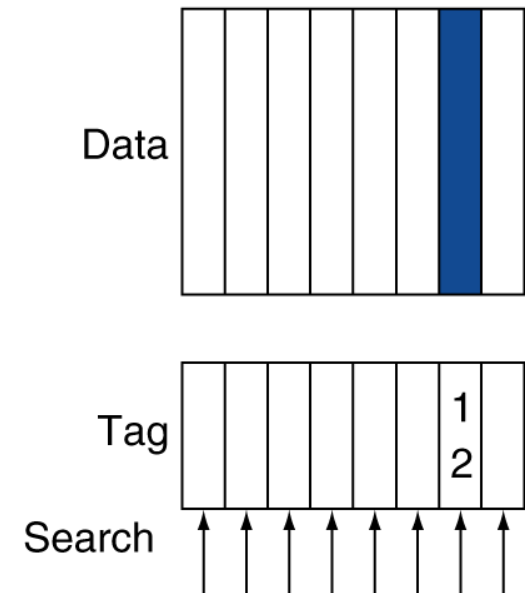
Direct mapped



Set associative



Fully associative



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				

Four-way set associative

Set	Tag	Data	Tag	Data	Tag	Data	Tag	Data
0								
1								

Eight-way set associative (fully associative)

[illegible]

Associativity Example

- Compare 4-block caches
 - Direct mapped, 2-way set associative, fully associative
 - 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

- 2-way set associative

Block address	Cache index	Hit/miss	Cache content after access			
			Set 0		Set 1	
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]		

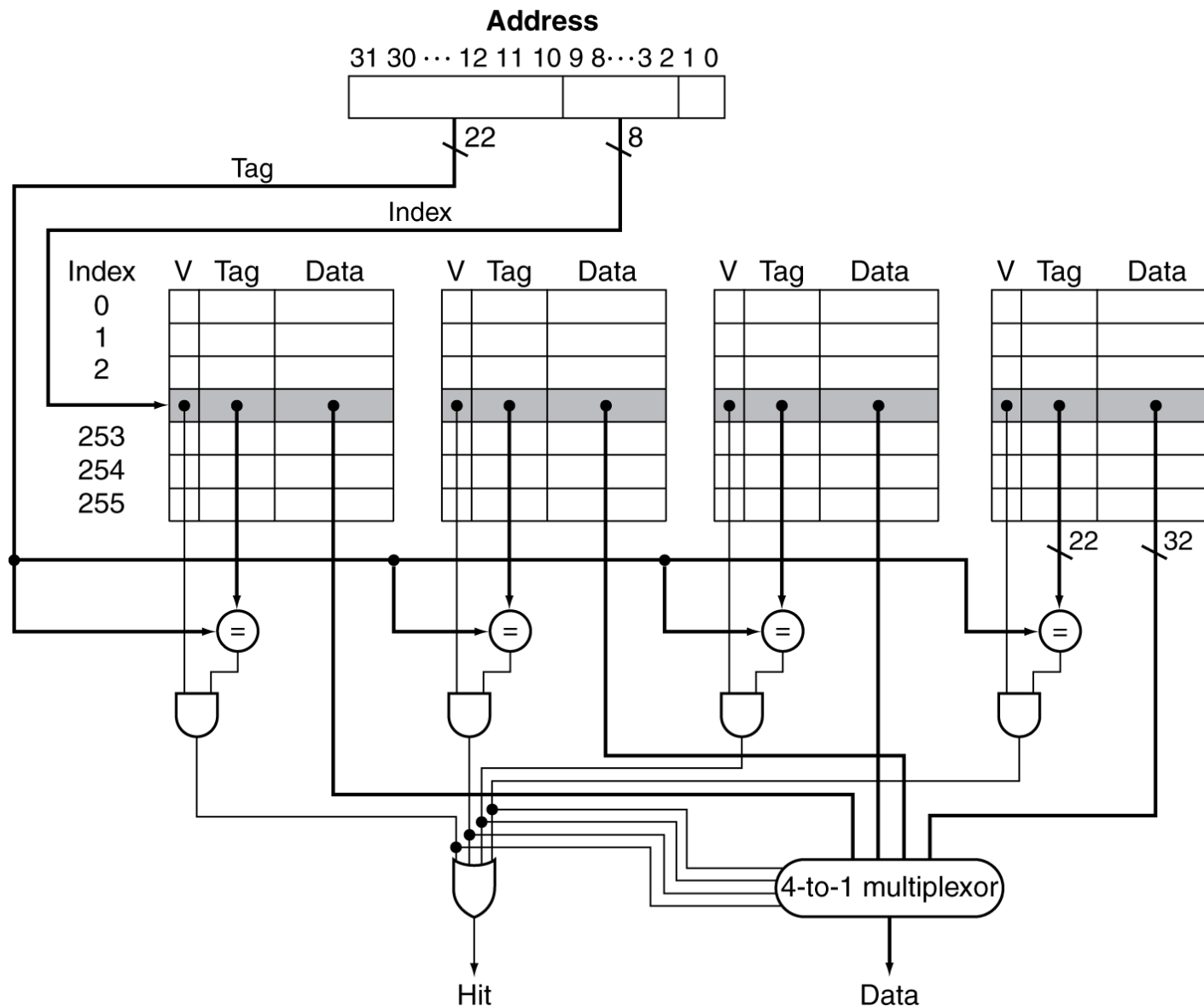
- Fully associative

Block address		Hit/miss	Cache content after access			
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
- 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%

Set Associative Cache Organization



4-way set associative - 4KB cache, 1 word/block, so $m=0$

Set Associative Cache Organization

Example:

4K blocks cache, 4-word block, and 32-bit address.

Find the total number of tag bits if cache is: direct-mapped, 2-way set and 4-way set associative and *fully* associative

Direct-mapped: 4 words=16 bytes so 4 bits, $32-4=28$ (for index and tag). Index $n = \log_2 4K = 12$, and tag=16 bit, hence total is $4K \times 16 = 64K$ tag bits

2-way-set: $4K \text{ block} / 2 = 2K$ sets, $n=11$, tag= 17 bits, the total is $2K \times 2 \times 17 = 34 \times 2K = 68K$ tag bits

4-way-set: 1K sets, $n=10$, tag=18, total number is $1K \times 4 \times 18 = 72 \times 1K = 72K$ tag bits

- Fully associative cache: one set with 4K blocks, tag=28, leading to $4K \times 1 \times 28 = 112K$ tag bits

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