

Cache Performance and Set Associative Cache

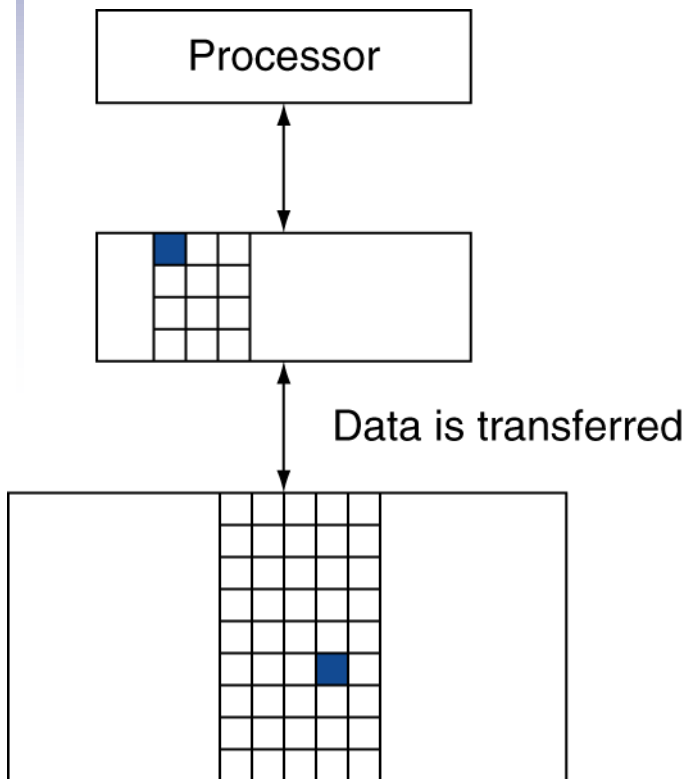
Lecture 12
CDA 3103
06-30-2014

Principle of Locality

- Programs access a small proportion of their address space at any time
- Temporal locality
 - Items accessed recently are likely to be accessed again soon
 - e.g., instructions in a loop, induction variables
- Spatial locality
 - Items near those accessed recently are likely to be accessed soon
 - E.g., sequential instruction access, array data

Memory Hierarchy Levels

- Block (aka line): unit of copying
 - May be multiple words
- If accessed data is present in upper level
 - Hit: access satisfied by upper level
 - Hit ratio: hits/accesses
- If accessed data is absent
 - Miss: block copied from lower level
 - Time taken: miss penalty
 - Miss ratio: misses/accesses
= $1 - \text{hit ratio}$
 - Then accessed data supplied from upper level

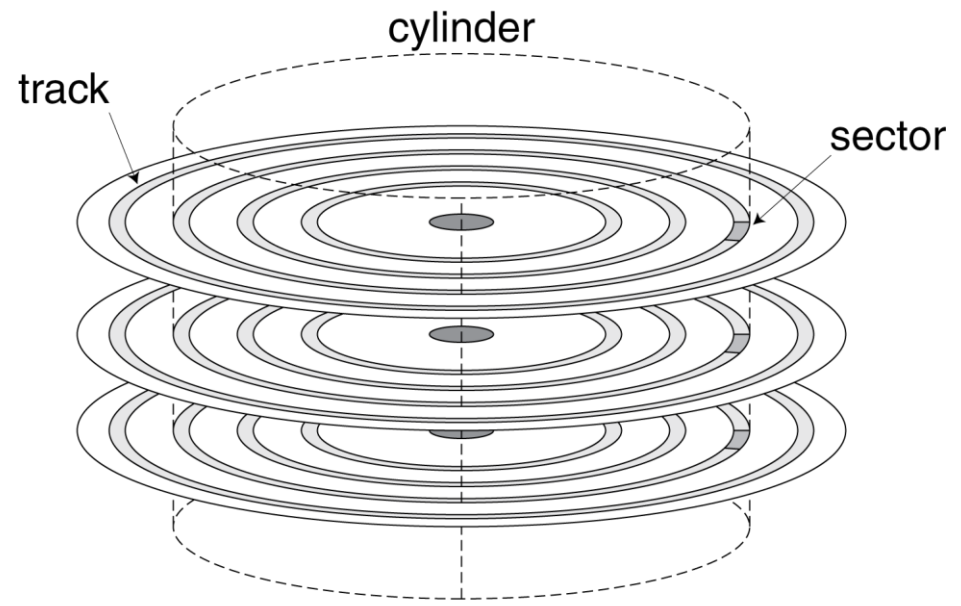


Memory Technology

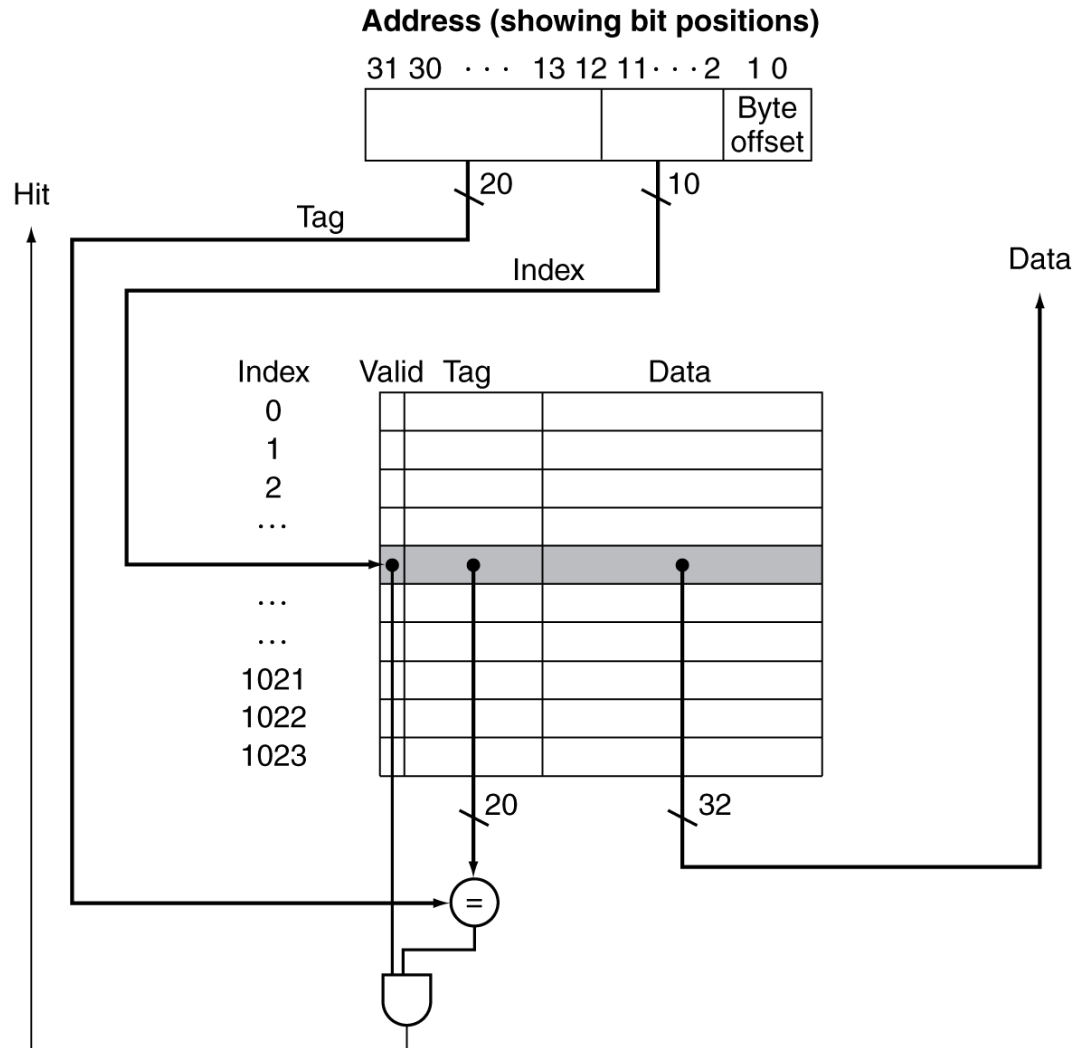
- Static RAM (SRAM)
 - 0.5ns – 2.5ns, \$2000 – \$5000 per GB
- Dynamic RAM (DRAM)
 - 50ns – 70ns, \$20 – \$75 per GB
- Magnetic disk
 - 5ms – 20ms, \$0.20 – \$2 per GB
- Ideal memory
 - Access time of SRAM
 - Capacity and cost/GB of disk

Disk Storage

- Nonvolatile, rotating magnetic storage



Address Subdivision



The number of bits in cache?

- $2^n \times (\text{block size} + \text{tag size} + \text{valid field size})$
- Cache size is 2^n blocks
- Block size is 2^m words (2^{m+2} words)
- Size of tag field $32 - (n + m + 2)$
- Therefore,
- $2^n \times (2^m \times 32 + 32 - (n + m + 2) + 1)$
- $= 2^n \times (2^m \times 32 + 31 - n - m)$

Question?

- How many total bits are required for a direct mapped cache with 16KiB of data and 4-word blocks, assuming 32 bit address?
- $2^n \times (2^m \times 32 + 31 - n - m)$

Anwer

- 16KiB = 4096 (2^{12} words)
- With Block size of 4 words (2^2) there are 1024 (2^{10}) blocks.
- Each block has 4 x 32 or 128 bits of data plus a tag which is 32 – 10 – 2 – 2 bits, plus a valid bit
- Thus total cache size is
- $2^{10} \times (4 \times 32 + (32 - 10 - 2 - 2) + 1) = 2^{10} \times 147 = 147 \text{ KibiBits}$

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
 - Early restart and critical-word-first can help

Block Size Tradeoff

■ Benefits of Larger Block Size

- **Spatial Locality**: if we access a given word, we're likely to access other nearby words soon
- Very applicable with Stored-Program Concept: if we execute a given instruction, it's likely that we'll execute the next few as well
- Works nicely in sequential array accesses too

■ Drawbacks of Larger Block Size

- Larger block size means **larger miss penalty**
 - on a miss, takes longer time to load a new block from next level
- If block size is too big relative to cache size, then there are too few blocks
 - Result: miss rate goes up

Extreme Example: One Big Block

Valid Bit



Tag

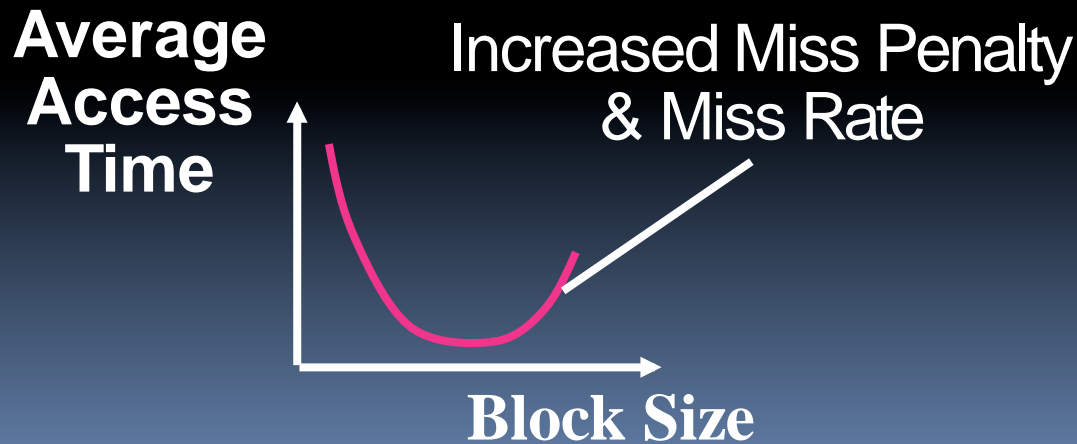
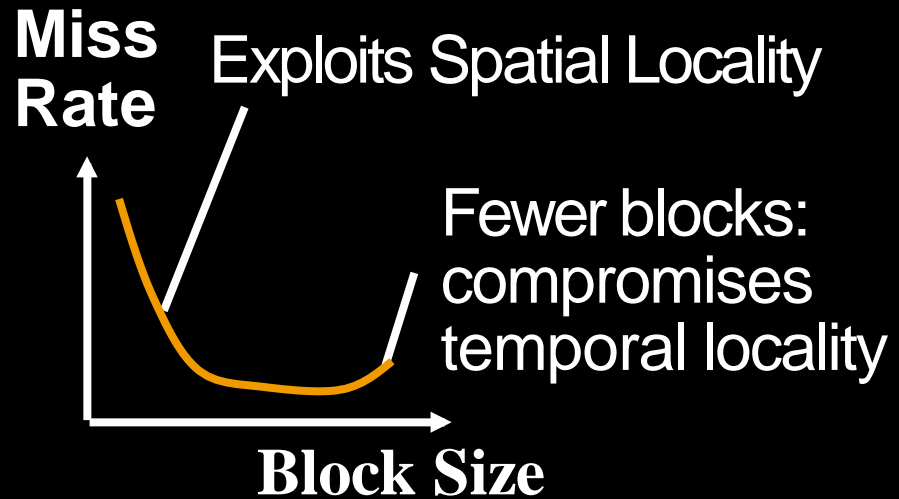
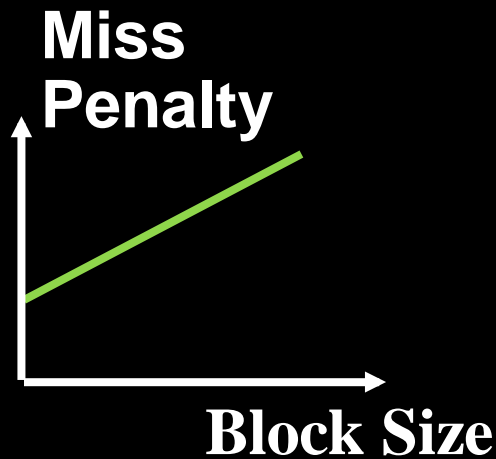


Cache Data



- **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 (force out) before use it again
 - Nightmare for cache designer: **Ping Pong Effect**

Block Size Tradeoff Conclusions



What to do on a write hit?

- **Write-through**

- update the word in cache block and corresponding word in memory

- **Write-back**

- update word in cache block
- allow memory word to be “stale”
- ⇒ add ‘dirty’ bit to each block indicating that memory needs to be updated when block is replaced
- ⇒ OS flushes cache before I/O...

- **Performance trade-offs?**

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

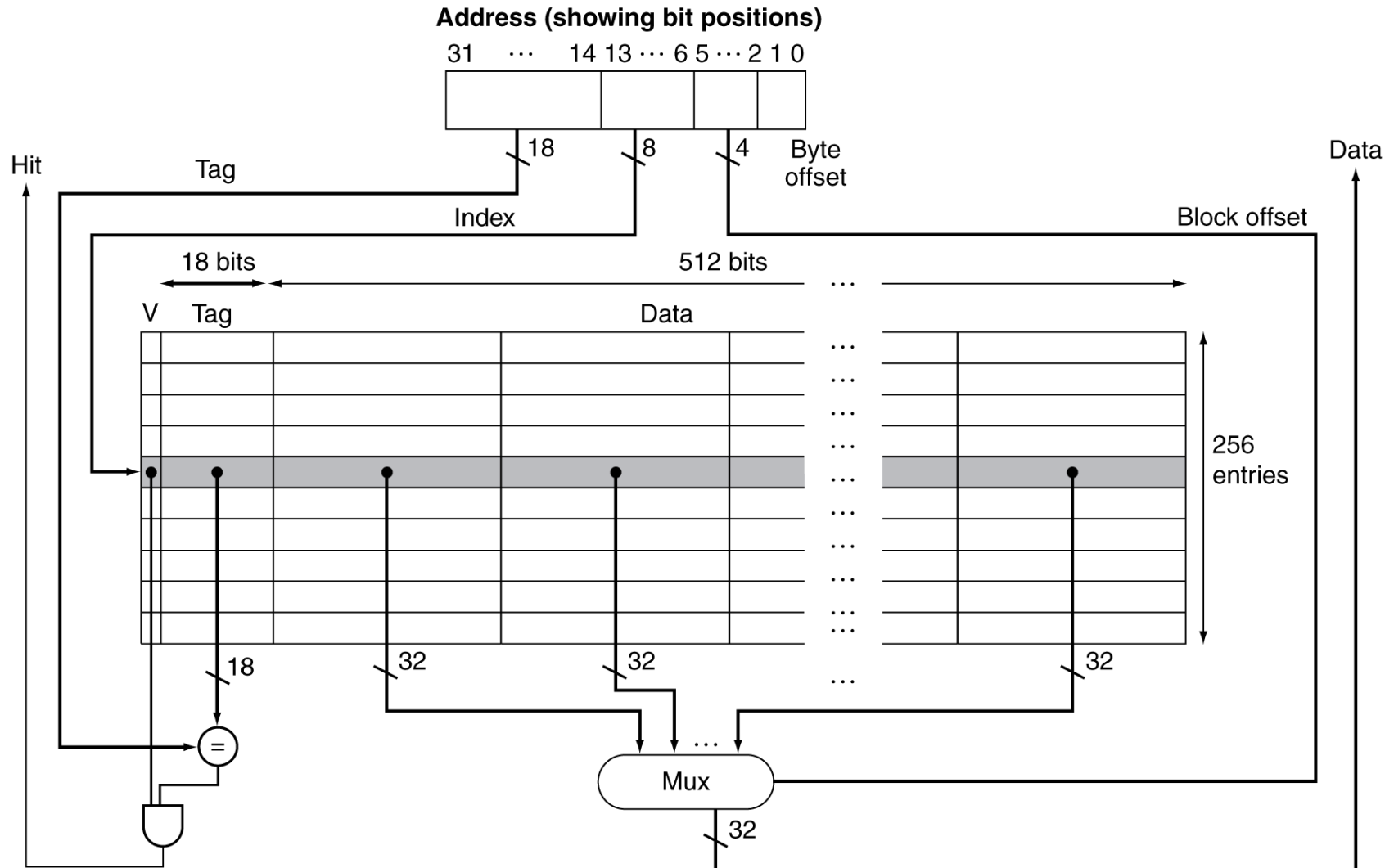
Write Allocation

- What should happen on a write miss?
- Alternatives for write-through
 - Allocate on miss: fetch the block
 - Write around: don't fetch the block
 - Since programs often write a whole block before reading it (e.g., initialization)
- For write-back
 - Usually fetch the block

Example: Intrinsicity FastMATH

- Embedded MIPS processor
 - 12-stage pipeline
 - Instruction and data access on each cycle
- Split cache: separate I-cache and D-cache
 - Each 16KB: 256 blocks × 16 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: Intrinsic FastMATH



Types of Cache Misses (1/2)

- “Three Cs” Model of Misses
- 1st 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, so won’t focus on these in this course

Pandora uses cache warm up

When should be cache performance measured?

Types of Cache Misses (2/2)

■ 2nd 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
- how do we lessen the effect of these?

■ **Dealing with Conflict Misses**

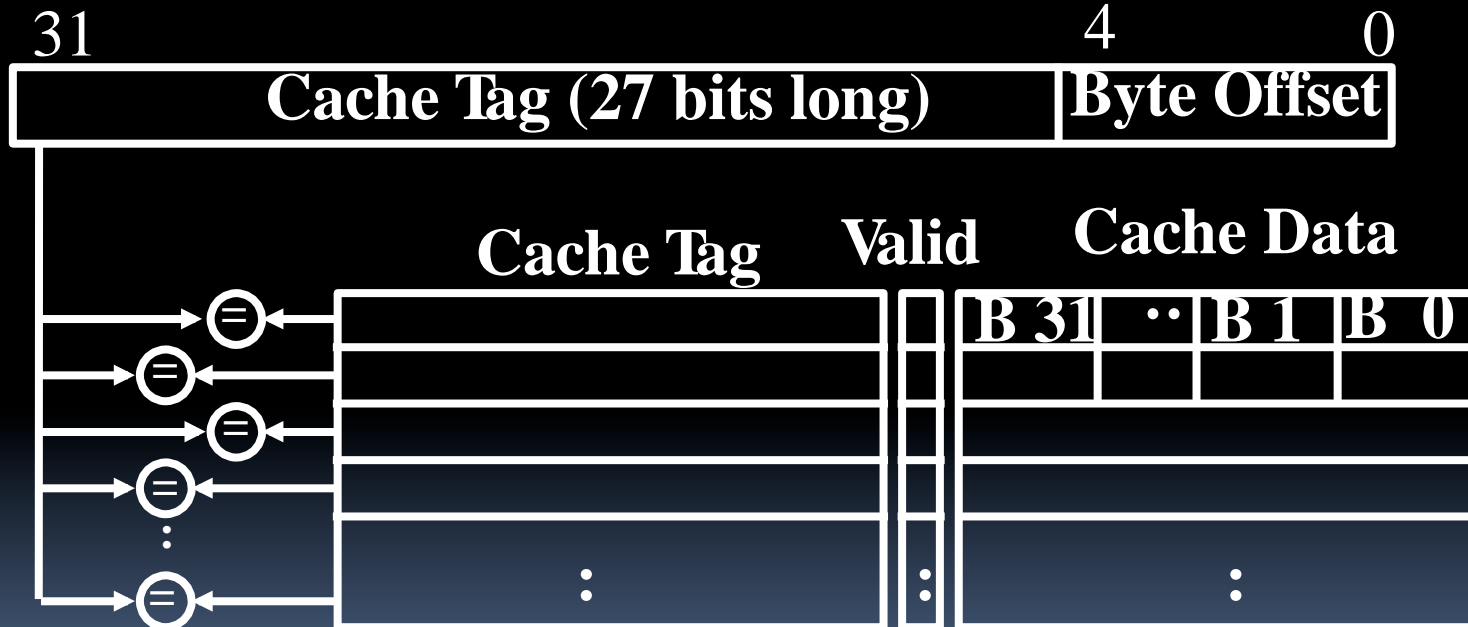
- Solution 1: Make the cache size bigger
 - Fails at some point
- Solution 2: Multiple distinct blocks can fit in the same cache Index?

Fully Associative Cache (1/3)

- **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 (2/3)

- Fully Associative Cache (e.g., 32 B block)
 - ▣ compare tags in parallel



Fully Associative Cache (3/3)

- **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 4B entries, we need 16K comparators: infeasible

Final Type of Cache Miss

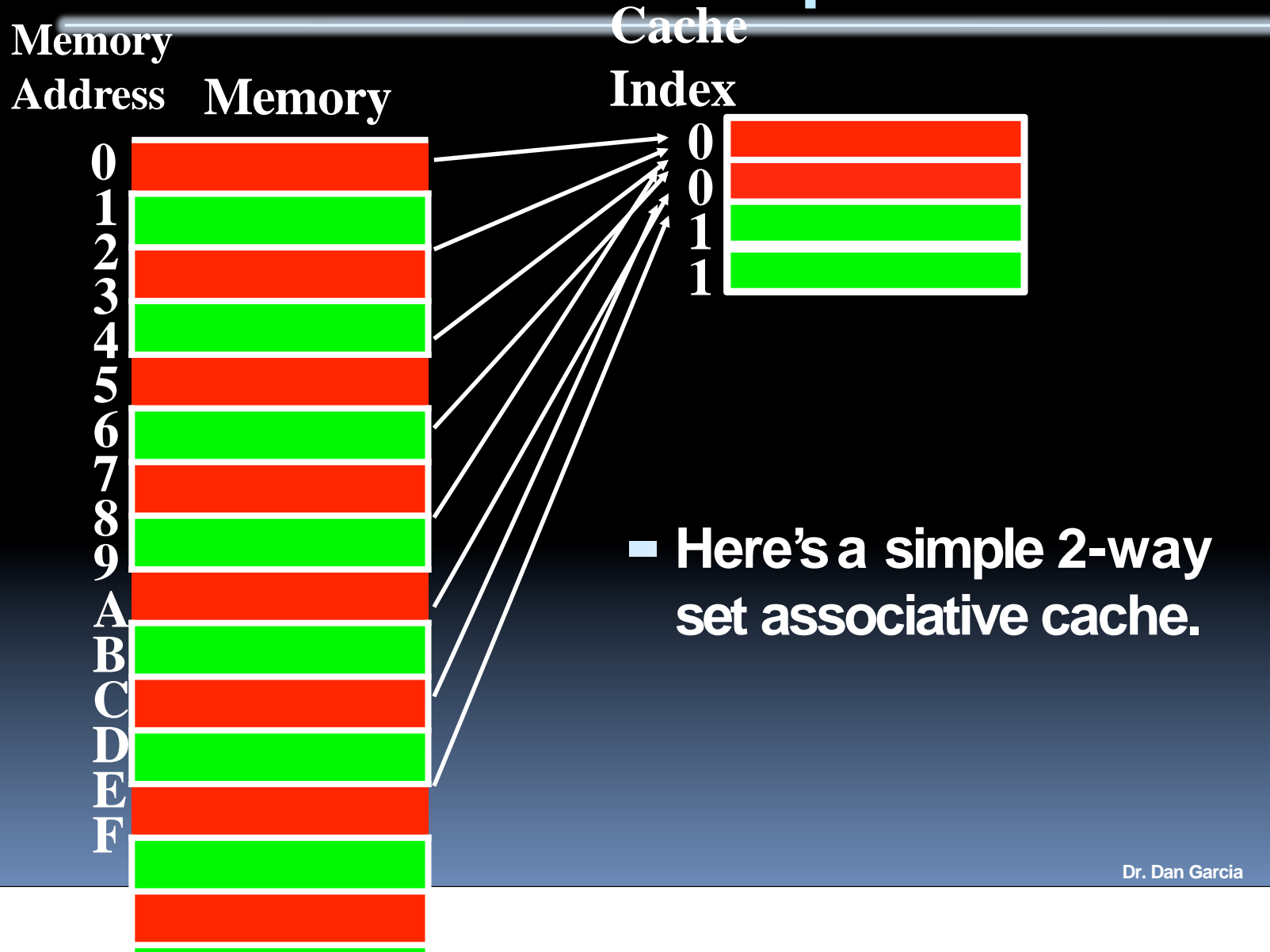
- 3rd C: **Capacity Misses**
 - miss that occurs because the cache has a limited size
 - miss that would not occur if we increase the size of the cache
 - sketchy definition, so just get the general idea
- **This is the primary type of miss for Fully Associative caches.**

N-Way Set Associative Cache (1/3)

- **Memory address fields:**
 - **Tag:** same as before
 - **Offset:** same as before
 - **Index:** points us to the correct “row” (called a set in this case)
- **So 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

Is the temporal or spatial locality exploited here?

Associative Cache Example



N-Way Set Associative Cache (2/3)

- **Basic Idea**

- ▣ cache is direct-mapped w/respect to sets
- ▣ each set is fully associative with N blocks in it

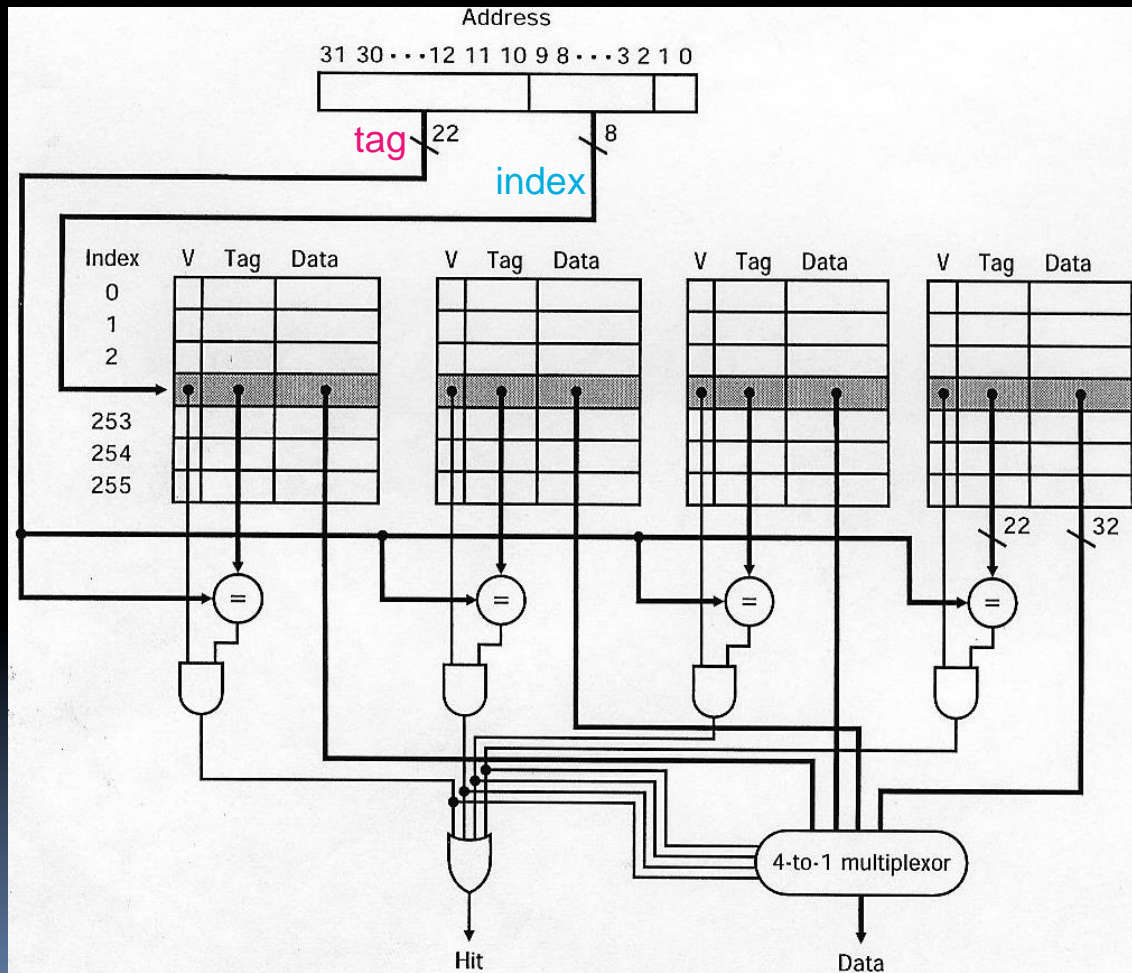
- **Given memory address:**

- ▣ Find correct set using Index value.
- ▣ Compare $\bar{T}ag$ with all $\bar{T}ag$ 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 (3/3)

- **What's so great about this?**
 - even a 2-way set assoc 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 assoc
 - it's **Fully Assoc** if it's M-way set assoc
 - so these two are just special cases of the more general set associative design

4-Way Set Associative Cache Circuit



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)

Tag	Data	Tag	Data	Tag	Data	Tag	Data	Tag	Data	Tag	Data	Tag	Data	Tag	Data

Associativity Example

- Compare 4-block caches
 - Direct mapped, 2-way set associative, fully associative
 - Block access sequence: 0, 8, 0, 6, 8
- For direct map
 - (Block address) modulo (Number of block in the cache)
- For set-associative
 - (Block address) modulo (Number of sets in the cache)

Direct-Mapped Cache

- Direct mapped

Block Address	Cache Block
0	$(0 \text{ modulo } 4) = 0$
6	$(6 \text{ modulo } 4) = 2$
0	$(8 \text{ modulo } 4) = 0$

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]		

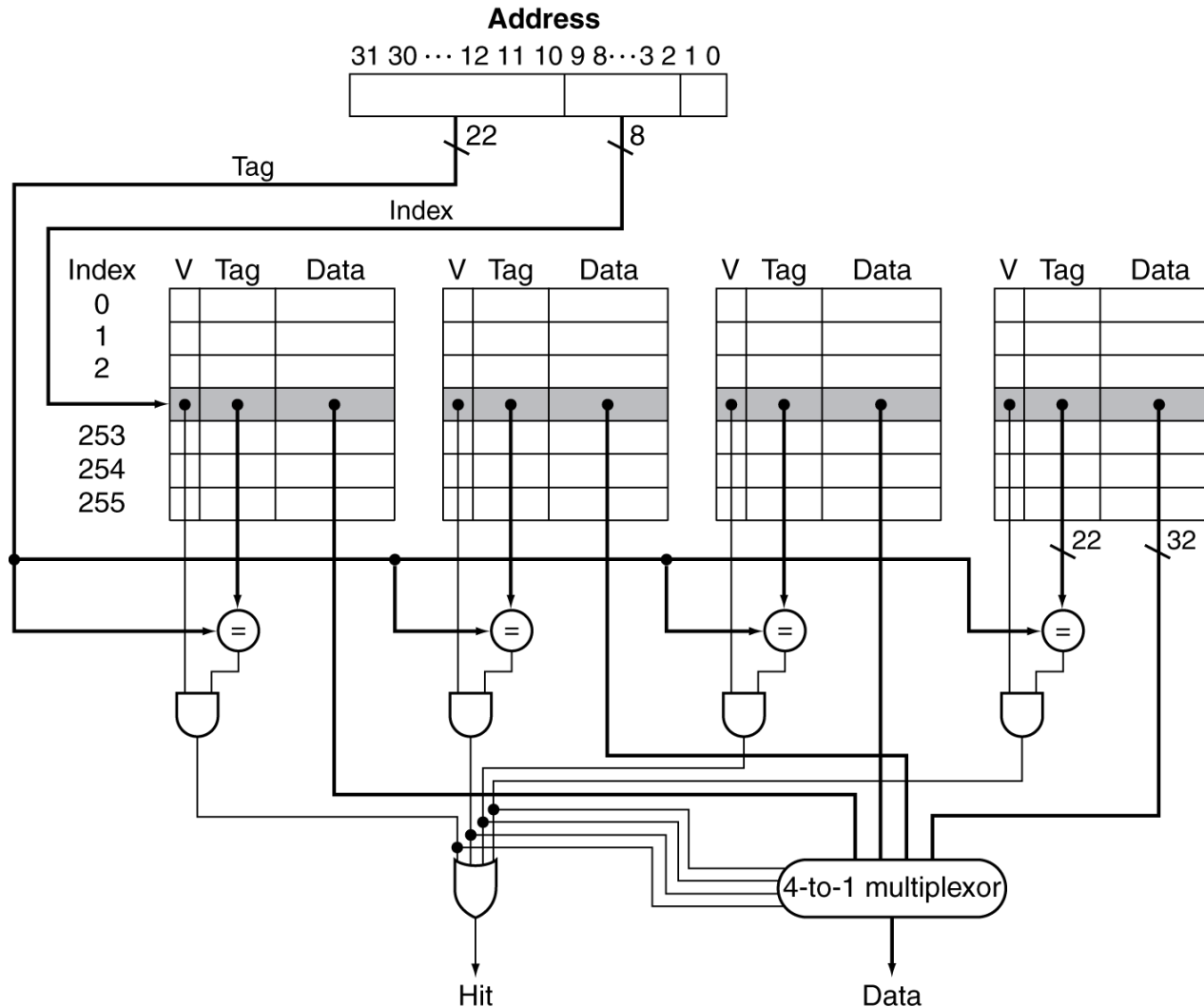
- Fully associative

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



Block Replacement Policy

- **Direct-Mapped Cache**

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

- **N-Way Set Assoc**

- 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

- **Question: 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: LRU

■ LRU (Least Recently Used)

- Idea: cache out block which has been accessed (read or write) least recently
- Pro: **temporal locality** \Rightarrow recent past use implies likely future use: in fact, this is a very effective policy
- Con: with 2-way set assoc, easy to keep track (one LRU bit); with 4-way or greater, requires complicated hardware and much time to keep track of this

Block Replacement Example

- We have a 2-way set associative cache with a four word total capacity and one word blocks. We perform the following word accesses (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 Example: LRU

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)

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

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

0: hit

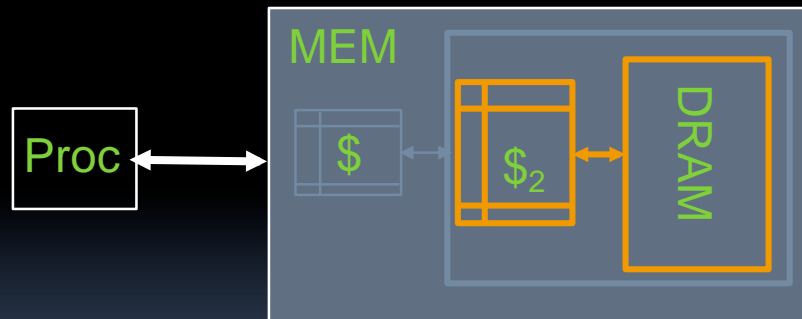
	loc 0	loc 1
set 0	0	lru
set 1		
set 0	lru 0	2
set 1		
set 0	0	lru 2
set 1		
set 0	0	lru 2
set 1	1	lru
set 0	lru 0	4
set 1	1	lru
set 0	0	lru 4
set 1	1	lru

Big Idea

- How to choose between associativity, block size, replacement & write policy?
- Design against a performance model
 - Minimize: Average Memory Access Time
 - = Hit Time
 - + Miss Penalty x Miss Rate
 - influenced by technology & program behavior
- Create the illusion of a memory that is large, cheap, and fast - on average
- How can we improve miss penalty?

Improving Miss Penalty

- When caches first became popular, Miss Penalty ~ 10 processor clock cycles
- Today 2400 MHz Processor (0.4 ns per clock cycle) and 80 ns to go to DRAM
 - 200 processor clock cycles!



Solution: another cache between memory and the processor cache: **Second Level (L2) Cache**

Peer Instruction

1. A 2-way set-associative cache can be outperformed by a direct-mapped cache.
2. Larger block size \Rightarrow lower miss rate

	1	2
a)	F	F
b)	F	T
c)	T	F
d)	T	T

Peer Instruction Answer

1. Sure, consider the caches from the previous slides with the following workload: 0, 2, 0, 4, 2
2-way: 0m, 2m, 0h, 4m, 2m;
DM: 0m, 2m, 0h, 4m, 2h

2. Larger block size \Rightarrow lower miss rate, true until a certain point, and then the ping-pong effect takes over

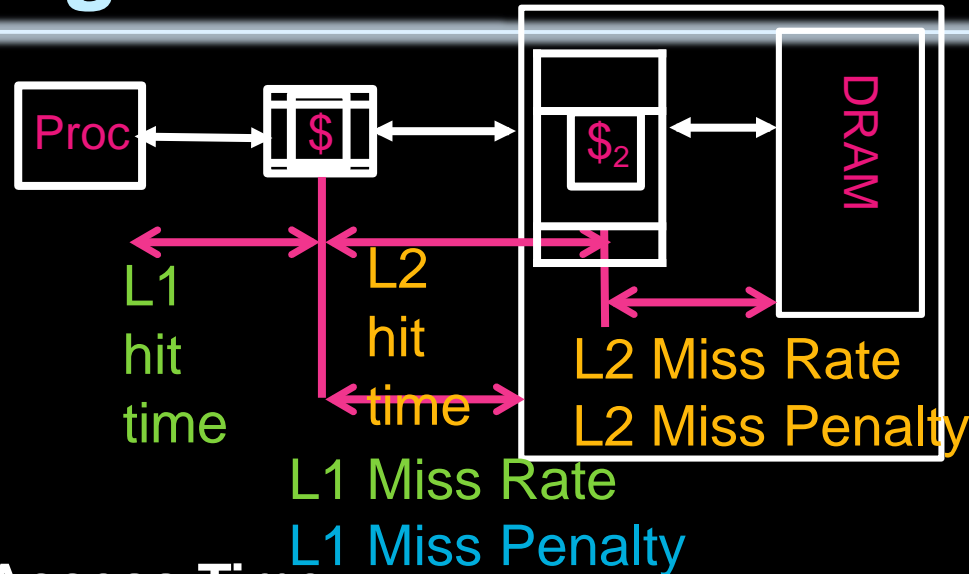
1. A 2-way set-associative cache can be outperformed by a direct-mapped cache.
2. Larger block size \Rightarrow lower miss rate

	12
a)	FF
b)	FT
c)	TF
d)	TT

And in Conclusion...

- We've discussed memory caching in detail. Caching in general shows up over and over in computer systems
 - Filesystem cache, Web page cache, Game databases / tablebases, Software memoization, Others?
- Big idea: if something is expensive but we want to do it repeatedly, do it once and cache the result.
- Cache design choices:
 - Size of cache: speed v. capacity
 - Block size (i.e., cache aspect ratio)
 - Write Policy (Write through v. write back)
 - Associativity choice of N (direct-mapped v. set v. fully associative)
 - Block replacement policy
 - 2nd level cache?
 - 3rd level cache?
- Use performance model to pick between choices, depending on programs, technology, budget, ...

Analyzing Multi-level cache hierarchy



Avg Mem Access Time =

$$\text{L1 Hit Time} + \text{L1 Miss Rate} * \text{L1 Miss Penalty}$$

L1 Miss Penalty =

$$\text{L2 Hit Time} + \text{L2 Miss Rate} * \text{L2 Miss Penalty}$$

Avg Mem Access Time =

$$\text{L1 Hit Time} + \text{L1 Miss Rate} * (\text{L2 Hit Time} + \text{L2 Miss Rate} * \text{L2 Miss Penalty})$$

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

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

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

Question

- Assume the miss rate of an instruction cache is 2% and the miss rate of the data cache is 4%.
If a processor has CPI of 2 without any memory stalls and the miss penalty is 100 cycles for all misses,
- Determine how much faster a processor would run with perfect cache that never missed?

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, l-cache miss rate = 5%
 - $AMAT = 1 + 0.05 \times 20 = 2\text{ns}$
 - 2 cycles per instruction

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 Considerations

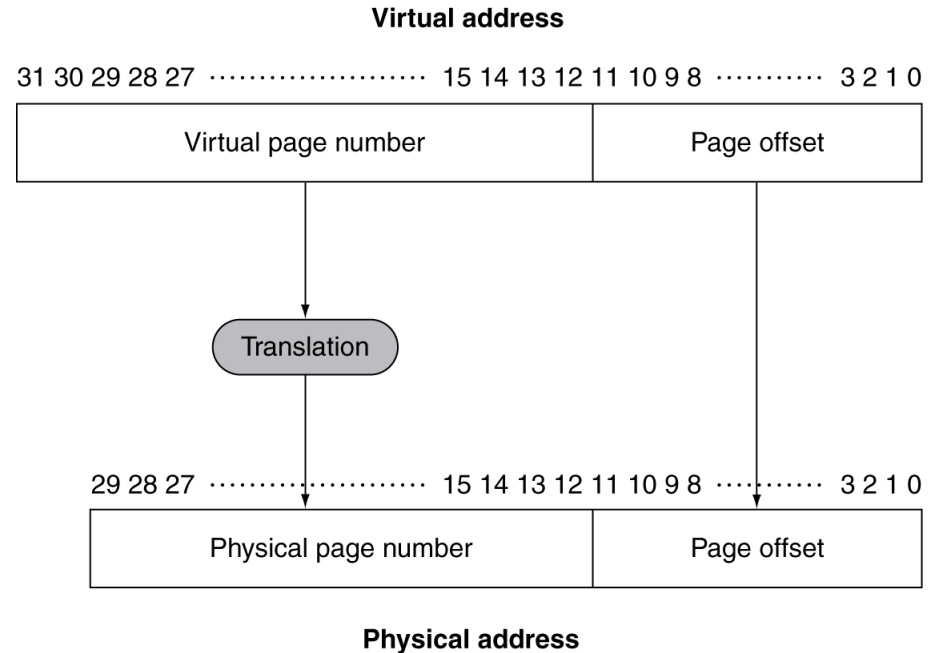
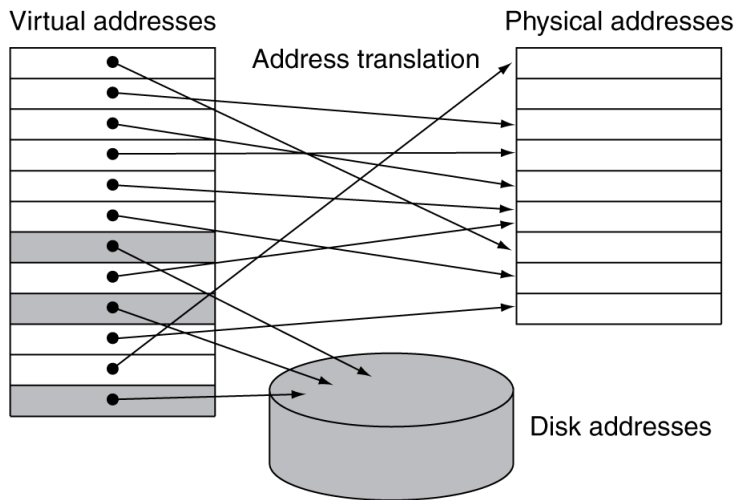
- Primary cache
 - Focus on minimal hit time
- L-2 cache
 - Focus on low miss rate to avoid main memory access
 - Hit time has less overall impact
- Results
 - L-1 cache usually smaller than a single cache
 - L-1 block size smaller than L-2 block size

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

Address Translation

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



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)

The Memory Hierarchy

The BIG Picture

- Common principles apply at all levels of the memory hierarchy
 - Based on notions of caching
- At each level in the hierarchy
 - Block placement
 - Finding a block
 - Replacement on a miss
 - Write policy

Finding a Block

Associativity	Location method	Tag comparisons
Direct mapped	Index	1
n-way set associative	Set index, then search entries within the set	n
Fully associative	Search all entries	#entries
	Full lookup table	0

- Hardware caches
 - Reduce comparisons to reduce cost
- Virtual memory
 - Full table lookup makes full associativity feasible
 - Benefit in reduced miss rate

Concluding Remarks

- Fast memories are small, large memories are slow
 - We really want fast, large memories ☹️
 - Caching gives this illusion 😊
- Principle of locality
 - Programs use a small part of their memory space frequently
- Memory hierarchy
 - L1 cache ↔ L2 cache ↔ ... ↔ DRAM memory
↔ disk
- Memory system design is critical for multiprocessors