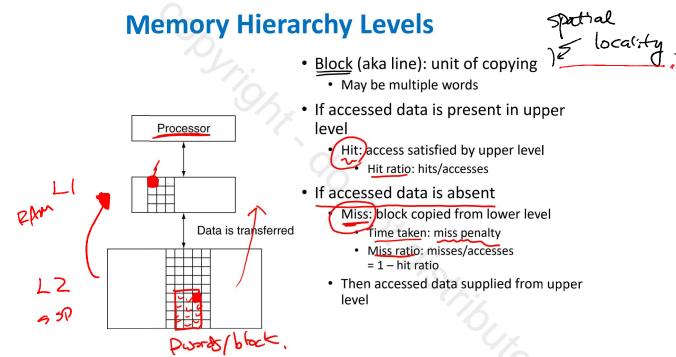
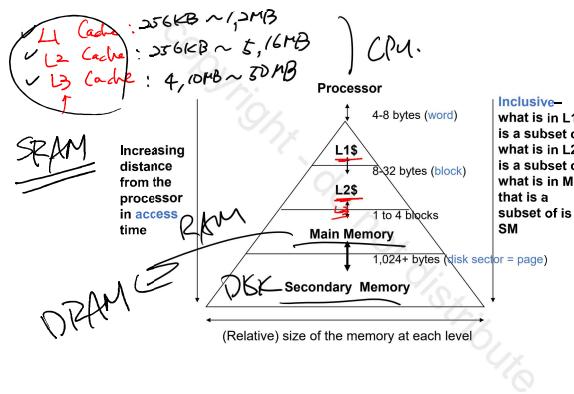


### Taking Advantage of Locality

- Memory hierarchy
- Store everything on disk — **OPT. 1 spatial**
- Copy recently accessed (and nearby) items from disk to smaller DRAM memory
  - Main memory — **OPT. 2**
- Copy more recently accessed (and nearby) items from DRAM to smaller SRAM memory
  - Cache memory attached to CPU — **OPT. 3**

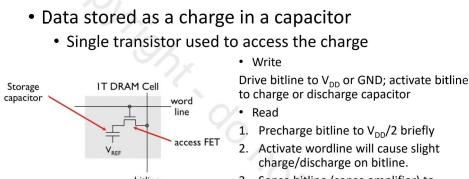


### Memory Technology

- 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

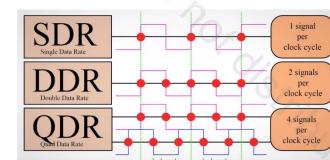
### DRAM Technology

- Data stored as a charge in a capacitor
  - Single transistor used to access the charge
    - Write: Drive bitline to  $V_{DD}$  or GND; activate bitline to charge or discharge capacitor
    - Read: Sense bitline (sense amplifier) to determine binary 0/1
- Must periodically be refreshed due to Read problem
  - Capacitor loses charge, thus needs to be written back
  - Sense amplifier charge bitline to  $V_{DD}$ /GND



### Advanced DRAM Organization

- Bits in a DRAM are organized as a rectangular array
- Double data rate (DDR) DRAM
  - Transfer on rising and falling clock edges
- Quad data rate (QDR) DRAM
  - Separate DDR inputs and outputs



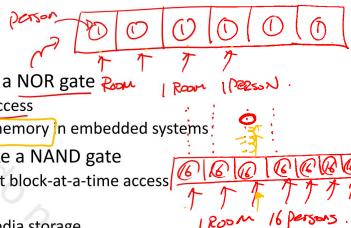
### Flash Storage

- Nonvolatile semiconductor storage
  - 100x – 1000x faster than disk
  - Smaller, lower power, more robust
  - But more \$/GB (between disk and DRAM)



## Flash Types

- NOR flash:** bit cell like a NOR gate
  - Random read/write access
  - Used for instruction memory in embedded systems
- MAND flash:** bit cell like a NAND gate
  - Denser (bits/area), but block-at-a-time access
  - Cheaper per GB
  - Used for USB keys, media storage, ...
- Flash bits wears out after 1000's of accesses
  - Not suitable for direct RAM or disk replacement
  - Wear leveling: remap data to less used blocks



## Disk Storage

- Nonvolatile, rotating magnetic storage



10

11

## Disk Access Example

- Given
  - 512B sector, 15,000rpm, 4ms average seek time, 100MB/s transfer rate, 0.2ms controller overhead, idle disk
- Average read time = Queue + Seek + Rotational + Transfer + Controller Delay
  - 4ms seek time  
+  $\frac{1}{2} / (15,000/60) = 2\text{ms}$  rotational latency  
+  $512 / 100\text{MB/s} = 0.005\text{ms}$  transfer time  
+ 0.2ms controller delay  
= 6.2ms
- If actual average seek time is 1ms
  - Average read time = 3.2ms

$$\frac{1}{2} \times \frac{15000}{60} = 2\text{ms}$$

Queue + Seek + Rotational + Transfer + Controller Delay  
 4ms + 1ms + 2ms + 0.005ms + 0.2ms = 6.2ms

## Disk Performance Issues

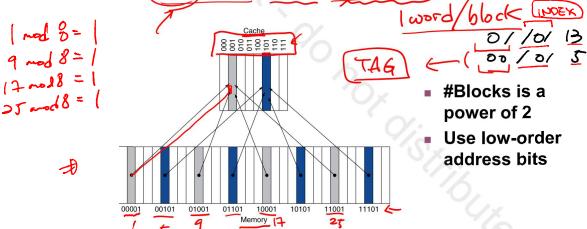
- Manufacturers quote average seek time
  - Based on all possible seeks
  - Locality and OS scheduling lead to smaller actual average seek times
- Smart disk controller allocate physical sectors on disk
  - Present logical sector interface to host
  - SCSI, ATA, SATA
- Disk drives include caches
  - Prefetch sectors in anticipation of access
  - Avoid seek and rotational delay

13

14

## Direct Mapped Cache Vs. Associative Cache.

- Location determined by address
- Direct mapped: only one choice
  - (Block address) modulo (#Blocks in cache)



15

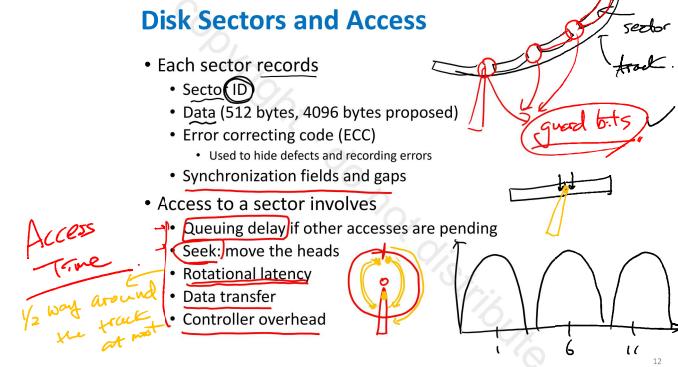
## Tags and Valid Bits

- How do we know which particular block is stored in a cache location?
  - Store block address as well as the data
  - Actually, only need the high-order bits
  - Called the tag
- What if there is no data in a location?
  - Valid bit: 1 = present, 0 = not present
  - Initially 0

16

## Disk Sectors and Access

- Each sector records
  - Sector ID
  - Data (512 bytes, 4096 bytes proposed)
  - Error correcting code (ECC)
    - Used to hide defects and recording errors
  - Synchronization fields and gaps
- Access to a sector involves
  - Queuing delay if other accesses are pending
  - Seek: move the heads
  - Rotational latency
  - Data transfer
  - Controller overhead

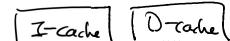


12

## Cache Memory

**SPEED UP!**  
Caching puts all your storage for accessing fast, & data

- Cache memory
  - The level of the memory hierarchy closest to the CPU
- Given accesses  $X_1, \dots, X_{n-1}, X_n$



$X_4$
$X_1$
$X_{n-2}$
$X_{n-1}$
$X_2$
$X_n$
$X_3$

$X_4$
$X_1$
$X_{n-2}$
$X_{n-1}$
$X_2$
$X_n$
$X_3$

a. Before the reference to  $X_n$ b. After the reference to  $X_n$ 

- How do we know if the data is present?
- Where do we look?

## Cache Example

- 8-blocks, 1 word/block, direct mapped
- Initial state

Index	V	Tag	Data	
			Valid	
000	N = 0			
001	N			
010	N			
011	N			
100	N			
101	N			
110	N			
111	N			

*cold start*

18

## Cache Example

Word addr	Binary addr	Hit/miss	Cache block
26	11 010	Miss	010

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

20

## Cache Example

Word addr	Binary addr	Hit/miss	Cache block
22	10 110	Hit ✓	110
26	11 010	Hit	010

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

21

## Cache Example

Word addr	Binary addr	Hit/miss	Cache block
16	10 000	Miss ✓	000
3	00 011	Miss ✓	011
16	10 000	Hit ✓	000

Index	V	Tag	Data
<b>000</b>	<b>Y=1</b>	<b>10</b>	<b>Mem[10000]</b>
001	N		
<b>010</b>	<b>Y</b>	<b>11</b>	<b>Mem[11010]</b>
<b>011</b>	<b>Y</b>	<b>00</b>	<b>Mem[00011]</b>
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

22

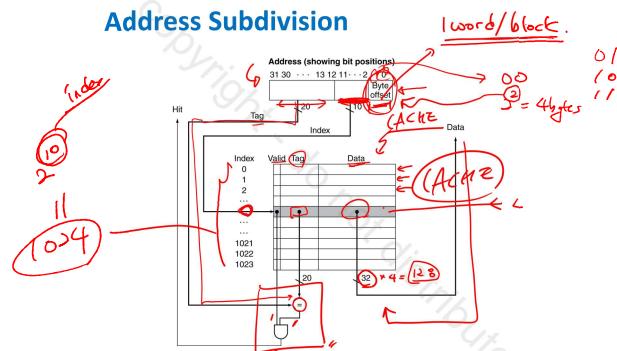
## Cache Example

Word addr	Binary addr	Hit/miss	Cache block
18	10 010	Miss	010

Index	V	Tag	Data
000	Y	10	Mem[10000]
<b>010</b>	<b>Y=1</b>	<b>10</b>	<b>Mem[10010]</b>
011	Y	00	Mem[00011]
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

23

## Address Subdivision



24

## 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

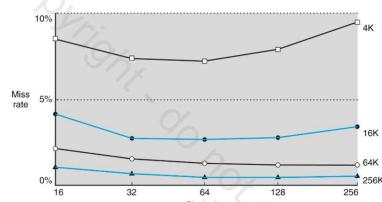


FIGURE 5.11 Miss rate versus block size. Note that the miss rate actually goes up if the block size is too large relative to the cache size. Each line represents a cache of different size. (This figure is independent of associativity, discussed soon.) Unfortunately, SPEC CPU2000 traces would take too long if block size were included, so this data is based on SPEC92.

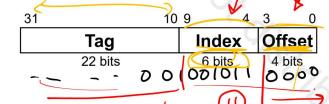
25

## Associative Caches

### Vs. Direct Mapped

- 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)

Four way set associative

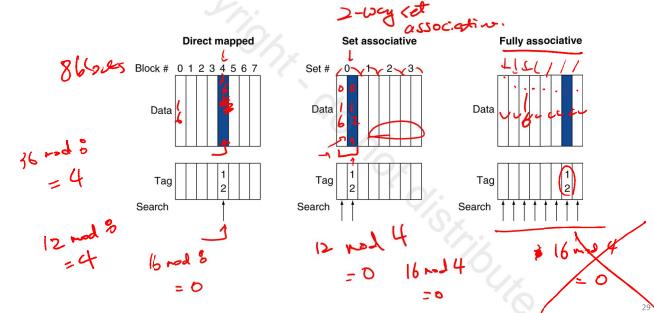


26

27

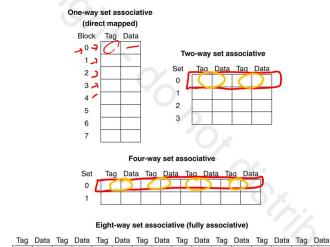
28

## Associative Cache Example



## Spectrum of Associativity

- For a cache with 8 entries



## Associativity Example

- Compare 4-block caches

- Direct mapped, 2-way set associative, fully associative
- Block access sequence: 0, 8, 0, 6, 8  
 $0 \downarrow 8 \downarrow 0 \downarrow 6 \downarrow 8 \pmod{4}$
- 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[0]	Mem[8]			
0	0	hit	Mem[0]	Mem[8]			
6	0	miss	Mem[0]	Mem[8]	Mem[6]		
8	0	miss	Mem[8]	Mem[6]			

0 hits

## Associativity Example

- 2-way set associative

Block address	Cache index	Hit/miss		Cache content after access	
		Set 0	Set 1	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[8]	Mem[6]
8	0	miss	Mem[8]	Mem[6]	

① LRU  
② LFU  
③ Random

- Fully associative

Block address		Hit/miss		Cache content after access	
		0	8	0	8
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]

6 hits

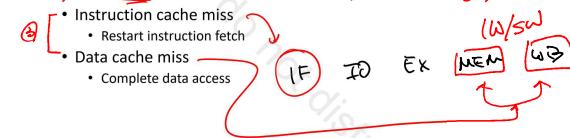
## Replacement Policy

- 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

32

## Cache Misses

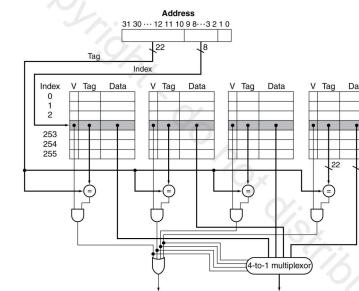
- On cache hit, CPU proceeds normally
- On cache miss
  - Stall the CPU pipeline
  - Fetch block from next level of hierarchy
    - Instruction cache miss
      - Restart instruction fetch
    - Data cache miss
      - Complete data access



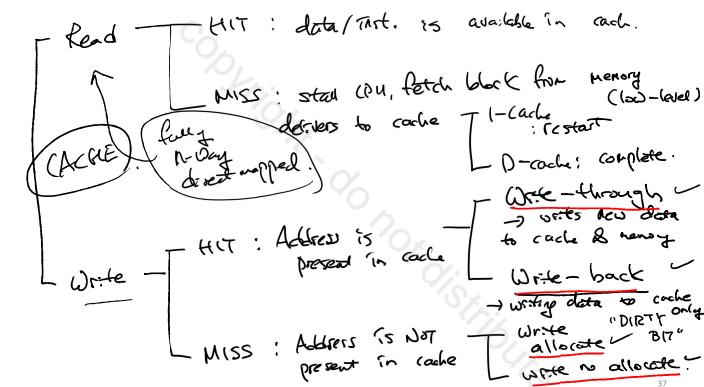
35

33

## Set Associative Cache Organization



34



36

37

## 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

38

## 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
- 

39

## Example: Intrinsic 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  $\times$  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%

41

## Main Memory Supporting Caches

- Use DRAMs for main memory
  - Fixed width (e.g., 1 word)
  - Connected by fixed-width clocked bus
    - Bus clock is typically slower than CPU clock
- Example cache block read
  - 1 bus cycle for address transfer ✓
  - 15 bus cycles per DRAM access
  - 1 bus cycle per data transfer ✓
- For 4-word/block, 1-word-wide DRAM
  - Miss penalty =  $1 + 4 \times 15 + 4 \times 1 = 65$  bus cycles
  - Bandwidth = 16 bytes / 65 cycles = 0.25 B/cycle

42

## 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 CPI is  $5.44/2 = 2.72$  times faster
  - Miss rate 0%

$$\text{CPU time} = \text{IC} \times \text{CPI} \times \text{cc.}$$

↑ 5 CPI  $\times$  IC

↓

1-CACHE miss + D-CACHE miss

↑ ↑

2 2

$0.02 \times 100 = 2$        $0.36 \times 100 = 36$

$= 2 + 36 = 38$

$= 5.44$

$\frac{\text{Perf}}{\text{Perf A}} = \frac{5.44}{2} = 2.72$

44

## Average Access Time

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

$$1 + 20 \times 0.05$$

$$= 2 + 1\text{ns}$$

$$= 2\text{ns}.$$

45

## 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
- 

40

## Measuring Cache Performance

- Components of CPU time  $= (\text{CPU exec clock cycles}) + (\text{Memory stall cycles}) \times \text{Clock Cycle Time}$ 
  - Program execution cycles
    - Includes cache hit time
  - Memory stall cycles
    - Mainly from cache misses
- With simplifying assumptions:
 
$$\begin{aligned} \text{Read stall cycles} &= \frac{\text{Reads}}{\text{Prog}} \times \text{Misses} \times \text{Miss Penalty} \\ \text{Memory stall cycles} &= \frac{\text{Writes}}{\text{Prog}} \times \text{Misses} \times \text{Miss Penalty} \\ &= \frac{\text{Memory accesses}}{\text{Program}} \times \text{Miss rate} \times \text{Miss penalty} \\ &= \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Misses}}{\text{Instruction}} \times \text{Miss penalty} \end{aligned}$$

43

## 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

46