

## Chapter 5

**Large and Fast: Exploiting  
Memory Hierarchy**

# Recap

---

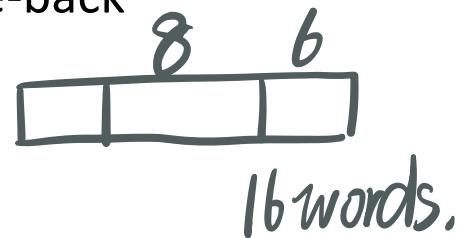
- Memory hierarchy
- Storage technologies
- Direct mapped cache
- Write policies

# Write Policies Summary

- If that memory location is in the cache?
  - ◆ Send it to the cache
  - ◆ Should we also send it to memory right away?  
*(write-through policy)*
  - ◆ Wait until we kick the block out *(write-back policy)*
- If it is not in the cache?
  - ◆ Allocate the line (put it in the cache)?  
*(write allocate policy)*
  - ◆ Write it directly to memory without allocation?  
*(no write allocate policy)*

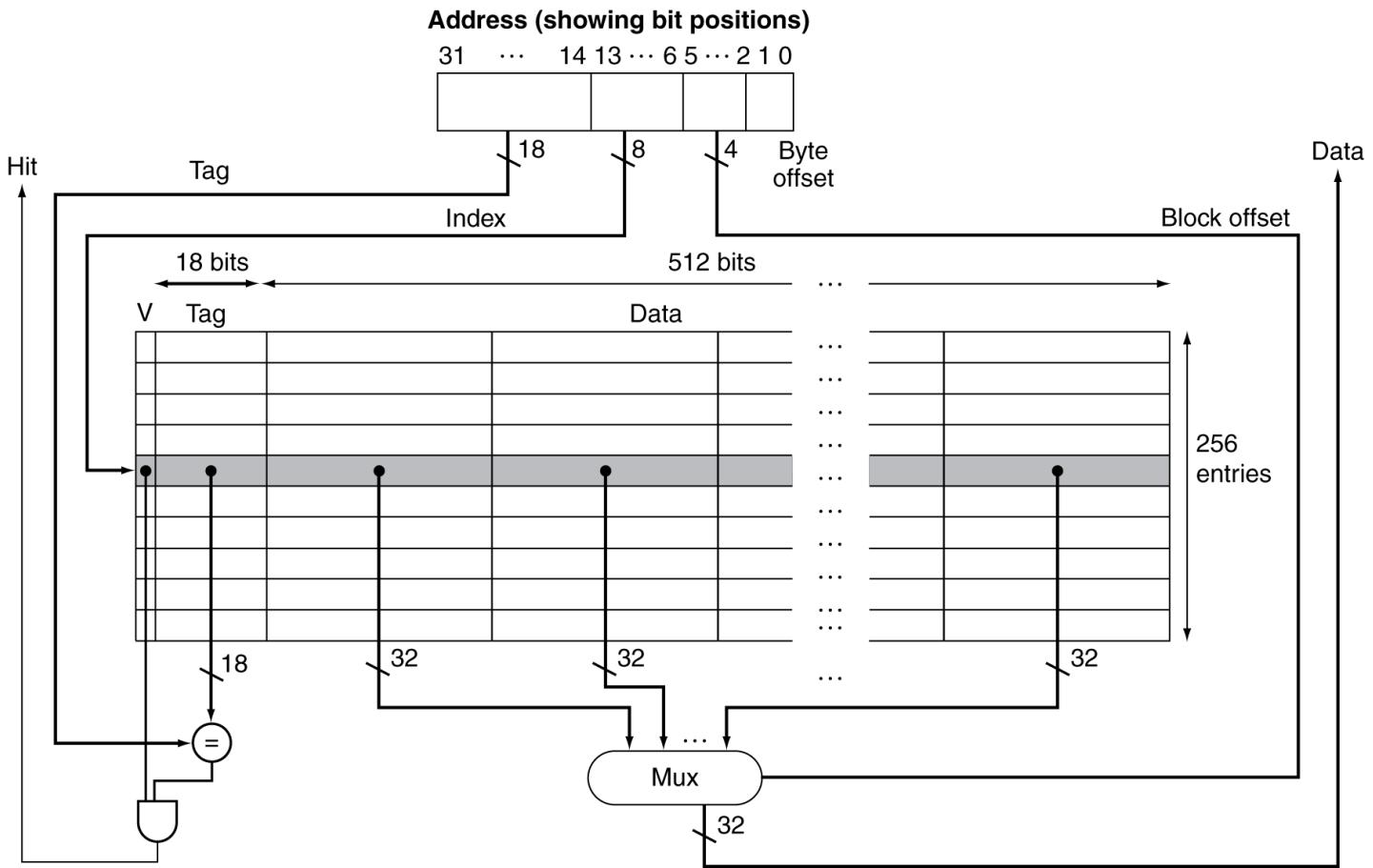
# Example: Intrinsity FastMATH

- Embedded MIPS processor
  - ◆ 12-stage pipeline
  - ◆ Instruction and data access on each cycle
- Split cache: separate I-cache and D-cache
  - ◆ Each 16KB:  $2^8 \times 2^4 \times 2^2$ .  
Each 16KB: 256 blocks  $\times$  16 words/block
  - ◆ D-cache: write-through or write-back
- SPEC2000 miss rates
  - ◆ I-cache: 0.4%      *randomness.*
  - ◆ D-cache: 11.4%
  - ◆ Weighted average: 3.2%



一个地址线是一个Byte .

# 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

- Calculate actual CPI, given that
  - ◆ 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 (assume  $N$  ins. In total)
  - ◆ I-cache:  $N \times 0.02 \times 100/N = 2$
  - ◆ D-cache:  $N \times 0.36 \times 0.04 \times 100/N = 1.44$
- Actual CPI =  $2 + 2 + 1.44 = 5.44$ 
  - ◆ Ideal CPU is  $5.44/2 = 2.72$  times faster

# Average Access Time

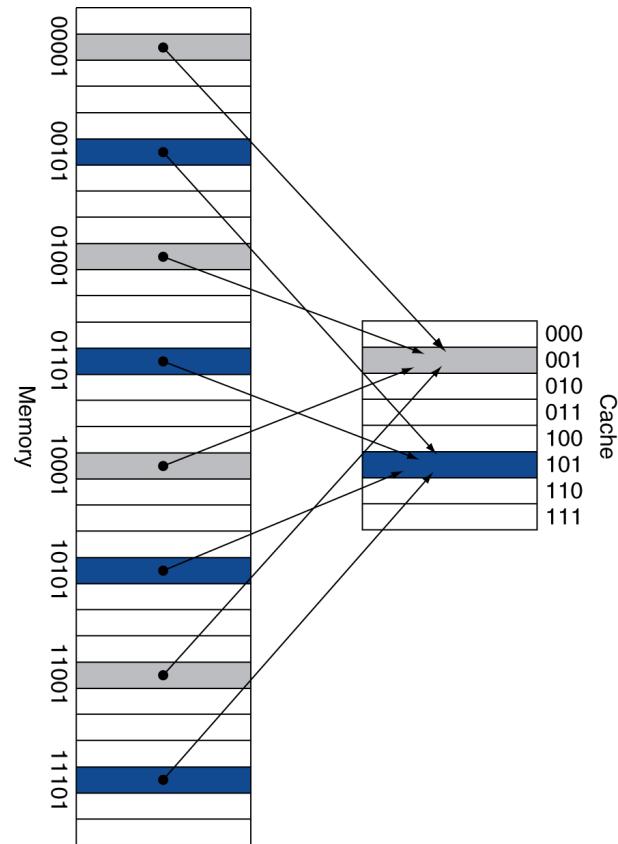
- Hit time is also important for performance
- Average memory access time (AMAT)
  - ◆  $\text{AMAT} = \text{Hit time} + \text{Miss rate} \times \text{Miss penalty}$
- Example
  - ◆ CPU with 1ns clock, hit time = 1 cycle, miss penalty = 20 cycles, I-cache miss rate = 5%  
$$95\% \times 1 + 5\% \times (1+20)$$
  - ◆  $\text{AMAT} = 1 + 0.05 \times 20 = 2\text{ns}$ 
    - 2 cycles per instruction

# Performance Summary

- When CPU performance increased
  - ◆ Miss penalty becomes more significant
  - ◆ CPI=2, Miss=3.44, % of memory stall:  $3.44/5.44=63\%$
  - ◆ CPI=1, Miss=3.44, % of memory stall:  $3.44/4.44=77\%$
- 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

# Recall: Direct Mapped Cache

- Location determined by address
- Direct mapped: only one choice
  - ◆ Capacity of cache is not fully exploited
  - ◆ Miss rate is high



# 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</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		

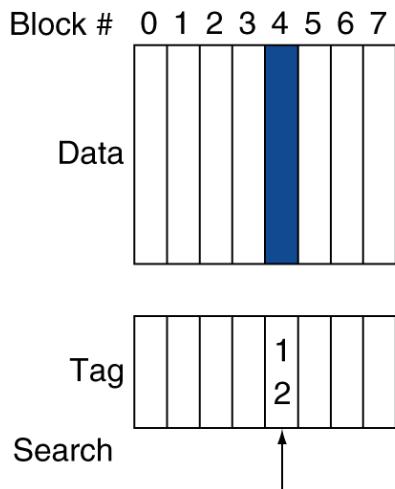
# Cache Example

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

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 Cache Example

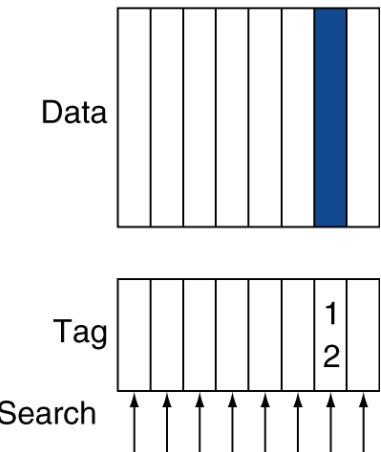
Direct mapped



Set associative



Fully associative



比較次數增加。

# 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)

# Spectrum of Associativity

- For a cache with 8 blocks

## 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

# 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

0-8-0-6-8

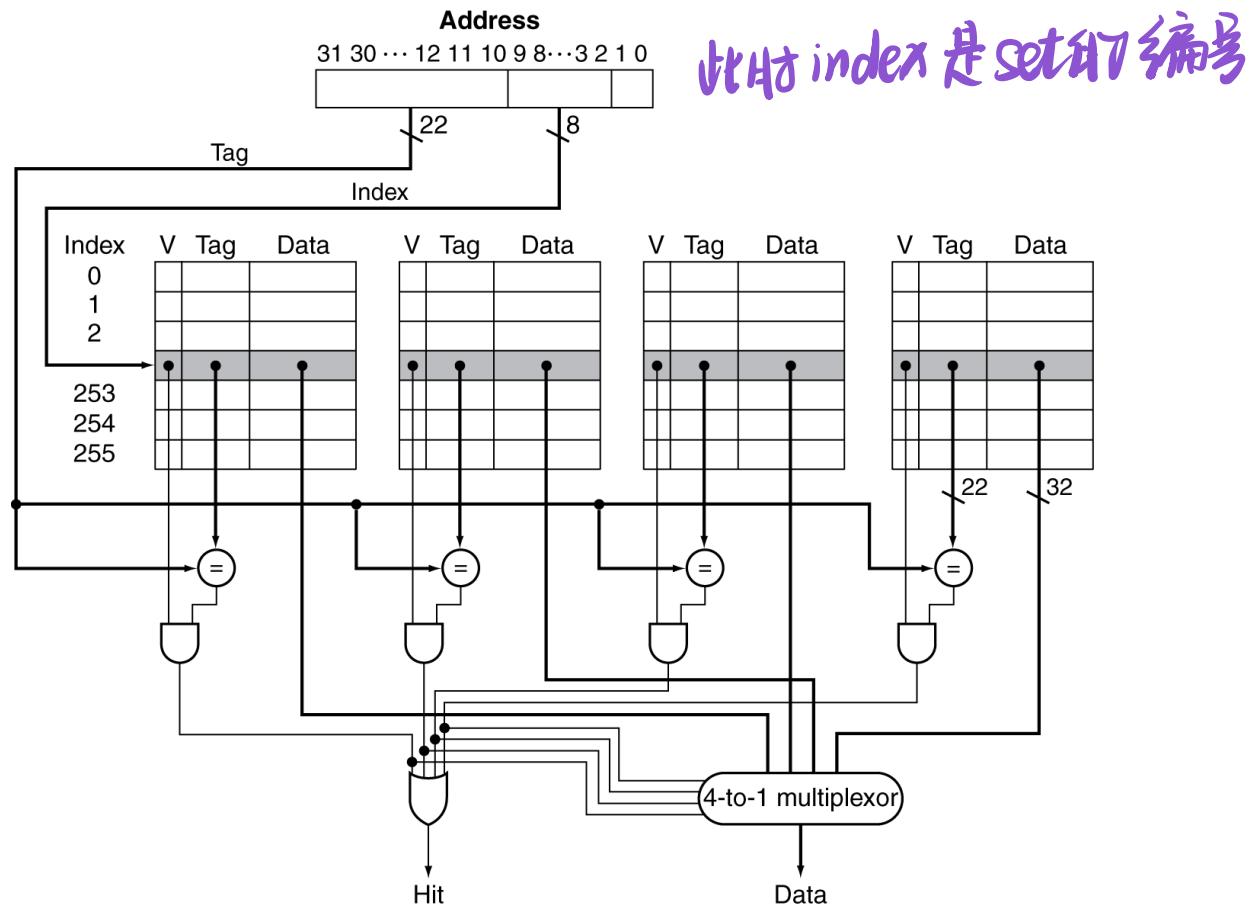
- 2-way set associative

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			
			Mem[0]	Mem[8]	Mem[8]	Mem[6]
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]	

# Set Associative Cache Organization

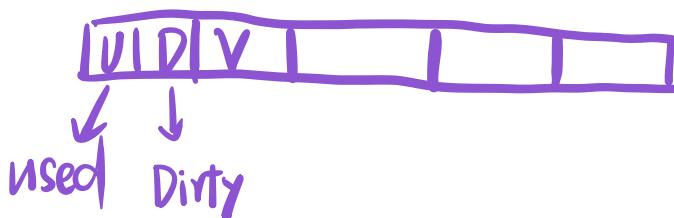


# How Much Associativity

- 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% (highlighted)
  - ◆ 8-way: 8.1%

# 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



# 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

reg → cache 1

# Multilevel Cache Example

- Given
  - ◆ CPU base CPI = 1, clock rate = 4GHz
  - ◆ Miss rate/instruction = 2%
  - ◆ Main memory access time = 100ns
- With just primary cache
  - ◆ Miss penalty =  $100\text{ns}/0.25\text{ns} = 400$  cycles
  - ◆ Effective CPI =  $1 + 0.02 \times 400 = 9$

# Example (cont.)

- Now add L-2 cache

- Access time = 5ns
- Global miss rate to main memory =  $0.5\%$

- Primary miss with L-2 hit

- Penalty =  $5\text{ns}/0.25\text{ns} = 20 \text{ cycles}$

- Primary miss with L-2 miss

- Extra penalty = 400 cycles

$$\text{CPI} = 1 + 0.02 \times 20 + 0.005 \times 400 = 3.4$$

$$\text{Performance ratio} = 9/3.4 = 2.6$$

if  $L_2 \text{ missrate} = 95\%$

$$\text{CPI} = 1 + L_1 \text{ miss} (L_2 \text{ hit time} + L_2 \text{ miss} \times \text{mem time})$$

# Multilevel Cache Considerations

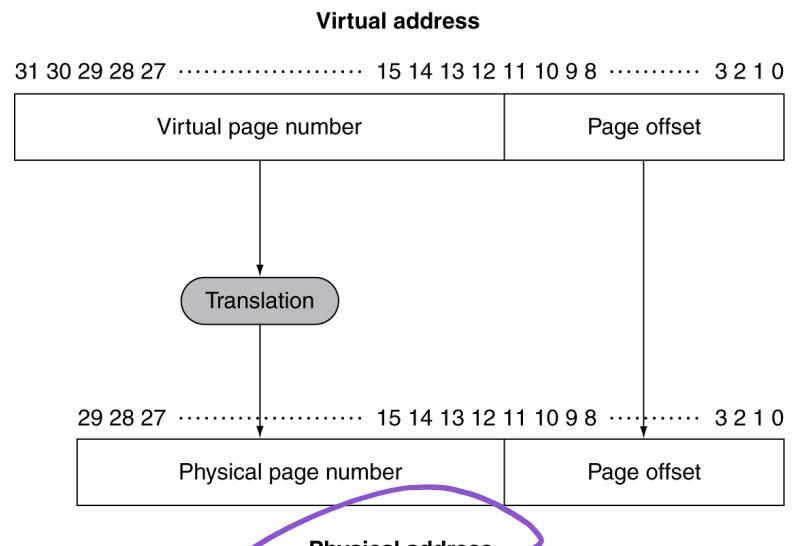
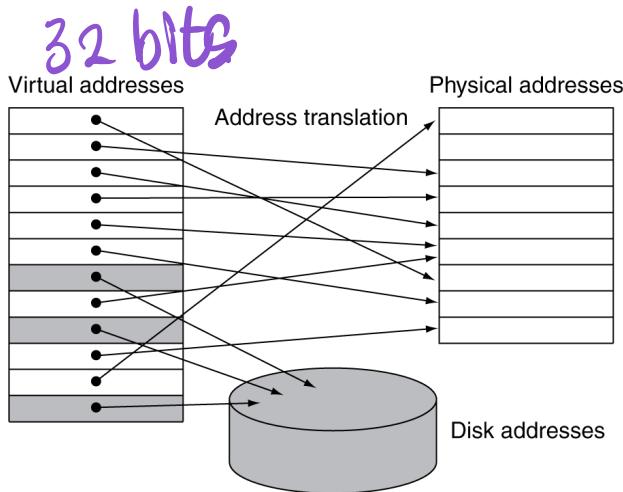
- Primary cache
  - ◆ Focus on minimal hit time
- L-2 cache
  - ◆ Focus on low miss rate to avoid main memory access
  - ◆ Hit time has less overall impact
- Results
  - ◆ L-1 cache usually smaller than a single cache
  - ◆ L-1 block size smaller than L-2 block size

# Virtual Memory 32 bit

- 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 “miss” is called a page fault

# Address Translation

- Fixed-size pages (e.g., 4K)  $\frac{\text{byte}}{(10+2)} = 12$



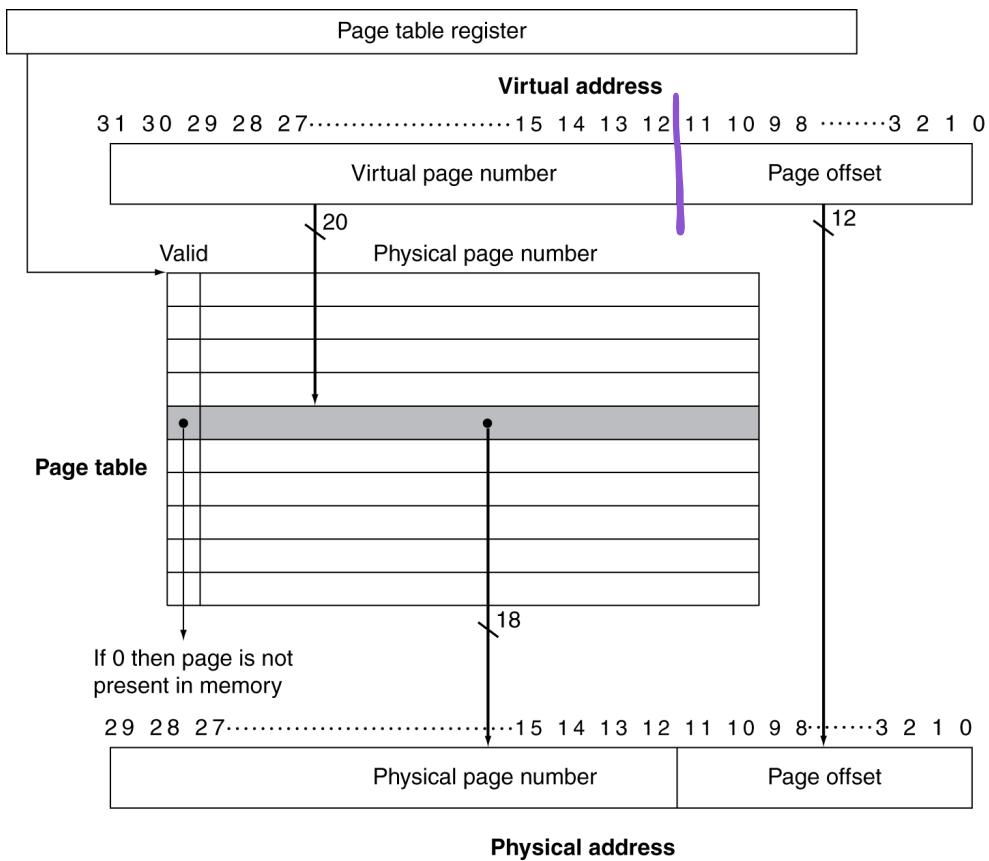
# Page Fault Penalty

- 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

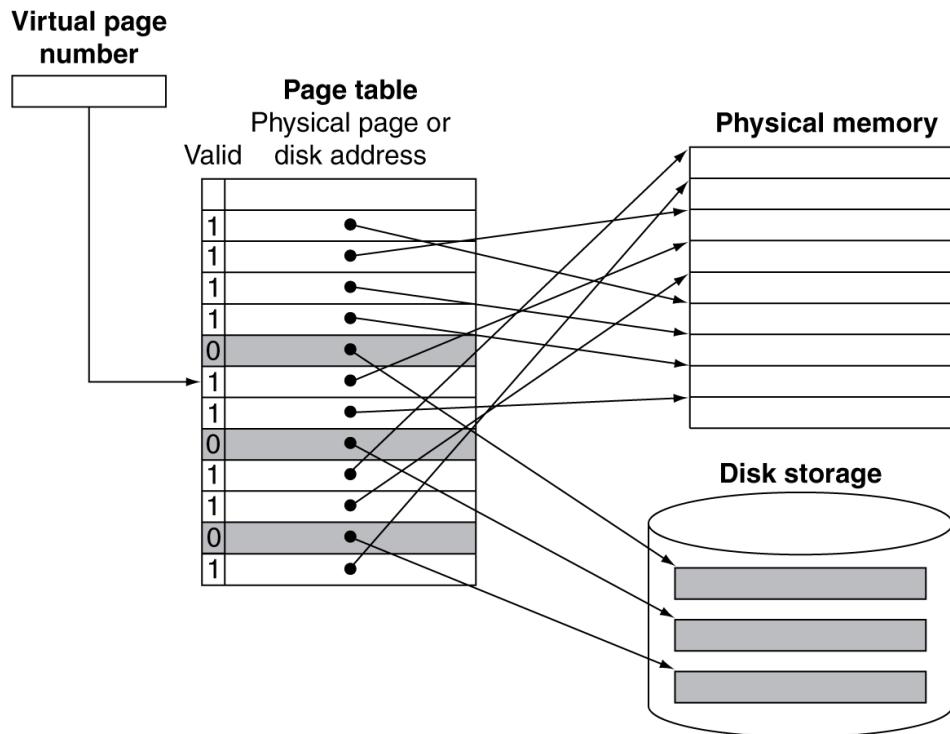
# Page Tables

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

# Translation Using a Page Table



# Mapping Pages to Storage



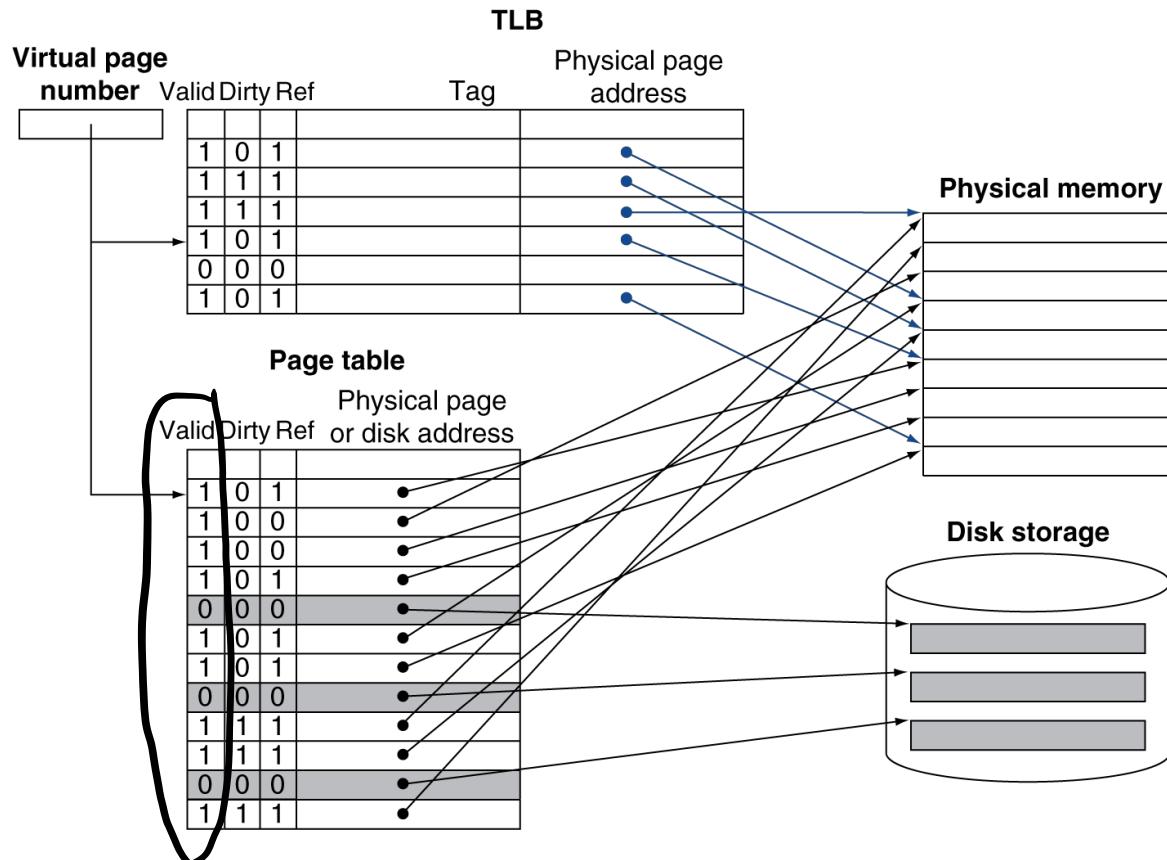
# Replacement and Writes

- 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
  - ◆ Use write-back, because write through is impractical
  - ◆ Dirty bit in PTE set when page is written

# Fast Translation Using a TLB

- Since page table is in memory, every memory access by a program requires two memory accesses
  - ◆ One to access the page table entry
  - ◆ Then the actual memory access
- Can we move the page table to CPU?
  - ◆ Yes, use a fast cache in CPU to store recently used PTEs, because access to page tables has good locality
  - ◆ Called a Translation Look-aside Buffer (TLB)  
*cache for pagetable*
  - ◆ Typical: 16–512 PTEs, 0.5–1 cycle for hit, 10–100 cycles for miss, 0.01%–1% miss rate  
*pagetable entry*
  - ◆ Misses could be handled by hardware or software

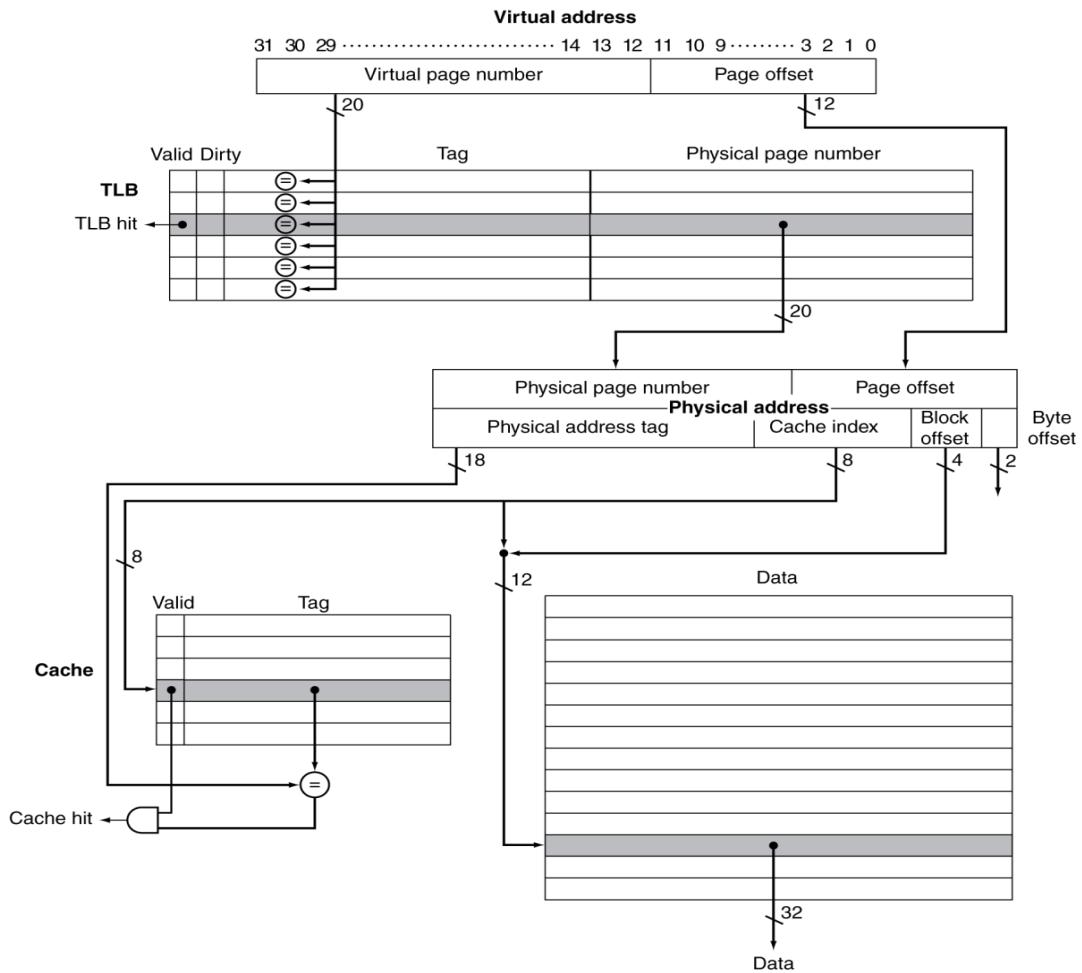
# Fast Translation Using a TLB



# 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
  - ◆ Then restart the faulting instruction

# TLB and Cache Interaction



# Memory Protection

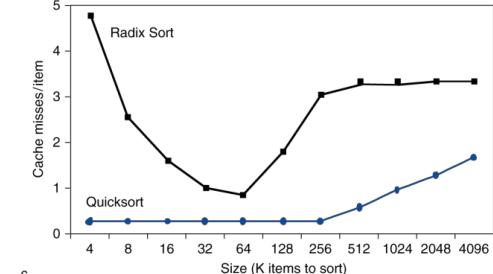
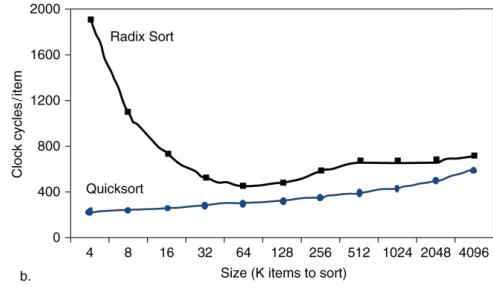
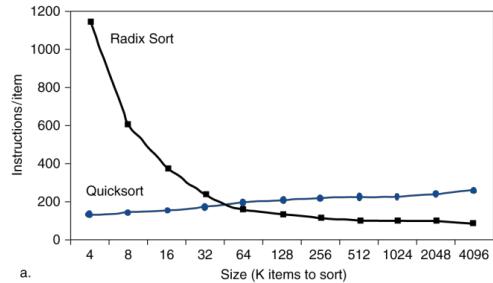
- Different tasks can share parts of their virtual address spaces
  - ◆ But need to protect against errant access
  - ◆ Requires OS assistance
- Hardware support for OS protection
  - ◆ Privileged supervisor mode (aka kernel mode)
  - ◆ Privileged instructions
  - ◆ Page tables and other state information 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

# Interactions with Software

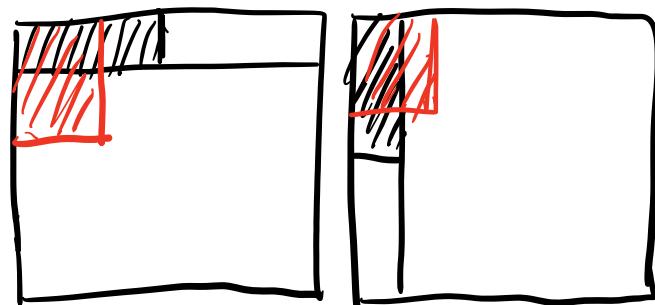
- Compare two algorithms:  
Radix sort & Quicksort
- When size is large,
  - ◆ Radix sort has less instructions
  - ◆ But quicksort has less clock cycles
  - ◆ Because miss rate of radix sort is higher



# Software Optimization via Blocking

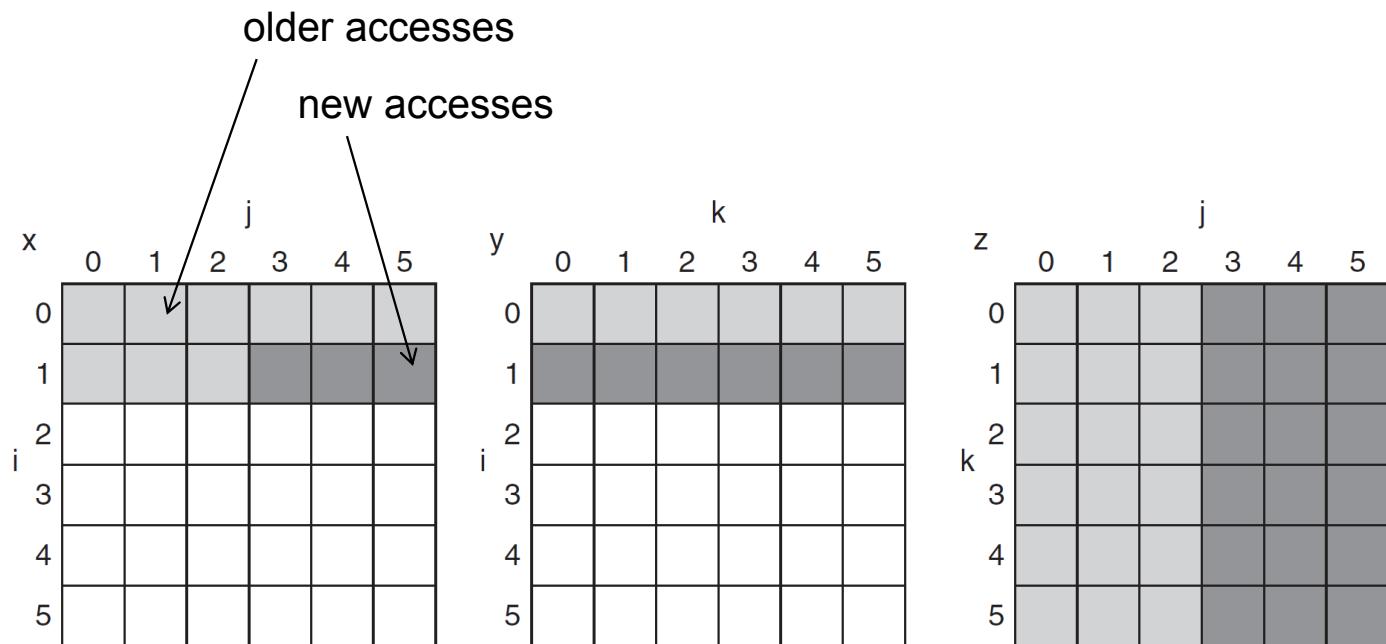
- Goal: maximize accesses to data before it is replaced
- Consider inner loops of DGEMM:

```
for (int j = 0; j < n; ++j)
{
    double cij = C[i+j*n];
    for( int k = 0; k < n; k++ )
        cij += A[i+k*n] * B[k+j*n];
    C[i+j*n] = cij;
}
```

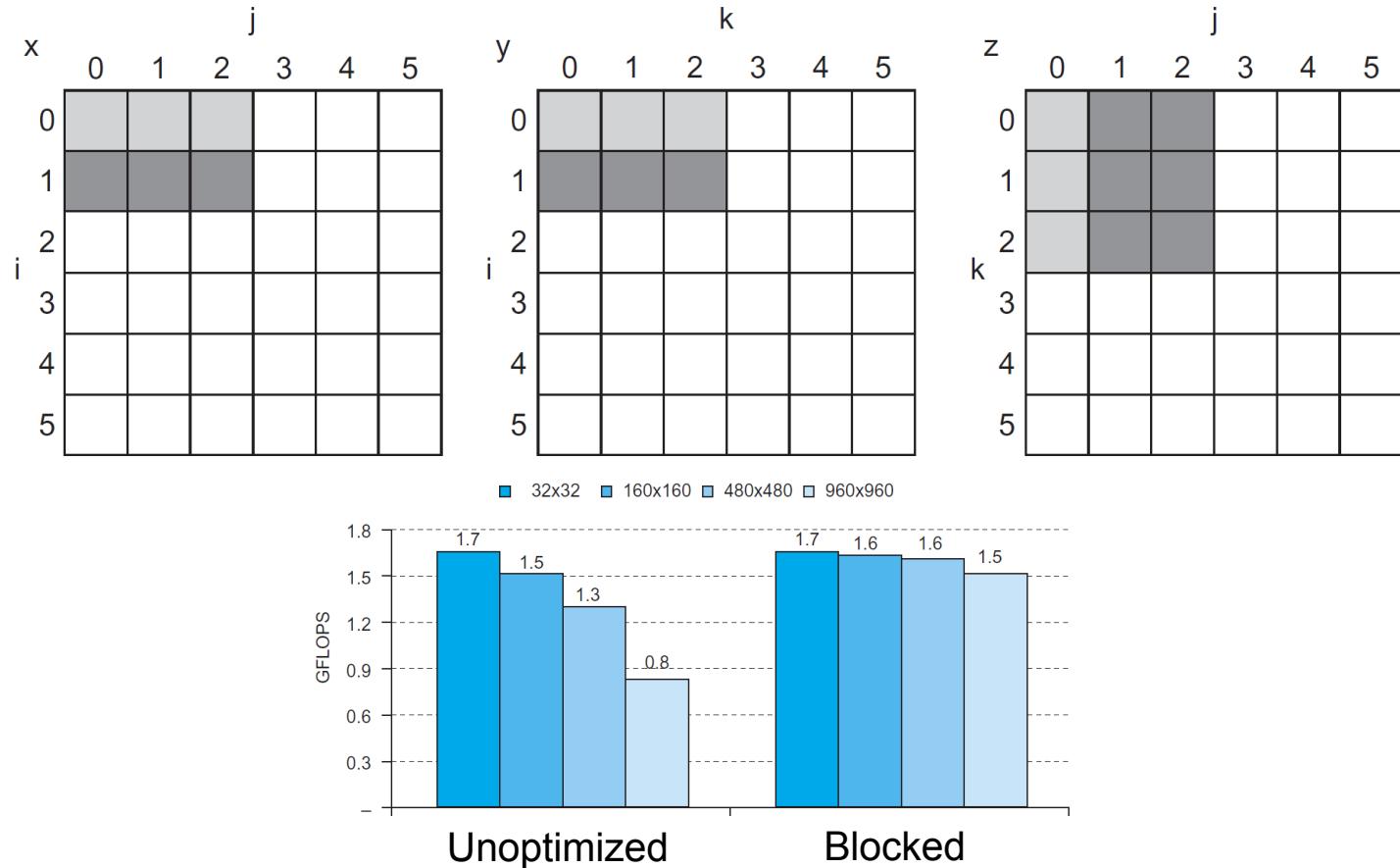


# DGEMM Access Pattern

- C, A, and B arrays



# Blocked DGEMM Access Pattern



# 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

# Block Placement

- Determined by associativity
  - ◆ Direct mapped (1-way associative)
    - One choice for placement
  - ◆ n-way set associative
    - n choices within a set
  - ◆ Fully associative
    - Any location
- Higher associativity reduces miss rate
  - ◆ Increases complexity, cost, and access time

# Finding a Block

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

- Virtual memory
  - ◆ Full table lookup makes full associativity feasible
  - ◆ Benefit in reduced miss rate
- Cache and TLB
  - ◆ Set-associative, some cache uses direct map

# Replacement

- Choice of entry to replace on a miss
  - ◆ Least recently used (LRU)
    - Complex and costly hardware for high associativity
  - ◆ Random
    - Close to LRU, easier to implement
- Virtual memory
  - ◆ LRU approximation with hardware support
- Cache
  - ◆ Both LRU and random is ok

# Write Policy

- Write-through
  - ◆ Update both upper and lower levels
  - ◆ Simplifies replacement, but may require write buffer
- Write-back
  - ◆ Update upper level only
  - ◆ Update lower level when block is replaced
  - ◆ Need to keep more state
- Virtual memory
  - ◆ Only write-back is feasible, given disk write latency

# Sources of Misses

- Compulsory misses (aka cold start misses)
  - ◆ First access to a block
- Capacity misses
  - ◆ Due to finite cache size
  - ◆ A replaced block is later accessed again
- Conflict misses (aka collision misses)
  - ◆ In a non-fully associative cache
  - ◆ Due to competition for entries in a set
  - ◆ Would not occur in a fully associative cache of the same total size

# Cache Design Trade-offs

Design change	Effect on miss rate	Negative performance effect
Increase cache size		
Increase associativity		
Increase block size		

# Cache Design Trade-offs

Design change	Effect on miss rate	Negative performance effect
Increase cache size	Decrease capacity misses	May increase access time
Increase associativity		
Increase block size		

# Cache Design Trade-offs

Design change	Effect on miss rate	Negative performance effect
Increase cache size	Decrease capacity misses	May increase access time
Increase associativity	Decrease conflict misses	May increase access time
Increase block size		

# Cache Design Trade-offs

Design change	Effect on miss rate	Negative performance effect
Increase cache size	Decrease capacity misses	May increase access time
Increase associativity	Decrease conflict misses	May increase access time
Increase block size	Decrease compulsory misses	<p>Increases <u>miss penalty</u>. For very large block size, <u>more write</u>, may increase miss rate due to pollution. <u>conflict miss rate ↑</u></p>

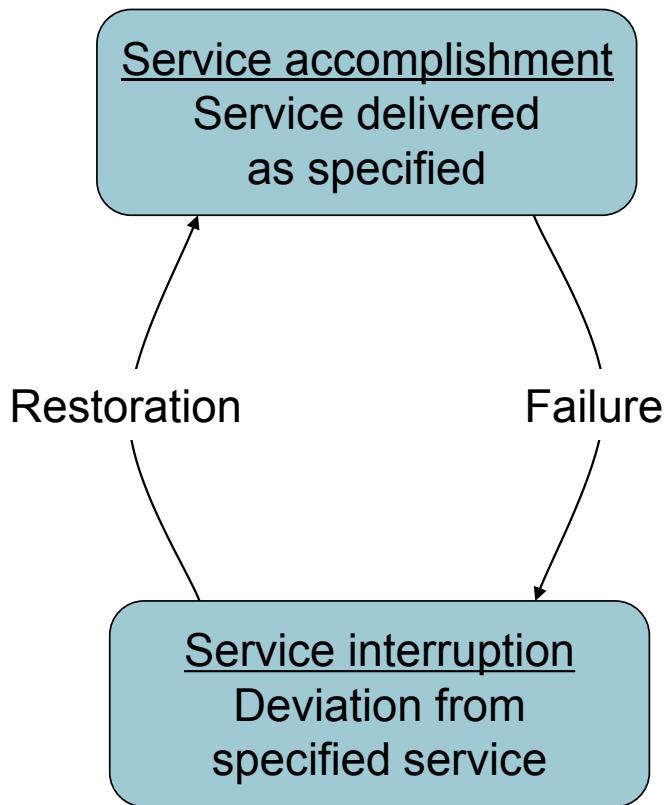
# Multilevel On-Chip Caches

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

# 2-Level TLB Organization

Characteristic	ARM Cortex-A8	Intel Core i7
Virtual address	32 bits	48 bits
Physical address	32 bits	44 bits
Page size	Variable: 4, 16, 64 KiB, 1, 16 MiB	Variable: 4 KiB, 2/4 MiB
TLB organization	<p>1 TLB for instructions and 1 TLB for data</p> <p>Both TLBs are <u>fully associative, with 32 entries</u>, round robin replacement</p> <p>TLB misses handled in hardware</p>	<p>1 TLB for instructions and 1 TLB for data per core</p> <p>Both L1 TLBs are <u>four-way set associative, LRU replacement</u></p> <p>L1 I-TLB has 128 entries for small pages, 7 per thread for large pages</p> <p>L1 D-TLB has 64 entries for small pages, 32 for large pages</p> <p>The <u>L2 TLB</u> is <u>four-way set associative, LRU replacement</u></p> <p>The L2 TLB has 512 entries</p> <p>TLB misses handled in hardware</p>

# Dependability 可信度



- Fault: failure of a component
  - ◆ May or may not lead to system failure

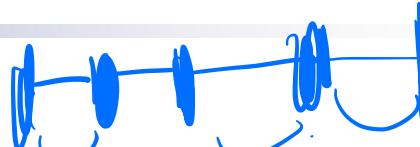


# Dependability Measures

- Reliability: mean time to failure (MTTF)
- Service interruption: mean time to repair (MTTR)
- Mean time between failures
  - ◆  $MTBF = MTTF + MTTR$
- Availability =  $MTTF / (MTTF + MTTR)$
- Improving Availability
  - ◆ Increase MTTF: fault avoidance, fault tolerance, fault forecasting
  - ◆ Reduce MTTR: fault detection, fault diagnosis and fault repair

# The Hamming SEC Code

三



## Hamming distance

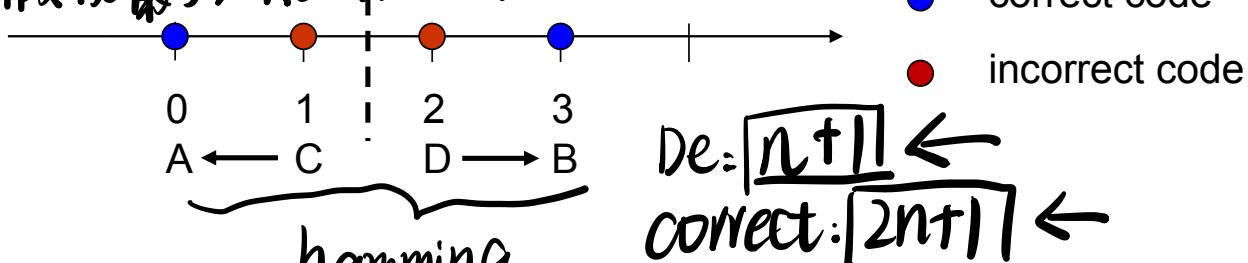
- ◆ Number of bits that are different between two bit patterns
- ◆ E.g. use 111 to represent 1, use 000 to represent 0, hamming distance (d) is 3, d=3. 比特位不同的位数.

## Minimum distance = 2 provides single bit error detection

- ◆ E.g. odd-parity code:  $10 \rightarrow 101$ ,  $11 \rightarrow 110$ ,  $d = 2$

## Minimum distance = 3 provides single error correction(SEC), 2-bit/ double error detection (DED)

假设最多只能错1位。



De:  $n+1$  ←  
correct:  $2n+1$  ←

# Encoding SEC

n位纠错

- To calculate Hamming code:

- Number bits from 1 on the left
- All bit positions that are a power of 2 are parity bits (bit 1 2 4 8 are parity bits)
- Each parity bit checks certain data bits:

Bit position	1	2	3	4	5	6	7	8	9	10	11	12
Encoded date bits	0	1	1	1	0	0	1	0	1	0	1	0
Parity bit coverage	p1	X		X		X		X		X		X
p2		X	X			X	X			X	X	
p4				X	X	X						X
p8								X	X	X	X	X

# Decoding SEC

- Value of parity bits indicates which bits are in error
  - Use numbering from encoding procedure
  - E.g.
    - Parity bits = 0000 indicates no error
    - Parity bits = 0101 indicates bit 10 was flipped

Bit position	1	2	3	4	5	6	7	8	9	10	11	12
Encoded date bits	0	1	1	1	0	0	1	0	1	1	1	0
Parity bit coverage	p1	X		X	X		X		X		X	
p2		X	X			X	X			X	X	
p4				X	X	X	X					X
p8								X	X	X	X	X

✓ 0  
✗ 1  
✓ 0  
✗ 1

# SEC/DED Code

1 bit error

- Add an additional parity bit for the whole word ( $p_n$ )
- Make Hamming distance = 4
- Decoding:
  - ◆ Let  $H$  = SEC parity bits
    - $H$  even,  $p_n$  even, no error
    - $H$  odd,  $p_n$  odd, correctable single bit error ✓
    - $H$  even,  $p_n$  odd, error in  $p_n$  bit
    - $H$  odd,  $p_n$  even, double error occurred ✓
- Note: ECC DRAM uses SEC/DED with 8 bits protecting each 64 bits



ECC Error Correcting Code

# Summary

- Cache Performance
  - ◆ Mainly depends on miss rate and miss penalty
- To improve cache performance:
  - ◆ Fully associative cache
  - ◆ Set-associative cache
  - ◆ Replacement policy
  - ◆ Multilevel cache
- Dependability
  - ◆ MTTF, MTTR, reliability, availability
  - ◆ Hamming code: SEC/DED code

# Virtual Machines

- Host computer emulates guest operating system and machine resources
  - ◆ Improved isolation of multiple guests
  - ◆ Avoids security and reliability problems
  - ◆ Aids sharing of resources
- Virtualization has some performance impact
  - ◆ Feasible with modern high-performance computers
- Examples
  - ◆ IBM VM/370 (1970s technology!)
  - ◆ VMWare
  - ◆ Microsoft Virtual PC

# Virtual Machine Monitor

- Maps virtual resources to physical resources
  - ◆ Memory, I/O devices, CPUs
- Guest code runs on native machine in user mode
  - ◆ Traps to VMM on privileged instructions and access to protected resources
- Guest OS may be different from host OS
- VMM handles real I/O devices
  - ◆ Emulates generic virtual I/O devices for guest

# Example: Timer Virtualization

- In native machine, on timer interrupt
  - ◆ OS suspends current process, handles interrupt, selects and resumes next process
- With Virtual Machine Monitor
  - ◆ VMM suspends current VM, handles interrupt, selects and resumes next VM
- If a VM requires timer interrupts
  - ◆ VMM emulates a virtual timer
  - ◆ Emulates interrupt for VM when physical timer interrupt occurs

# Instruction Set Support

- User and System modes
- Privileged instructions only available in system mode
  - ◆ Trap to system if executed in user mode
- All physical resources only accessible using privileged instructions
  - ◆ Including page tables, interrupt controls, I/O registers
- Renaissance of virtualization support
  - ◆ Current ISAs (e.g., x86) adapting

# 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  $\leftrightarrow$  L2 cache  $\leftrightarrow \dots \leftrightarrow$  DRAM memory  
 $\leftrightarrow$  disk
- Virtual Memory and TLB

# Practice

---

- Exercise 5.7, 12, 5.13.