# Chapter 5

## Large and Fast: Exploiting Memory Hierarchy

# Principle of Locality

- Programs access a small proportion of their address space at any time

-

  - Items accessed recently are likely to be accessed again soon

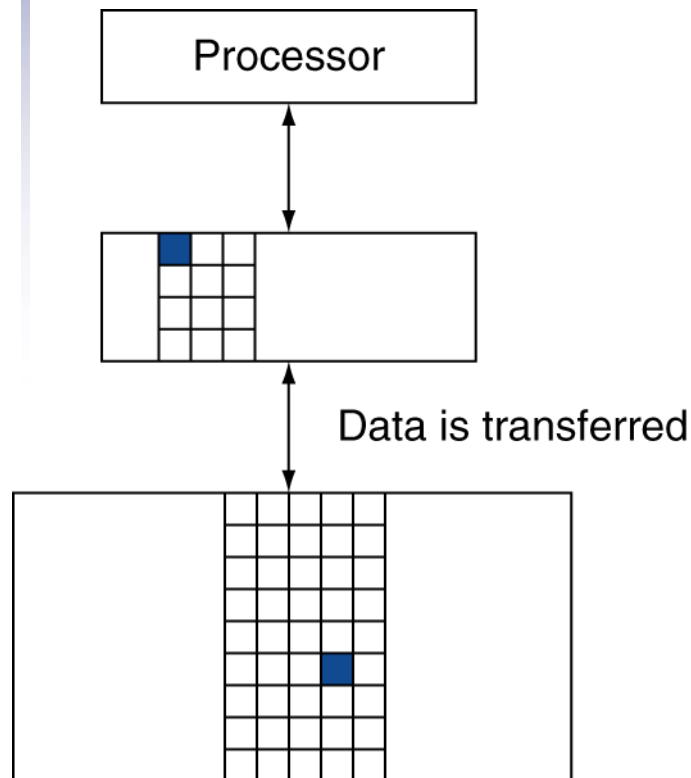  - e.g., instructions in a loop, induction variables

-

  - Items near those accessed recently are likely to be accessed soon

  - E.g., sequential instruction access, array data

# Taking Advantage of Locality

- Memory hierarchy
- Store everything on [          ]
- Copy recently accessed (and nearby) items from disk to smaller DRAM memory
  - [                    ]
- Copy more recently accessed (and nearby) items from DRAM to smaller SRAM memory
  - [                    ] attached to CPU

# Memory Hierarchy Levels

Processor

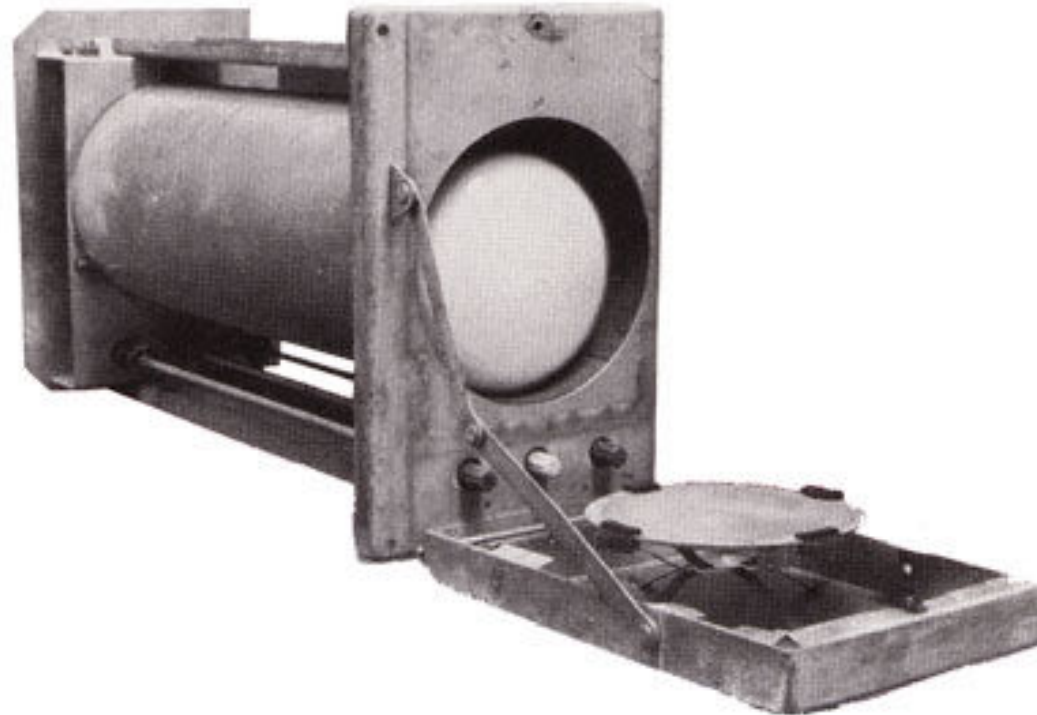Data is transferred

- Block (aka line): unit of copying
  - May be multiple words
- If accessed data is present in upper level
  - [    ] access satisfied by upper level
    - Hit ratio: hits/accesses
- If accessed data is absent
  - [    ] block copied from lower level
    - Time taken: miss penalty
    - Miss ratio: misses/accesses = 1 – 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
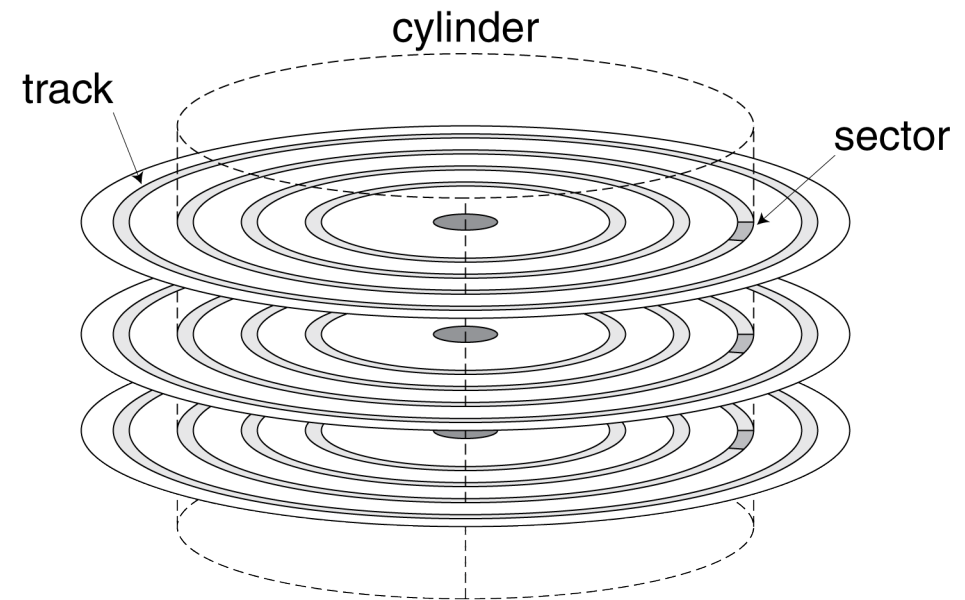    - Capacity and cost/GB of disk

# Memory Technology

- IBM 701 adopted 72 3-inch **Williams-Kilburn tubes** to realize 2048 36-bit words.

# Disk Storage

- Non-volatile, rotating magnetic storage

cylinder

track

sector

# Disk Sectors and Access

Gap  Sync  Address Mark

←— ECC: 40 X10 bit symbols = 50 Bytes

- ## 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
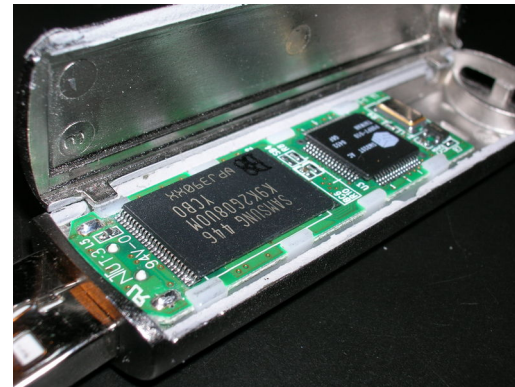
# Disk Access Example

- ## Given
  - 512B sector, 15,000rpm, 4ms average seek time, 100MB/s transfer rate, 0.2ms controller overhead, idle disk

- ## Average read time
  - 4ms seek time
    + ½ / (15,000/60) = 2ms rotational latency
    + 512 / 100MB/s = 0.005ms transfer time
    + 0.2ms controller delay
    = 6.2ms

- ## If actual average seek time is 1ms
  - Average read time = 3.2ms

# 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 allocates physical sectors on disk
  - Present logical sector interface to host
  - SCSI (Small Computer Systems Interface), ATA (Advanced Technology Attachment), SATA (Serial ATA)
- Disk drives include caches (i.e., buffers)
  - Prefetch sectors in anticipation of access
  - Avoid seek and rotational delay

# Flash Storage

- Non-volatile semiconductor storage
    - 100× – 1000× 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
- NAND flash: bit cell like a NAND gate
  - Denser (bits/area), but block-at-a-time (seq.) access
  - Cheaper per GB
  - Used for USB keys, media storage, …
- Flash bits wears out after 10000 ~ 100000 writes
  - Not suitable for direct RAM or disk replacement
  - Wear leveling: remap data to less used blocks

# DRAM Organization

- Bits in a DRAM are organized as a rectangular array
    - DRAM accesses an entire row
    - Burst mode: supply successive words from a row with reduced latency

- Double data rate (DDR) DRAM
    - Transfer on rising and falling clock edges

- Quad data rate (QDR) DRAM
    - Separate DDR inputs and outputs

# DRAM Generations

| Year | Capacity | $/GB |
|------|----------|------|
| 1980 | 64Kbit | $1500000 |
| 1983 | 256Kbit | $500000 |
| 1985 | 1Mbit | $200000 |
| 1989 | 4Mbit | $50000 |
| 1992 | 16Mbit | $15000 |
| 1996 | 64Mbit | $10000 |
| 1998 | 128Mbit | $4000 |
| 2000 | 256Mbit | $1000 |
| 2004 | 512Mbit | $250 |
| 2007 | 1Gbit | $50 |

# Cache Memory

- Cache memory
  - The level of the memory hierarchy closest to the CPU

- Given accesses $X_1, \ldots, X_{n-1}, X_n$

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

a. Before the reference to $X_n$

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

b. After the reference to $X_n$

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

# Direct Mapped Cache

- Location determined by address
- Direct mapped: only one choice
  - (Block address) modulo (#Blocks in cache)
  - *N*-way associative: *N* choices

Cache

| tag | index | offset |
|-----|-------|--------|

| tag | index | W.O. | B.O. |
|-----|-------|------|------|

- #Blocks is a power of 2
- Use low-order address bits

00001  00101  01001  01101  10001  10101  11001  11101

Memory

block address: address to a block

# 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

# Cache Example

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

- Initial state

| Index | V | Tag | Data |
|-------|---|-----|------|
| 000 | N | | |
| 001 | N | | |
| 010 | N | | |
| 011 | N | | |
| 100 | N | | |
| 101 | N | | |
| 110 | N | | |
| 111 | N | | |

# Cache Example

| Word addr | Binary addr | Hit/miss | Cache block |
|-----------|-------------|----------|-------------|
| 22 | 10 110 | Miss | 110 |

| Index | V | Tag | Data |
|-------|---|-----|------|
| 000 | N | | |
| 001 | N | | |
| 010 | N | | |
| 011 | N | | |
| 100 | N | | |
| 101 | N | | |
| 110 | Y | 10 | Mem[10110] |
| 111 | N | | |

# Cache Example

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

| Index | V | Tag | Data |
|-------|---|-----|------|
| 000 | N | | |
| 001 | N | | |
| 010 | Y | 11 | Mem[11010] |
| 011 | N | | |
| 100 | N | | |
| 101 | N | | |
| 110 | Y | 10 | Mem[10110] |
| 111 | N | | |

# 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 | | |
| 010 | Y | 11 | Mem[11010] |
| 011 | N | | |
| 100 | N | | |
| 101 | N | | |
| 110 | Y | 10 | Mem[10110] |
| 111 | N | | |

# 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 |
|-------|---|-----|------|
| 000 | Y | 10 | Mem[10000] |
| 001 | N | | |
| 010 | Y | 11 | Mem[11010] |
| 011 | Y | 00 | Mem[00011] |
| 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 | | |
| **010** | **Y** | **10** | **Mem[10010]** |
| 011 | Y | 00 | Mem[00011] |
| 100 | N | | |
| 101 | N | | |
| 110 | Y | 10 | Mem[10110] |
| 111 | N | | |

# Address Subdivision

**Address (showing bit positions)**

31 30 · · · 13 12 11 · · · 2 1 0

Tag /20  Index /10  Byte offset

Hit

Data

Index | Valid | Tag | Data
0
1
2
…

…
…
1021
1022
1023

/20  /32

= 

1K word data

- 32-bit byte address, a direct-mapped cache, cache size $2^n$ blocks, block data size $2^m$ words ($2^{m+2}$ bytes)
  - Tag size: $32 - (n+m+2)$
  - Total number of bits:

$$2^n \times (2^m \times 32 + (32-n-m-2)+1) =$$
$$2^n \times (2^m \times 32 + 31 - n - m)$$
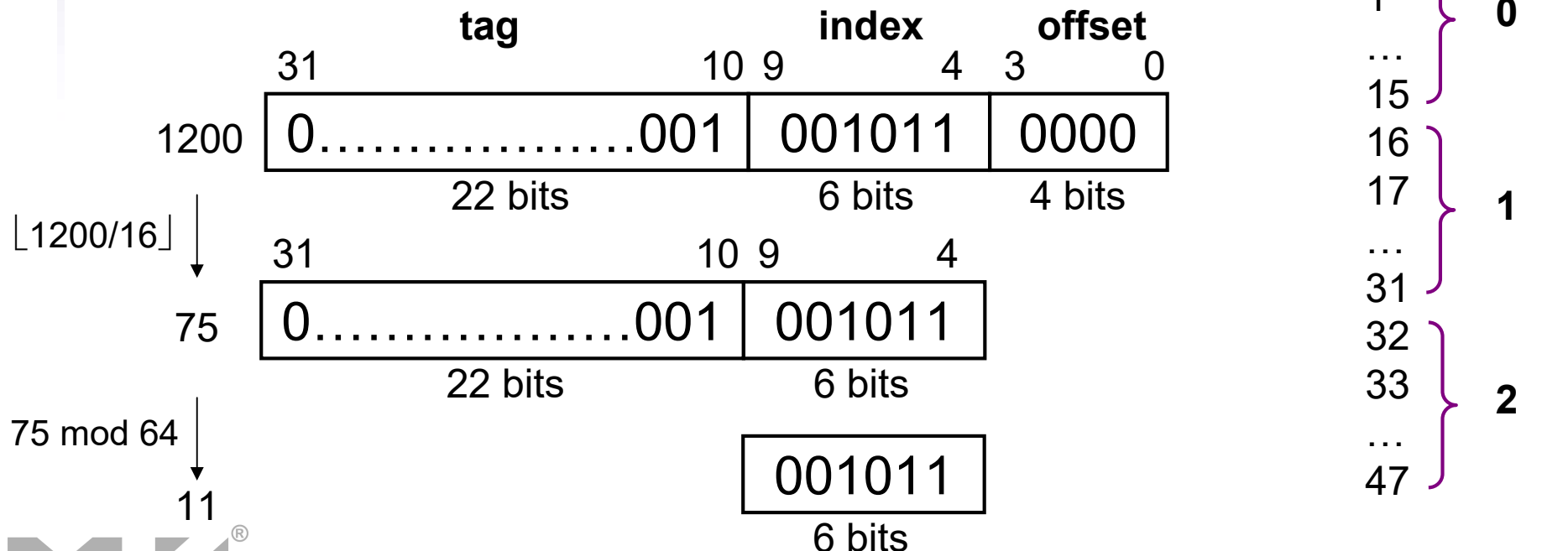
- 16KB data, 4-word blocks, assume a 32-bit address

$$2^{10} \times (4 \times 32 + (32-10-2-2)+1) =$$
$$2^{10} \times 147 = 147\,Kbits$$

# Example: Larger Block Size

- 64 blocks, 16 bytes/block (i.e., 4 words/block)
  - To what block number does address 1200 map?
- Block address = $\lfloor 1200/16 \rfloor$ = 75
- Block number = 75 modulo 64 = 11

| | tag | index | offset |
|---|---|---|---|

| | 31 | 10 9 | 4 3 | 0 |
|---|---|---|---|---|

1200

| 0..................001 | 001011 | 0000 |
|---|---|---|

22 bits      6 bits      4 bits

$\lfloor 1200/16 \rfloor$

| | 31 | 10 9 | 4 |
|---|---|---|---|

75

| 0.................001 | 001011 |
|---|---|

22 bits      6 bits

75 mod 64

11

| 001011 |
|---|

6 bits

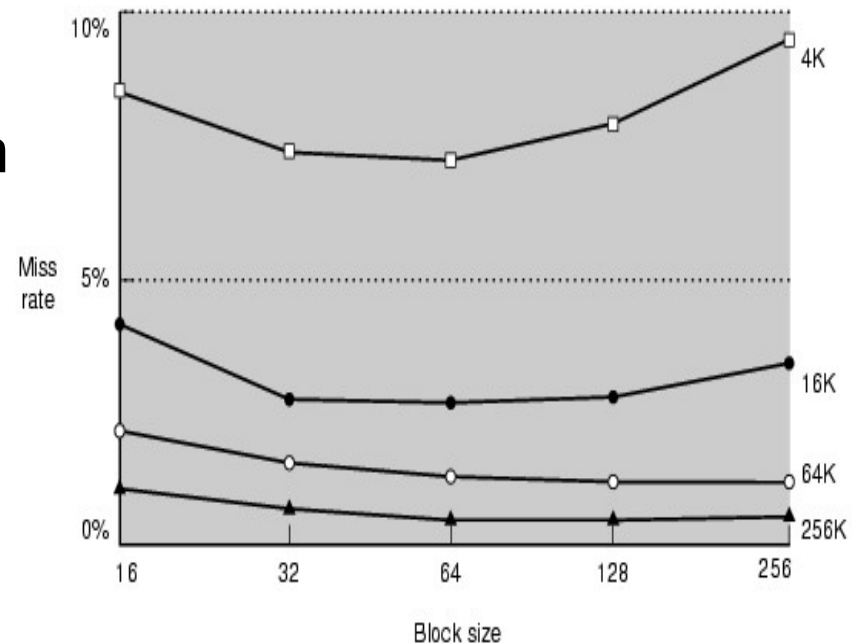| Byte no | Block no |
|---|---|
| 0 | |
| 1 | 0 |
| ... | |
| 15 | |
| 16 | |
| 17 | 1 |
| ... | |
| 31 | |
| 32 | |
| 33 | 2 |
| ... | |
| 47 | |

# 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 ☐
  - Can override benefit of reduced miss rate

# 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

# 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 + 0.1×100 = 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
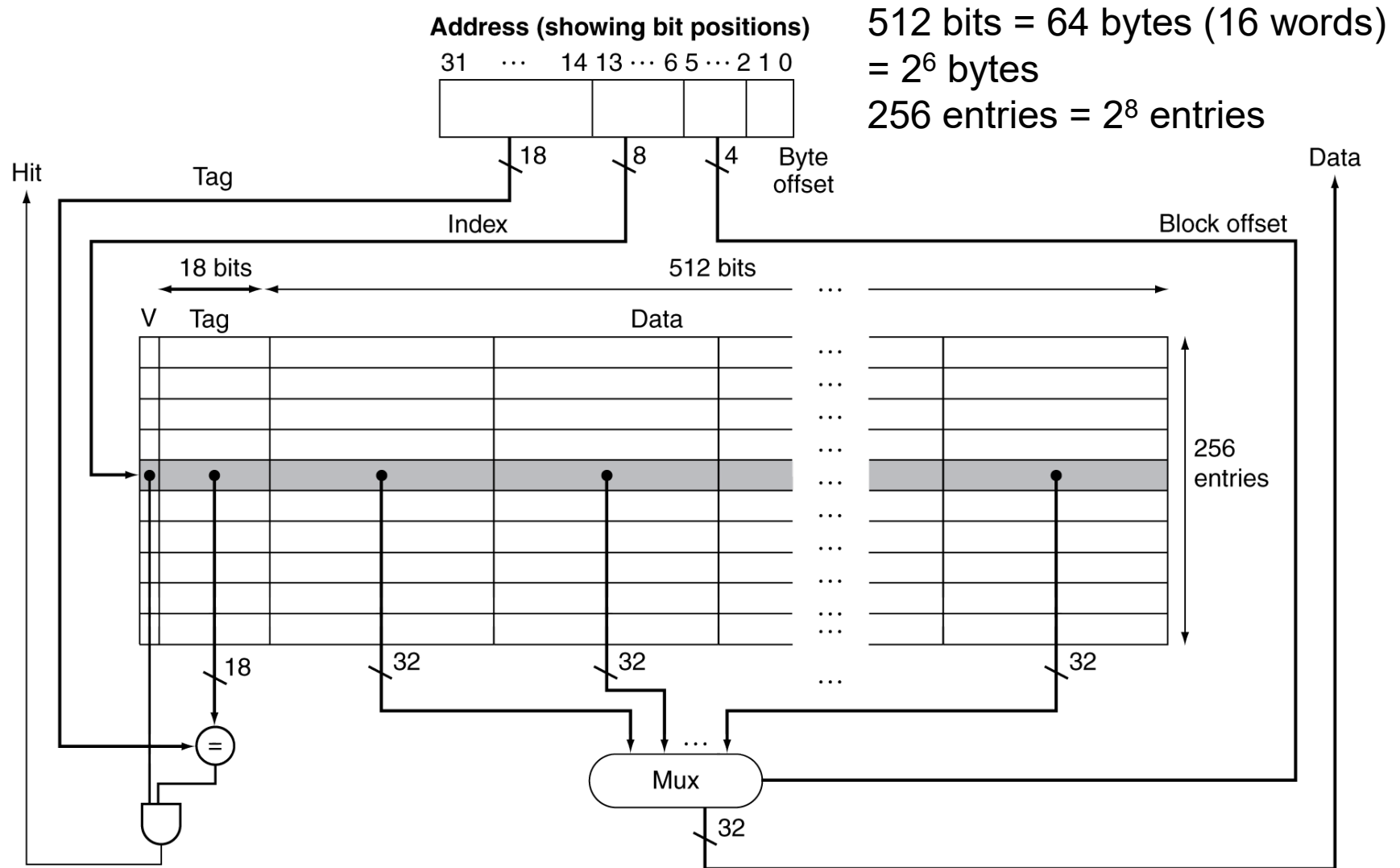    - Write it back to memory

# Write Allocation

- What should happen on a write miss?

- For write-through

  - Write around: don't fetch the block and write directly into memory (*write no-allocate*)

  - Allocate on miss: fetch the block into cache (*write allocate*)

# 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: 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: Intrinsity FastMATH

**Address (showing bit positions)**

31 ⋯ 14 13 ⋯ 6 5 ⋯ 2 1 0

512 bits = 64 bytes (16 words)
= $2^6$ bytes
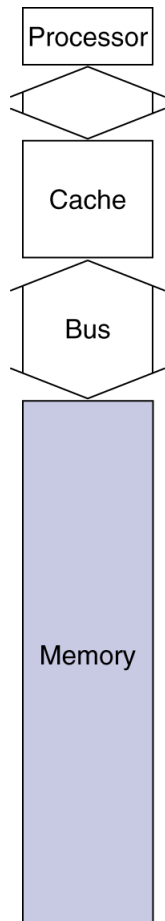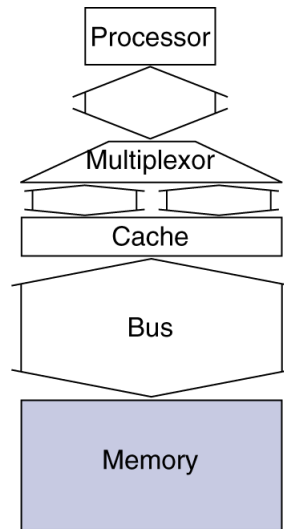256 entries = $2^8$ entries

# 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 cache block, 1-word-wide DRAM
  - Miss penalty = 1 + 4×15 + 4×1 = 65 bus cycles
  - Bandwidth = 16 bytes / 65 cycles = 0.25 B/cycle

# Increasing Memory Bandwidth



Processor — Cache — Bus — Memory

a. One-word-wide memory organization

Processor — Multiplexor — Cache — Bus — Memory

b. Wider memory organization

Processor — Cache — Bus — Memory bank 0 — Memory bank 1 — Memory bank 2 — Memory bank 3

c. Interleaved memory organization

- **4-word wide memory (DRAM)**
  - Miss penalty =
  - Bandwidth =
- **4-bank interleaved memory (DRAM)**
  - Miss penalty =
  - Bandwidth =

# Measuring Cache Performance

- Components of CPU time
  - Program execution cycles
    - Includes cache hit time
  - Memory stall cycles
    - Mainly from cache misses
- With simplifying assumptions
  - In terms of cycle count

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 <u>per instruction</u>**
  - I-cache:
  - D-cache:
- **Actual CPI =**
  - Ideal CPU is _____ times faster

# Average Access Time

- Hit time is also important for performance
- Average memory access time (AMAT)
  - AMAT = Hit time + Miss rate × Miss penalty
    - In terms of time
  - Hit time: cache/hardware design, etc.
  - Miss rate: block size, etc.
  - Miss penalty: memory organization, etc.
- Example
  - CPU with 1ns clock, hit time = 1 cycle, miss penalty = 20 cycles, cache miss rate = 5%
  - AMAT = 1 + 0.05 × 20 = 2 *ns*
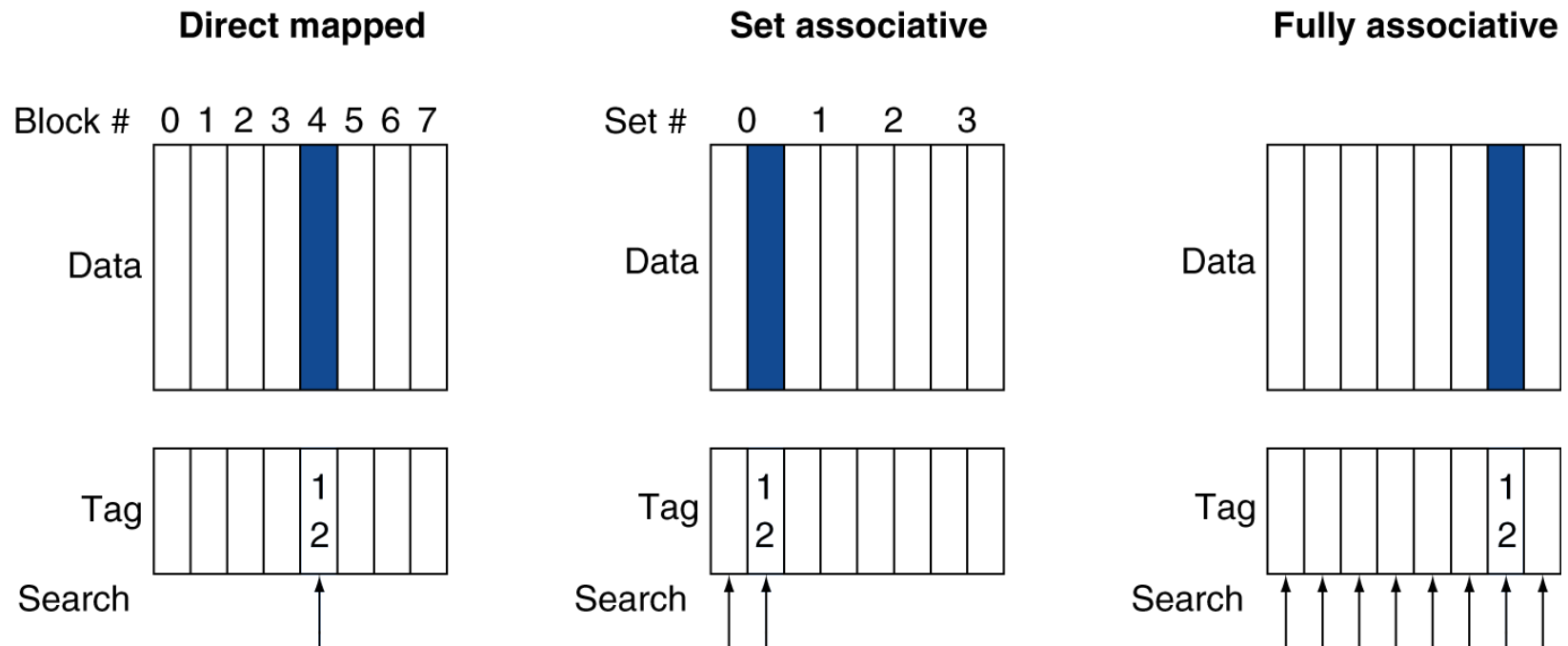    - AMAT per instruction = 2 cycles

# Performance Summary

- When CPU performance increased
  - Miss penalty becomes more significant
  - Performance is a function of **CPI** and **clock rate**

- 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

# 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

# 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
  - Increasing associativity shrinks index, expands tag
  - 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 (0000) | 0 (00) | miss | Mem[0] | | | |
| 8 (1000) | 0 (00) | miss | Mem[8] | | | |
| 0 (0000) | 0 (00) | miss | Mem[0] | | | |
| 6 (0110) | 2 (10) | miss | Mem[0] | | Mem[6] | |
| 8 (1000) | 0 (00) | 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 (0000) | 0 (0) | miss | Mem[0] | | | |
| 8 (1000) | 0 (0) | miss | Mem[0] | Mem[8] | | |
| 0 (0000) | 0 (0) | hit | Mem[0] | Mem[8] | | |
| 6 (0110) | 0 (0) | miss | Mem[0] | Mem[6] | | |
| 8 (1000) | 0 (0) | miss | Mem[8] | Mem[6] | | |

- ## Fully associative

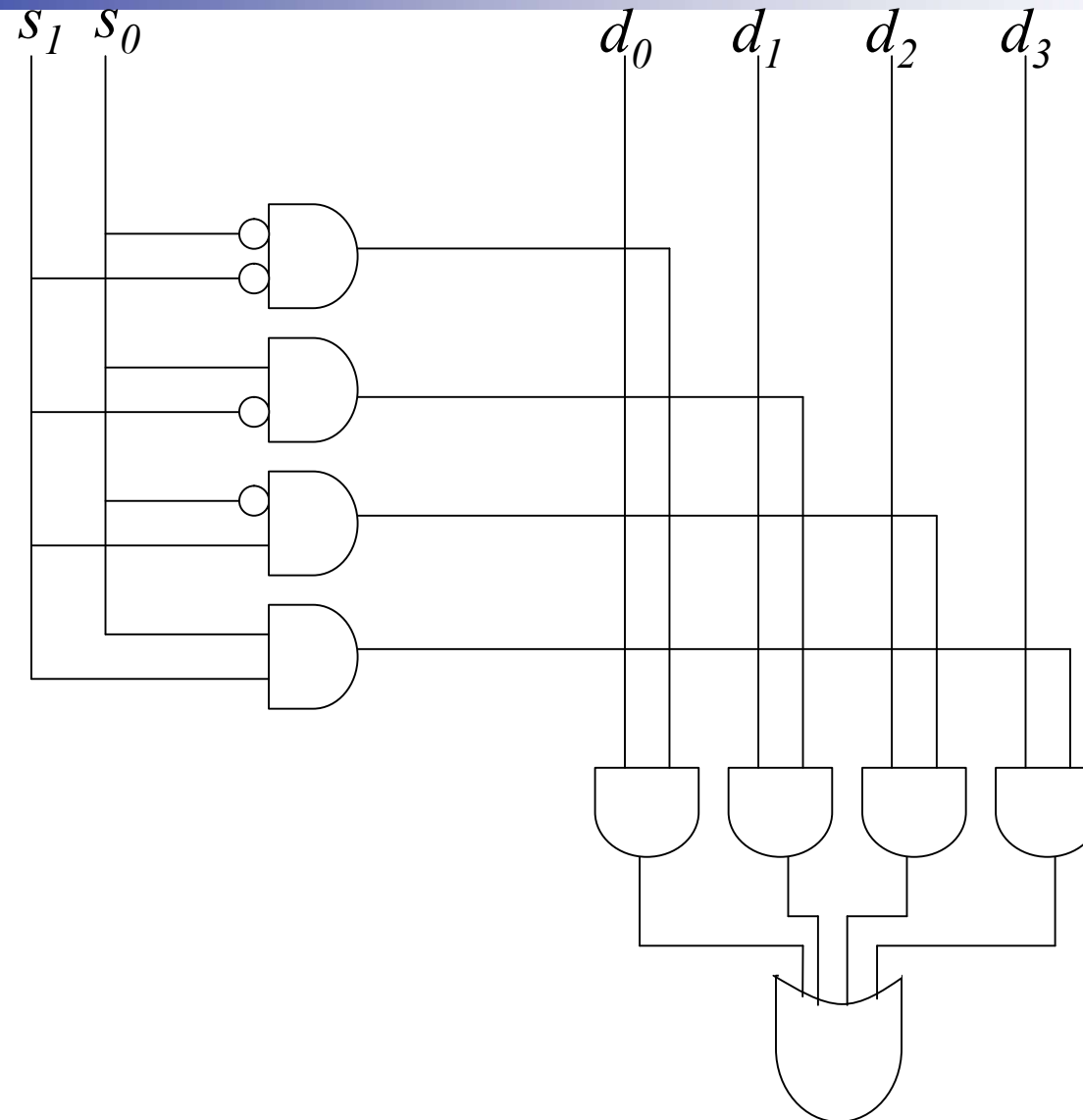| Block address | | Hit/miss | Cache content after access | | | |
|---|---|---|---|---|---|---|
| 0 (0000) | | miss | Mem[0] | | | |
| 8 (1000) | | miss | Mem[0] | Mem[8] | | |
| 0 (0000) | | hit | Mem[0] | Mem[8] | | |
| 6 (0110) | | miss | Mem[0] | Mem[8] | Mem[6] | |
| 8 (1000) | | hit | Mem[0] | Mem[8] | Mem[6] | |

# How Much Associativity

- Increased associativity decreases miss rate
  - But with diminishing returns
  - Extra search needed
- 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-to-1 Multiplexer

$s_1$  $s_0$                         $d_0$    $d_1$    $d_2$    $d_3$

# Replacement Policy

- Direct mapped: no choice
- Set associative
  - Prefer non-valid entry, if there is one
  - Otherwise, choose among entries in the set

  - [ ]

    - Choose the one unused for the longest time
      - Simple for 2-way, manageable for 4-way, too hard beyond that

  - [ ]

    - 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

# Multilevel Cache Example

- Given
  - CPU base CPI = 1, clock rate = 4GHz
  - Miss rate/instruction = 2%
  - Main memory access time = 100ns
- With just primary cache
  - Clock cycle $\qquad$ = 1/4GHz = 0.25ns
  - Miss penalty $\qquad$ = 100ns/0.25ns = 400 cycles
  - Effective CPI $\qquad$ = 1 + 0.02 × 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 = 5ns/0.25ns = 20 cycles
- Primary miss with L-2 miss
  - Extra penalty = 400 cycles
- CPI = 1 + 0.02 × 20 + 0.005 × 400 = 3.4
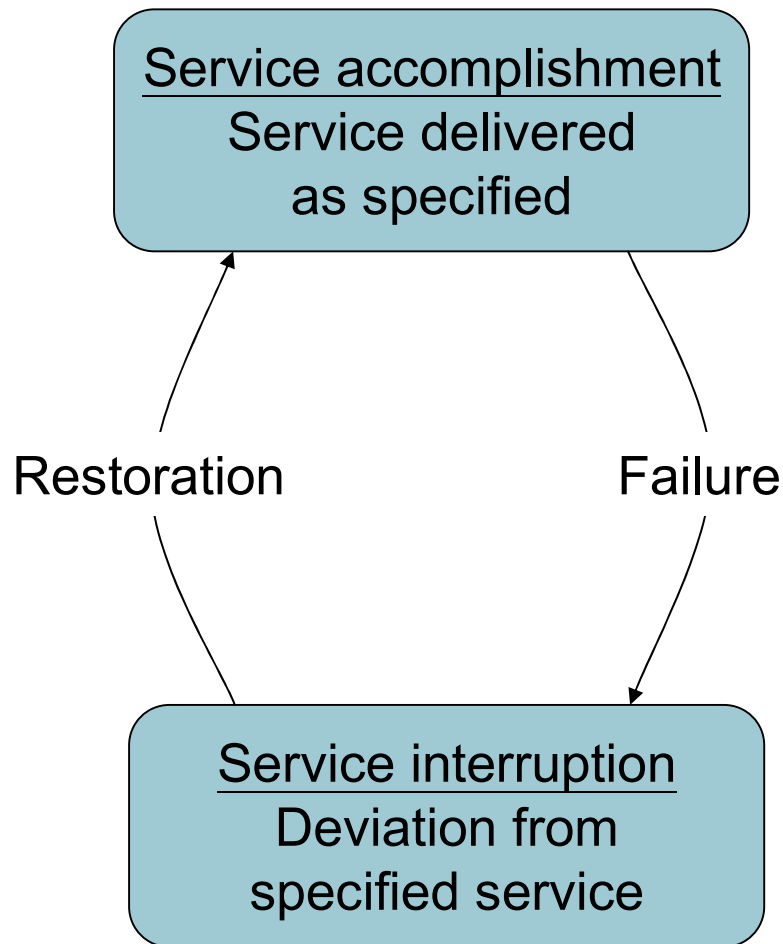- Performance ratio = 9/3.4 = 2.6

# 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

- Empirical results
  - L-1 cache usually smaller than a single cache
  - L-1 block size smaller than L-2 block size

# Interactions with Advanced CPUs

- Out-of-order CPUs can execute instructions during cache miss
  - Pending store stays in load/store unit
  - Dependent instructions wait in reservation stations
    - Independent instructions continue
- Effect of miss depends on program data flow
  - Much harder to analyze
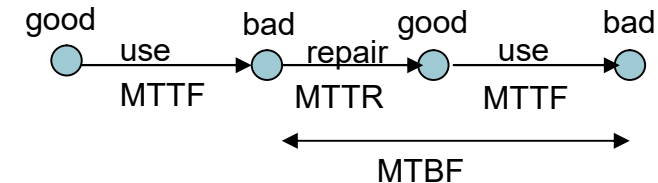  - Use system simulation

# Dependability



**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 prediction and correction
  - Reduce MTTR: improved tools and processes for (self-)diagnosis and (self-)repair

# Virtual Memory

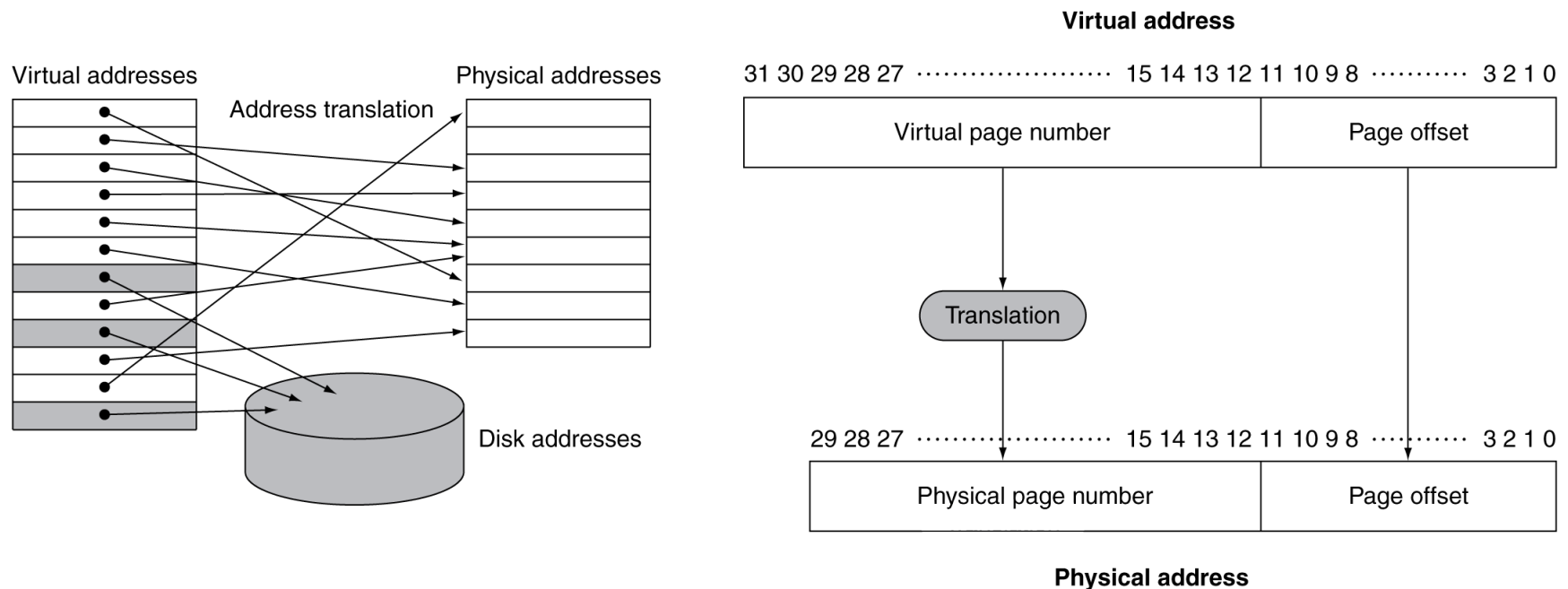- Use main memory as a "cache" for secondary (disk) storage
  - Managed jointly by CPU hardware and the operating system (OS)

- 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
  - VM translation "miss" is called a

# Address Translation

- Fixed-size pages (e.g., 4KB/page)

Virtual addresses                    Physical addresses

Address translation

Disk addresses

**Virtual address**

| 31 30 29 28 27 ·················· 15 14 13 12 11 10 9 8 ·········· 3 2 1 0 |
|---|---|
| Virtual page number | Page offset |

Translation

| 29 28 27 ·················· 15 14 13 12 11 10 9 8 ·········· 3 2 1 0 |
|---|---|
| Physical page number | Page offset |

**Physical address**

Since only a portion of pages exist in main memory,
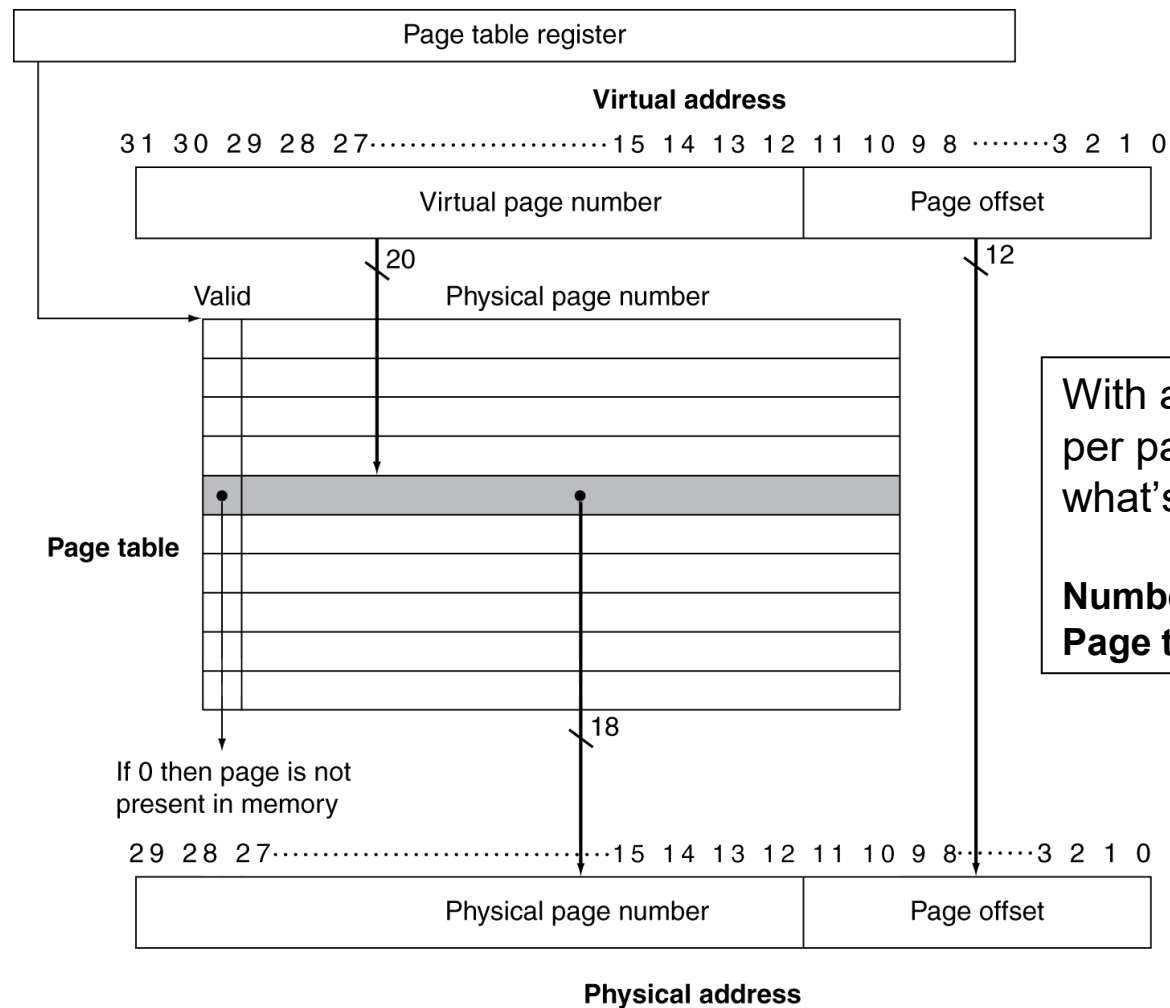| physical page number | < | virtual page number |

# 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

- **Stores placement information**
  - Array of page table entries (PTEs), indexed by virtual page number
  - <u>Page table register in CPU</u> points to <u>page table in physical memory</u>

- **If page is present in memory**
  - PTE stores the physical page number
  - Plus other status bits (referenced, dirty, …)

- **If page is not present**
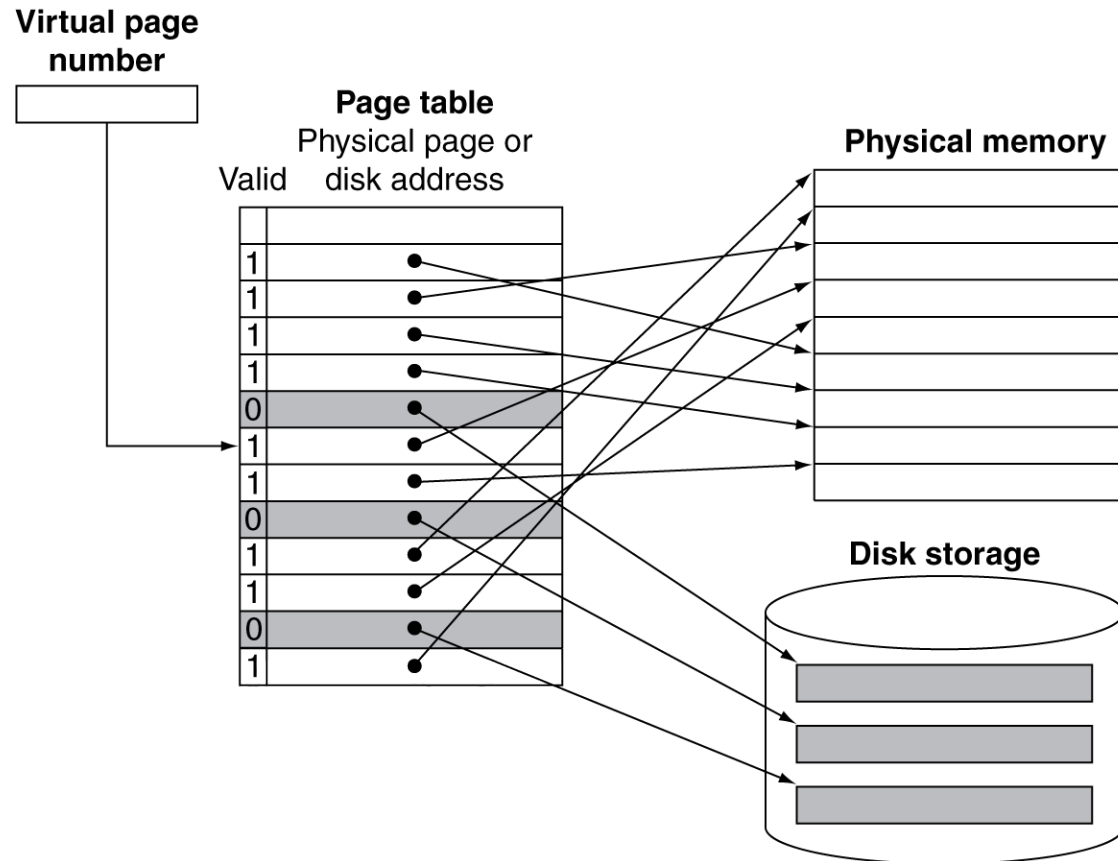  - PTE can refer to location in swap space on disk

# Translation Using a Page Table



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

Number of PTEs = $2^{32} / 2^{12} = 2^{20}$
Page table size = $2^{20} \times 4 = 4MB$

# 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/page at once, not individual locations
    - Write-through is impractical
    - Use write-back
    - Dirty bit in PTE set when page is written

# Page Table Problems

- ## Page table is too big
  - 4MB for 4KB per page and 4 bytes per PTE with a 32-bit virtual address

- ## Access to page table is too slow
  - One extra memory read

# Fast Translation Using a TLB

- Address translation would appear to require extra memory references
    - One to access the PTE
    - Then the actual memory access
- But access to page tables has good locality
    - So use a fast cache of PTEs within the CPU
    - Called a
    - Less entries than a real page table (in main memory)
    - 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
  - Then restart the faulting instruction

# TLB Miss Handler

- TLB miss indicates
  - Page present, but PTE not in TLB
  - Page not present
- Must recognize TLB miss before destination register overwritten (for "load")
  - Raise exception
- Handler copies PTE from memory to TLB
  - Then restarts instruction
  - If page not present, page fault will occur

# Page Fault Handler

- Use faulting virtual address to find PTE

- Locate page on disk

- Choose page to replace if there is no free memory page
  - If dirty, write to disk first

- Read page into memory and update page table

- Make process runnable again
  - Restart from faulting instruction

# TLB and Cache – Interaction

**Physically/virtually addressed cache**



- If cache tag uses physical address
  - Need to translate before cache lookup
  - Cache hit requires a TLB access and a cache access.
- Virtually addressed cache (rarely used)
  - remove TLB access
  - Complications due to aliasing
    - Different virtual addresses for shared physical address

# TLB and Cache – Flow

Virtual address

TLB access

TLB hit?
- No → TLB miss exception
- Yes → Physical address

Write?
- No → Try to read data from cache
- Yes → Write access bit on?

Write access bit on?
- No → Write protection exception
- Yes → Try to write data to cache

Try to read data from cache → Cache hit?
- No → Cache miss stall while read block
- Yes → Deliver data to the CPU

Try to write data to cache → Cache hit?
- No → Cache miss stall while read block
- Yes → Write data into cache, update the dirty bit, and put the data and the address into the write buffer

# 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

# 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 |
| n-way set associative | Set index, then search entries within the set | n |
| Fully associative | Search all entries | #entries |
| | Table lookup, e.g., page table | 0 |

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

# 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

# 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

# 3-Level Cache Organization

| 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 | 1 TLB for instructions and 1 TLB for data<br><br>Both TLBs are fully associative, with 32 entries, round robin replacement<br><br>TLB misses handled in hardware | 1 TLB for instructions and 1 TLB for data per core<br><br>Both L1 TLBs are four-way set associative, LRU replacement<br><br>L1 I-TLB has 128 entries for small pages, 7 per thread for large pages<br><br>L1 D-TLB has 64 entries for small pages, 32 for large pages<br><br>The L2 TLB is four-way set associative, LRU replacement<br><br>The L2 TLB has 512 entries<br><br>TLB misses handled in hardware |

# 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 | Decrease capacity misses | May increase access time |
| Increase associativity | Decrease conflict misses | May increase access time |
| Increase block size | Decrease compulsory misses | Increases miss penalty. For very large block size, may increase miss rate due to pollution. |

# TLB, Page Table and Cache

- The possible combinations of events in the TLB, virtual memory system, and physically indexed (tagged) cache.

| TLB | Page table | Cache | Possible ? If so, under what circumstance ? |
|-----|-----------|-------|---------------------------------------------|
| Hit | Hit | Miss | **Possible** |
| Miss | Hit | Hit | **Possible** |
| Miss | Hit | Miss | **Possible** |
| Miss | Miss | Miss | **Possible** |
| Hit | Miss | Miss | Impossible |
| Hit | Miss | Hit | Impossible |
| Miss | Miss | Hit | Impossible |

# Cache Coherence Problem

- Suppose two CPU cores share a physical address space
  - Write-through caches

| Time step | Event | CPU A's cache | CPU B's cache | Memory |
|---|---|---|---|---|
| 0 | | | | 0 |
| 1 | CPU A reads X | 0 | | 0 |
| 2 | CPU B reads X | 0 | 0 | 0 |
| 3 | CPU A writes 1 to X | 1 | 0 | 1 |

# Coherence Defined

- Informally: Reads return most recently written value

- Formally:

    - P writes X; P reads X (no intervening writes) $\Rightarrow$ read returns written value (program order)

    - $P_1$ writes X; $P_2$ reads X (sufficiently later) $\Rightarrow$ read returns written value (coherence)

        - c.f. CPU B reading X after step 3 in example

    - $P_1$ writes X, $P_2$ writes X $\Rightarrow$ all processors see writes in the same order

        - End up with the same final value for X

# Cache Coherence Protocols

- Operations performed by caches in multiprocessors to ensure coherence
    - Migration of data to local caches
        - Reduces bandwidth for shared memory
    - Replication of read-shared data
        - Reduces contention for access
- Snooping protocols
    - Each cache monitors bus reads/writes
- Directory-based protocols
    - Caches and memory record sharing status of blocks in a directory

# Invalidating Snooping Protocols

- Cache gets exclusive access to a block when it is to be written
  - Broadcasts an invalidate message on the bus
  - Subsequent read in another cache misses
    - Owning cache supplies updated value

| CPU activity | Bus activity | CPU A's cache | CPU B's cache | Memory |
|---|---|---|---|---|
| | | | | 0 |
| CPU A reads X | Cache miss for X | 0 | | 0 |
| CPU B reads X | Cache miss for X | 0 | 0 | 0 |
| CPU A writes 1 to X | Invalidate for X | 1 | | 0 |
| CPU B read X | Cache miss for X | 1 | 1 | 1 |

# Memory Consistency

- When are writes seen by other processors
  - "Seen" means a read returns the written value
  - Can't be instantaneously
- Assumptions
  - A write completes only when all processors have seen it
  - A processor does not reorder writes with other accesses
- Consequence
  - P writes location X then writes location Y
    $\Rightarrow$ all processors that see new Y also see new X
  - Processors can reorder reads, but not writes

# Pitfalls

- Byte (what we used) vs. word addressing
  - Example: 32-byte direct-mapped cache, 4-byte blocks
    - Byte 36 maps to block 1
    - Word 36 maps to block 4      0...0 0 0 1 0 0 1 0 1 1 0 0
  - Example: 256-byte cache and 32 bytes per block
    - Byte address 300 maps to block number 1

- Ignoring memory system effects when writing or generating code
  - Example: iterating over rows vs. columns of arrays
  - Large strides result in poor locality

```
for (i = 0; …)
   for (j = 0; …)
       x[i][j] = x[i][j]+y[i][j]
```

# Pitfalls

- In multiprocessor with shared L2 or L3 cache
  - Less associativity than cores results in conflict misses
  - More cores $\Rightarrow$ need to increase associativity
- Using AMAT to evaluate performance of out-of-order processors
  - Ignores effect of non-blocked accesses
  - Instead, evaluate performance by simulation

# 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**