

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

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 15

Taking Advantage of Locality

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

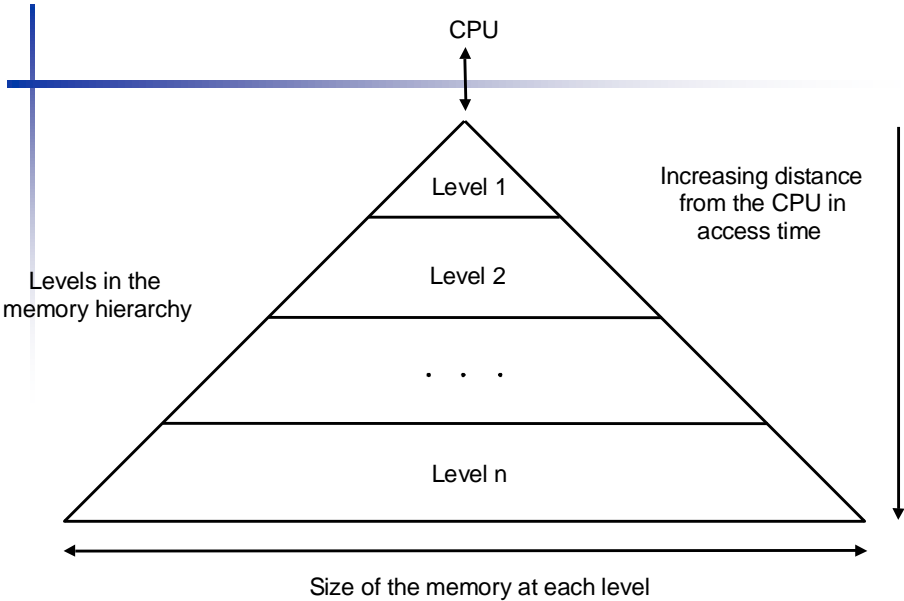
Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 16

Memory Hierarchy

The general idea is to build a hierarchy:

- at the top is a small, fast memory that is *close* to the processor.
- in the middle are larger, slower memories.
- At the bottom is massive memory with very slow access time.

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 17



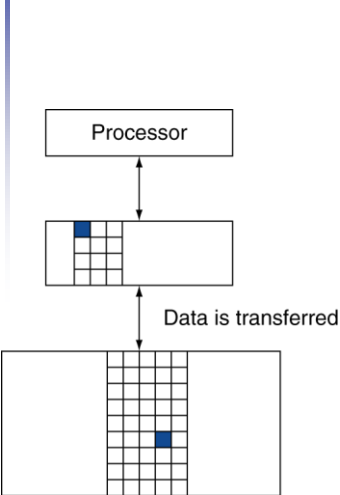
Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 18

Cache and Main Memory

- For now we will focus on a 2 level hierarchy:
 - cache (small, fast memory directly connected to the processor).
 - main memory (large, slow memory at level 2 in the hierarchy).

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 19

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 - hit ratio
 - Then accessed data supplied from upper level

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 20

Terminology

- **hit**: when the memory location accessed by the processor is in the cache (upper level).
- **miss**: when the memory location accessed by the process is not in the cache.
- **block**: the minimum unit of information transferred between the cache and the main memory. Typically measured in bytes or words.

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 21

Terminology (cont.)

- **hit rate**: the ratio of *hits* to total memory accesses.
- **miss rate**: $1 - \text{hit rate}$
- **hit time**: the time to access an element that is in the cache:
 - time to find out if it's in the cache.
 - time to transfer from cache to processor.

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 22

Terminology (cont.)

- *miss penalty*: the time to replace a block in the cache with a block from main memory and to deliver the element to the processor.
- *hit time* is small compared to *miss penalty* (otherwise we wouldn't bother with a memory hierarchy!)

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 23

Simple Cache Model

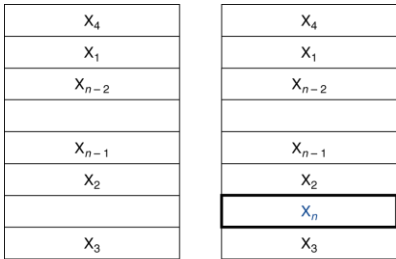
- Assume that the processor accesses memory one word at a time.
- A *block* consists of one word.
- When a word is referenced and is not in the cache, it is put in the cache (copied from main memory).

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 24

Cache Memory

§5.3 The Basics of Caches

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

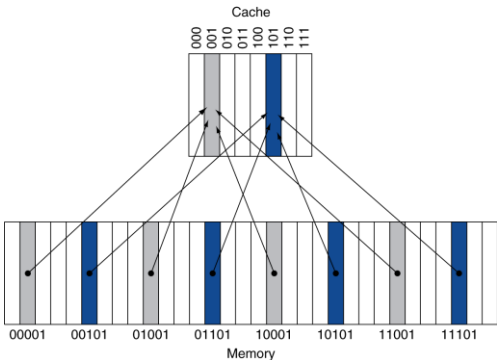


a. Before the reference to X_n b. After the reference to X_n

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

Direct Mapped Cache

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



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

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

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 27

Pop Quiz

- In order to determine which block is cached, we need to compare:
 - A: tag
 - B: valid bit
 - C: data
 - D: $A + B$
 - E: $B + C$

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 28

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		

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 29

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		

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 30

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		

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 31

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		

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 32

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		

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 33

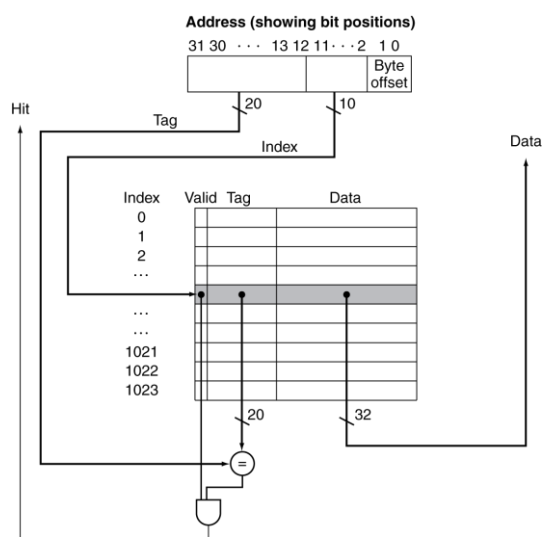
Cache Example

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		

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 34

Address Subdivision



Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 35

Pop Quiz

- Given the following:
 - 32 bit addresses (2^{32} byte memory, 2^{30} words)
 - 64 KB cache (16K words), each slot holds 1 word
 - Direct mapped cache
- How many bits are needed for each tag?
- A: 18 B: 16
- C: 14 D: 8

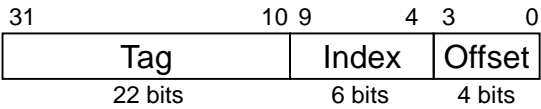
Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 36

Answer

- Memory has 2^{30} words
- Cache has $16K = 2^{14}$ slots (words)
- Each cache slot can hold any one of $2^{30} / 2^{14} = 2^{16}$ memory locations so the tag must be 16 bits

Example: Larger Block Size

- 64 blocks, 16 bytes/block
 - To what block number does address 1200 map?
- Block address = $\lfloor 1200/16 \rfloor = 75$
- Block number = $75 \bmod 64 = 11$



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

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 39

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

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 40

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

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 41

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

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 42

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

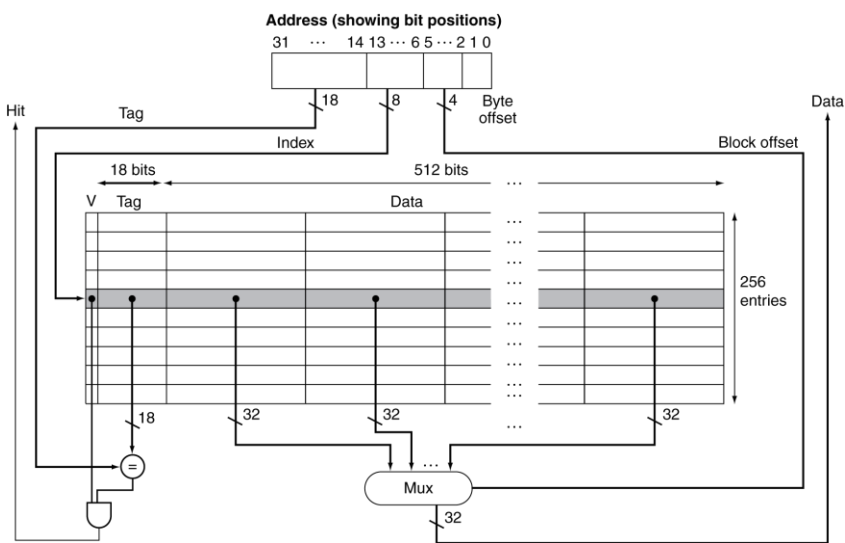
Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 43

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%

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 44

Example: Intrinsity FastMATH



Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 45

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 = $16 \text{ bytes} / 65 \text{ cycles} = 0.25 \text{ B/cycle}$

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 46