

# 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

§5.8 A Common Framework for Memory Hierarchies

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 99

# 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

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 100

## 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

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 101

## 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

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 102

## 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

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 103

## 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

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 104

## Pop Quiz

- In order to mitigate capacity misses in our cache, we can:      B
- A: Increase associativity
- B: Increase the size of our cache
- C: Increase the block size
- D: Increase the latency of our cache

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 105

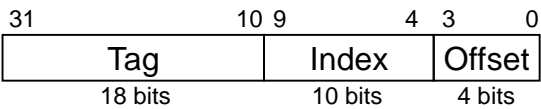
## Cache Design Trade-offs

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.

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 106

## Cache Control

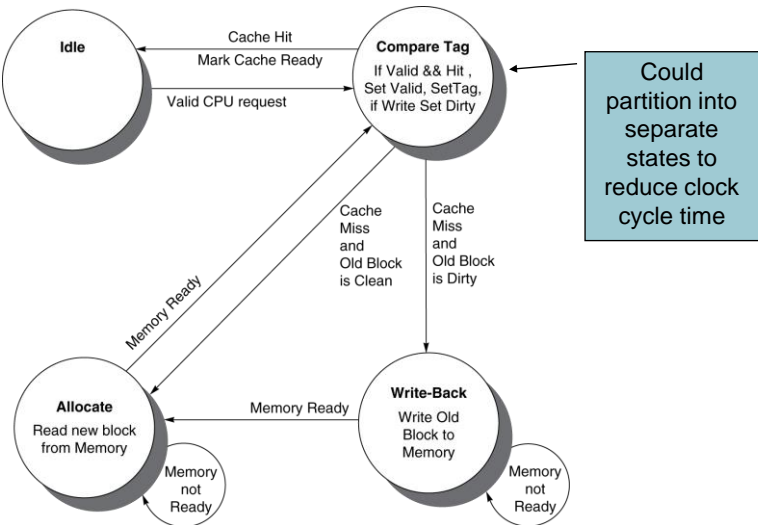
- Example cache characteristics
  - Direct-mapped, write-back, write allocate
  - Block size: 4 words (16 bytes)
  - Cache size: 16 KB (1024 blocks)
  - 32-bit byte addresses
  - Valid bit and dirty bit per block
  - Blocking cache
    - CPU waits until access is complete



Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 107

§5.9 Using a Finite State Machine to Control A Simple Cache

## Cache Controller FSM



Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 108

## 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

§5.10 Parallelism and Memory Hierarchies: Cache Coherence

## Coherence Defined

- Informally: Reads return most recently written value
- Formally:
  - P writes X; P reads X (no intervening writes)  
⇒ read returns written value
  - P<sub>1</sub> writes X; P<sub>2</sub> reads X (sufficiently later)  
⇒ read returns written value
    - c.f. CPU B reading X after step 3 in example
  - P<sub>1</sub> writes X, P<sub>2</sub> writes X  
⇒ 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

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 111

## 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

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 112

## 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 X then writes Y
    - ⇒ all processors that see new Y also see new X
  - Processors can reorder reads, but not writes

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 113

## Pitfalls

- Byte vs. word addressing
  - Example: 32-byte direct-mapped cache, 4-byte blocks
    - Byte 36 maps to block 1
    - Word 36 maps to block 4
- Ignoring memory system effects when writing or generating code
  - Example: iterating over rows vs. columns of arrays
  - Large strides result in poor locality

§5.15 Fallacies and Pitfalls

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 116



## 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

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 117

## 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  $\leftrightarrow$  ...  $\leftrightarrow$  DRAM memory  $\leftrightarrow$  disk
- Memory system design is critical for multiprocessors

§5.16 Concluding Remarks

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 119