+

William Stallings
Computer Organization
and Architecture
9th Edition
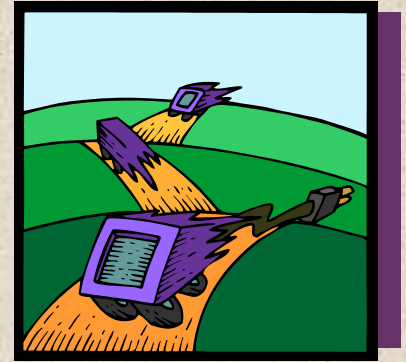
+ # Chapter 4

## Cache Memory

# Replacement Algorithms

- Once the cache has been filled, when a new block is brought into the cache, one of the existing blocks must be replaced

- For direct mapping there is only one possible line for any particular block and no choice is possible

- For the associative and set-associative techniques a replacement algorithm is needed

- To achieve high speed, an algorithm must be implemented in hardware

# The four most common replacement algorithms are:

- Least recently used (LRU)
  - Most effective
  - Replace that block in the set that has been in the cache longest with no reference to it
  - Because of its simplicity of implementation, LRU is the most popular replacement algorithm

- First-in-first-out (FIFO)
  - Replace that block in the set that has been in the cache longest
  - Easily implemented as a round-robin or circular buffer technique

- Least frequently used (LFU)
  - Replace that block in the set that has experienced the fewest references
  - Could be implemented by associating a counter with each line

# Write Policy

When a block that is resident in the cache is to be replaced there are two cases to consider:

There are two problems to contend with:

If the old block in the cache has not been altered then it may be overwritten with a new block without first writing out the old block

More than one device may have access to main memory

If at least one write operation has been performed on a word in that line of the cache then main memory must be updated by writing the line of cache out to the block of memory before bringing in the new block

A more complex problem occurs when multiple processors are attached to the same bus and each processor has its own local cache - if a word is altered in one cache it could conceivably invalidate a word in other caches

# Write Through and Write Back

- Write through
  - Simplest technique
  - All write operations are made to main memory as well as to the cache
  - The main disadvantage of this technique is that it generates substantial memory traffic and may create a bottleneck

- Write back
  - Minimizes memory writes
  - Updates are made only in the cache
  - Portions of main memory are invalid and hence accesses by I/O modules can be allowed only through the cache
  - This makes for complex circuitry and a potential bottleneck

# Line Size

When a block of data is retrieved and placed in the cache not only the desired word but also some number of adjacent words are retrieved
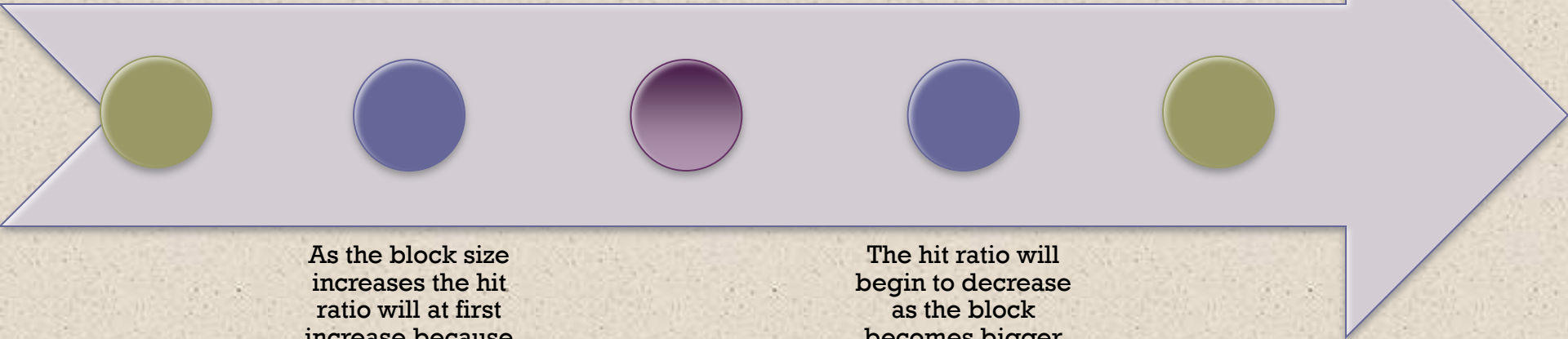
As the block size increases more useful data are brought into the cache

Two specific effects come into play:
- Larger blocks reduce the number of blocks that fit into a cache
- As a block becomes larger each additional word is farther from the requested word

As the block size increases the hit ratio will at first increase because of the principle of locality

The hit ratio will begin to decrease as the block becomes bigger and the probability of using the newly fetched information becomes less than the probability of reusing the information that has to be replaced

# Multilevel Caches

- As logic density has increased it has become possible to have a cache on the same chip as the processor

- The on-chip cache reduces the processor's external bus activity and speeds up execution time and increases overall system performance
  - When the requested instruction or data is found in the on-chip cache, the bus access is eliminated
  - On-chip cache accesses will complete appreciably faster than would even zero-wait state bus cycles
  - During this period the bus is free to support other transfers

- Two-level cache:
  - Internal cache designated as level 1 (L1)
  - External cache designated as level 2 (L2)

- Potential savings due to the use of an L2 cache depends on the hit rates in both the L1 and L2 caches

- The use of multilevel caches complicates all of the design issues related to caches, including size, replacement algorithm, and write policy

# Hit Ratio (L1 & L2)
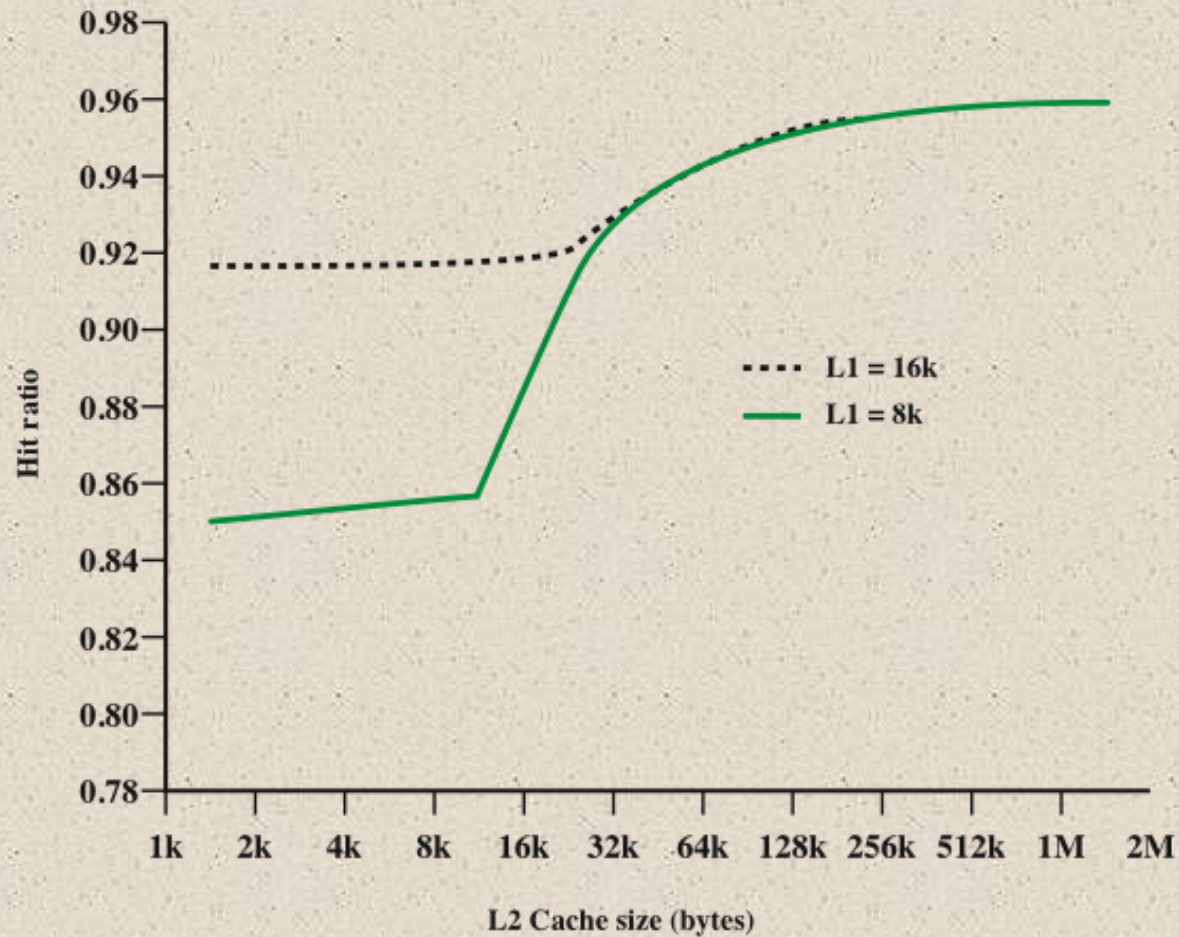# For 8 Kbyte and 16 Kbyte L1



Figure 4.17 Total Hit Ratio (L1 and L2) for 8 Kbyte and 16 Kbyte L1

# Unified Versus Split Caches

- Has become common to split cache:
    - One dedicated to instructions
    - One dedicated to data
    - Both exist at the same level, typically as two L1 caches

- Advantages of unified cache:
    - Higher hit rate
        - Balances load of instruction and data fetches automatically
        - Only one cache needs to be designed and implemented

- Trend is toward split caches at the L1 and unified caches for higher levels

- Advantages of split cache:
    - Eliminates cache contention between instruction fetch/decode unit and execution unit
        - Important in pipelining

# Summary

## Chapter 4

- Characteristics of Memory Systems
  - Location
  - Capacity
  - Unit of transfer
- Memory Hierarchy
  - How much?
  - How fast?
  - How expensive?
- Cache memory principles

## Cache Memory

- Elements of cache design
  - Cache addresses
  - Cache size
  - Mapping function
  - Replacement algorithms
  - Write policy
  - Line size
  - Number of caches
- Pentium 4 cache organization
- ARM cache organization