



Bahria University
Discovering Knowledge

Computer Architecture and Logic Design (CALD)

Lecture 04

Dr. Sorath Hansrajani

Assistant Professor

Department of Software Engineering

Bahria University Karachi Campus

Email: sorathhansrajani.bukc@bahria.edu.pk

Memory Systems



Characteristics of Memory Systems

Table 4.1 Key Characteristics of Computer Memory Systems

| | |
|--|---------------------------------|
| Location | Performance |
| Internal (e.g., processor registers, cache, main memory) | Access time |
| External (e.g., optical disks, magnetic disks, tapes) | Cycle time |
| Capacity | Transfer rate |
| Number of words | Physical Type |
| Number of bytes | Semiconductor |
| Unit of Transfer | Magnetic |
| Word | Optical |
| Block | Magneto-optical |
| Access Method | Physical Characteristics |
| Sequential | Volatile/nonvolatile |
| Direct | Erasable/nonerasable |
| Random | Organization |
| Associative | Memory modules |



Characteristics of Memory Systems

■ Location

- Refers to whether memory is internal or external to the computer
- Internal memory is often equated with main memory
- Processor requires its own local memory, in the form of registers
- Cache is another form of internal memory
- External memory consists of peripheral storage devices that are accessible to the processor via I/O controllers

■ Capacity

- Memory is typically expressed in terms of bytes

■ Unit of transfer

- For main memory, this is the number of bits read out of or written into memory at a time.
- For external memory, data are often transferred in much larger units than a word, and these are referred to as blocks



Method of Accessing Units of Data

■ Sequential Access

- The memory is accessed in a specific linear sequential manner, like accessing in a single Linked List.
- The access time depends on the location of the data.
- Examples: magnetic tapes, magnetic disk and optical memories.

■ Random Access

- Any location of the memory can be accessed randomly like accessing in Array.
- Physical locations are independent in this access method.
- Examples: RAM, and ROM



Method of Accessing Units of Data

■ Direct Access

- Individual blocks or records have a unique address based on physical location.
- Access is accomplished by direct access to reach a general vicinity plus sequential searching, counting or waiting to reach the final destination.
- Combination of Sequential and Random access methods.
- Example: Magnetic hard-disk

■ Associate Access

- A word is accessed rather than its address.
- special type of random access method.
- Example: Cache memory.

Capacity and Performance

The two most important characteristics of memory

Three performance parameters are used:

Access time (latency)

- Time between presenting the address and getting the valid data

Memory cycle time

- Time may be required for the memory to “recover” before next access.
- Cycle time = access time + recovery time

Transfer rate

- The rate at which data can be transferred into or out of a memory unit

+ Physical Types and Characteristics

■ The most common forms are:

- Semiconductor memory
- Magnetic surface memory
- Optical
- Magneto-optical

■ Several physical characteristics of data storage are important:

- Volatile memory
 - Information decays naturally or is lost when electrical power is switched off
- Nonvolatile memory
 - Once recorded, information remains without deterioration until deliberately changed
 - No electrical power is needed to retain information
- Magnetic-surface memories
 - Are nonvolatile
- Semiconductor memory
 - May be either volatile or nonvolatile
- Nonerasable memory
 - Cannot be altered, except by destroying the storage unit
 - Semiconductor memory of this type is known as read-only memory (ROM)

■ For random-access memory the organization is a key design issue

- Organization refers to the physical arrangement of bits to form words



+ Memory Hierarchy

- Design constraints on a computer's memory can be summed up by three questions:
 - How much, how fast, how expensive
- There is a trade-off among capacity, access time, and cost
 - Faster access time, greater cost per bit
 - Greater capacity, smaller cost per bit
 - Greater capacity, slower access time
- The way out of the memory dilemma is not to rely on a single memory component or technology, but to employ a memory hierarchy

+ Memory Hierarchy

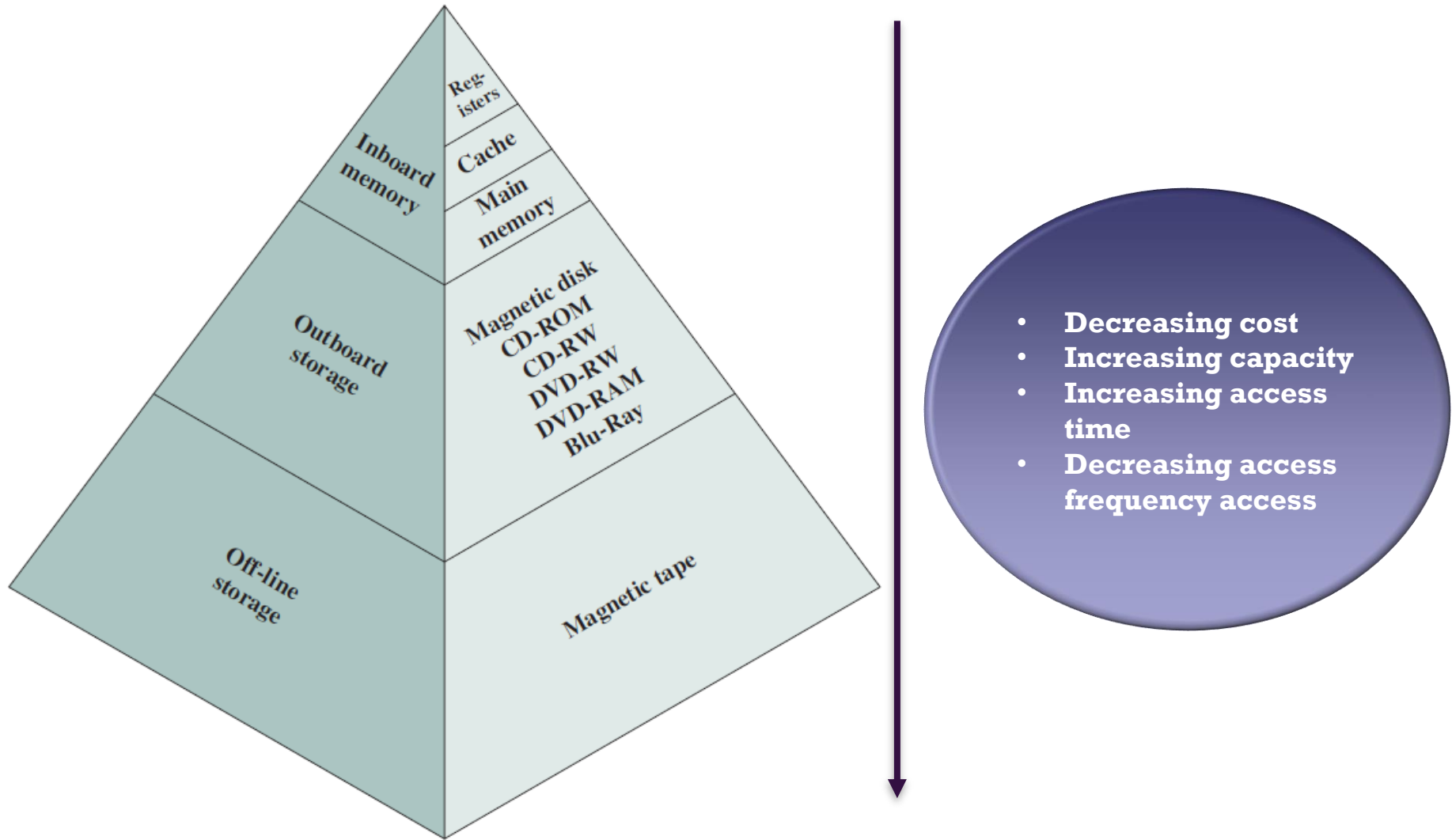
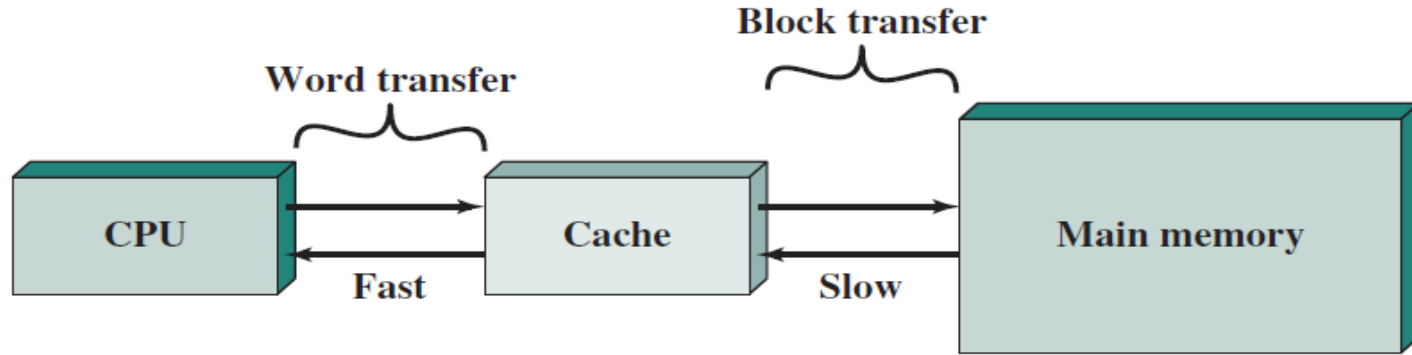


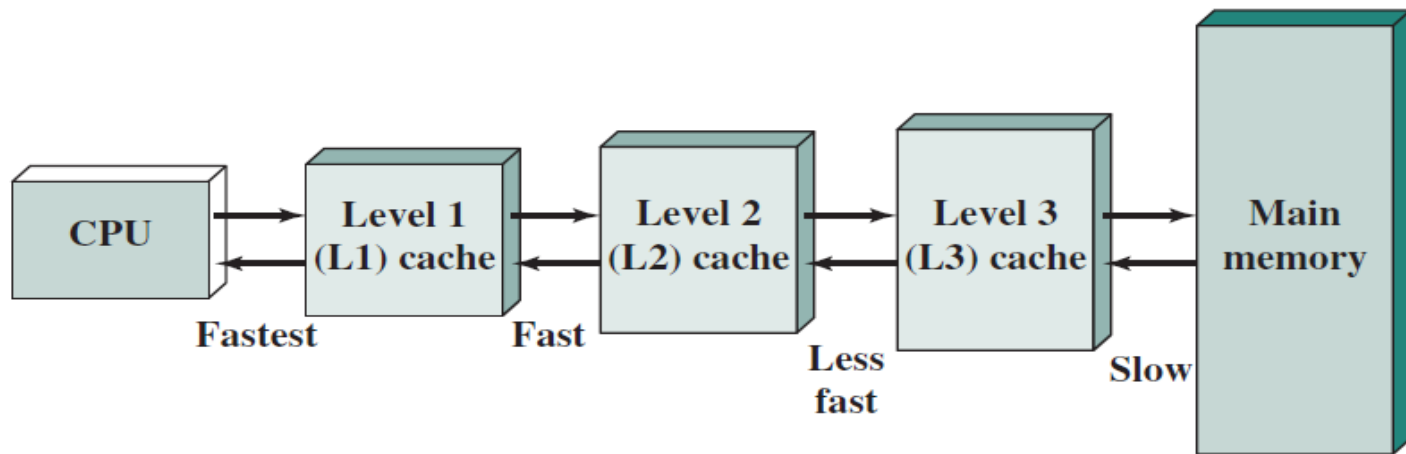
Figure 4.1 The Memory Hierarchy

+

Cache and Main Memory



(a) Single cache



(b) Three-level cache organization

Figure 4.3 Cache and Main Memory

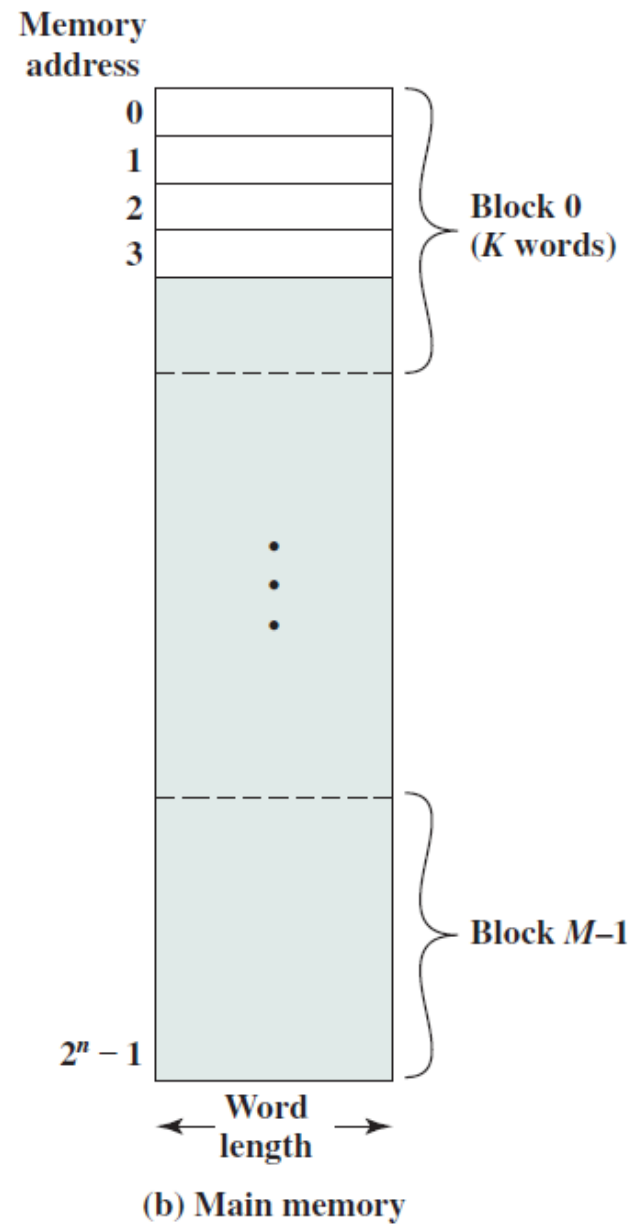
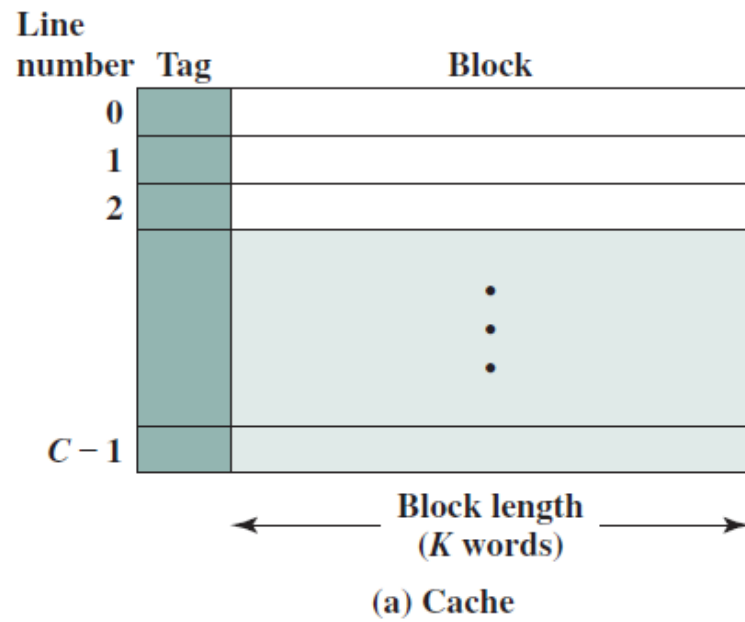


Figure 4.4 Cache/Main Memory Structure

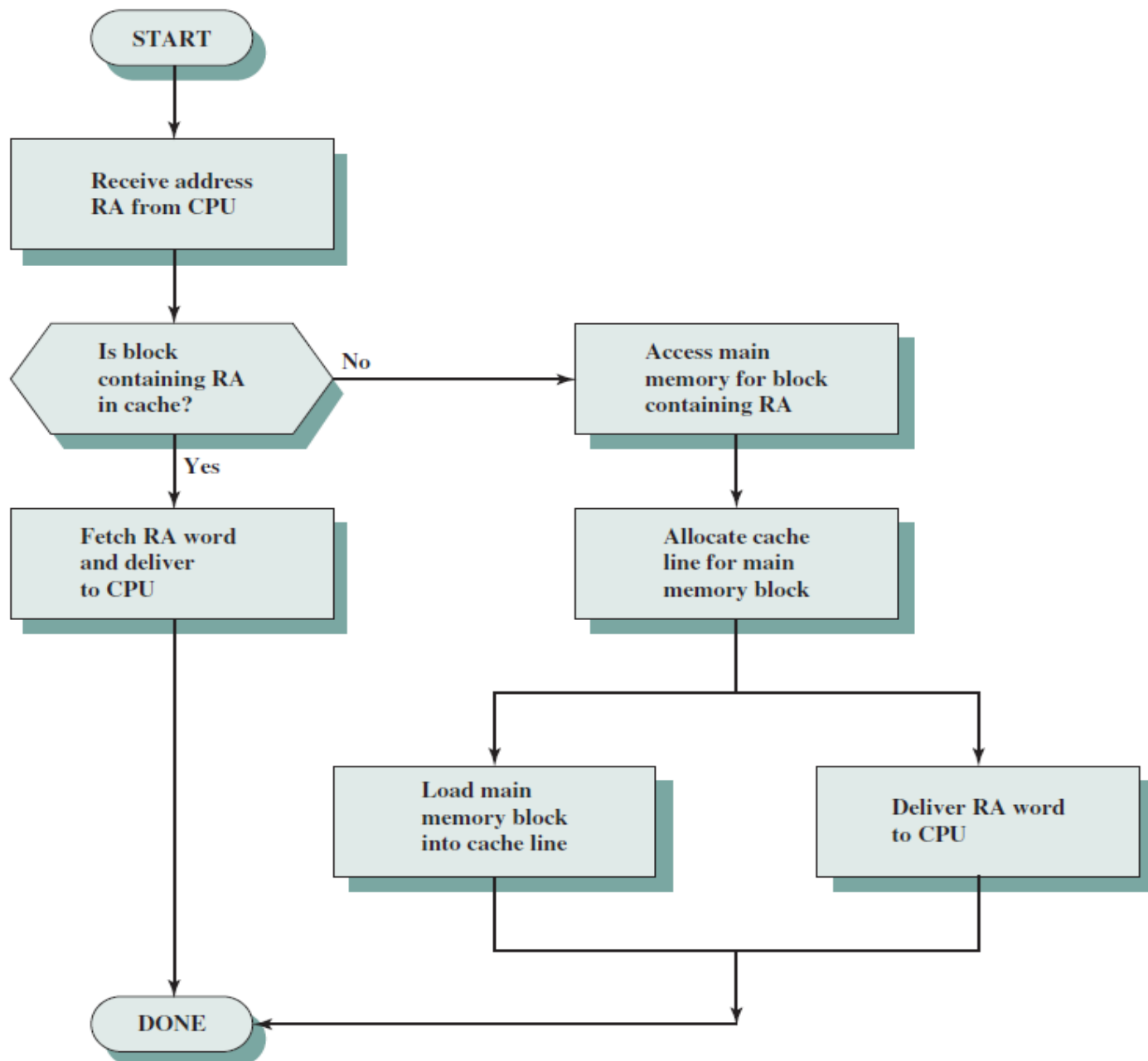


Figure 4.5 Cache Read Operation

| Processor | Type | Year of Introduction | L1 Cachea | L2 cache | L3 Cache |
|-----------------------|-----------------------------------|----------------------|-----------------------|----------------|--------------------------|
| IBM 360/85 | Mainframe | 1968 | 16 to 32 kB | — | — |
| PDP-11/70 | Minicomputer | 1975 | 1 kB | — | — |
| VAX 11/780 | Minicomputer | 1978 | 16 kB | — | — |
| IBM 3033 | Mainframe | 1978 | 64 kB | — | — |
| IBM 3090 | Mainframe | 1985 | 128 to 256 kB | — | — |
| Intel 80486 | PC | 1989 | 8 kB | — | — |
| Pentium | PC | 1993 | 8 kB/8 kB | 256 to 512 KB | — |
| PowerPC 601 | PC | 1993 | 32 kB | — | — |
| PowerPC 620 | PC | 1996 | 32 kB/32 kB | — | — |
| PowerPC G4 | PC/server | 1999 | 32 kB/32 kB | 256 KB to 1 MB | 2 MB |
| IBM S/390 G6 | Mainframe | 1999 | 256 kB | 8 MB | — |
| Pentium 4 | PC/server | 2000 | 8 kB/8 kB | 256 KB | — |
| IBM SP | High-end server/ supercomputer | 2000 | 64 kB/32 kB | 8 MB | — |
| CRAY MTA ^b | Supercomputer | 2000 | 8 kB | 2 MB | — |
| Itanium | PC/server | 2001 | 16 kB/16 kB | 96 KB | 4 MB |
| Itanium 2 | PC/server | 2002 | 32 kB | 256 KB | 6 MB |
| IBM POWER5 | High-end server | 2003 | 64 kB | 1.9 MB | 36 MB |
| CRAY XD-1 | Supercomputer | 2004 | 64 kB/64 kB | 1MB | — |
| IBM POWER6 | PC/server | 2007 | 64 kB/64 kB | 4 MB | 32 MB |
| IBM z10 | Mainframe | 2008 | 64 kB/128 kB | 3 MB | 24-48 MB |
| Intel Core i7 EE 990 | Workstaton/ server | 2011 | 6 ´ 32 kB/32 kB | 1.5 MB | 12 MB |
| IBM zEnterprise 196 | Mainframe/ Server | 2011 | 24 ´ 64 kB/ 128 kB | 24 ´ 1.5 MB | 24 MB L3 192 MB L4 |

Table 4.3

Cache Sizes of Some Processors

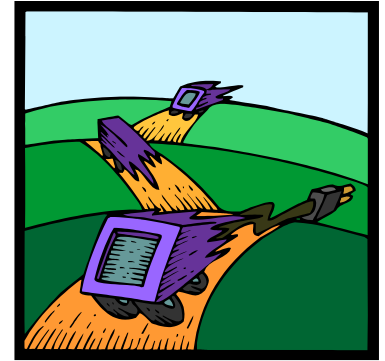
^a Two values separated by a slash refer to instruction and data caches.

^b Both caches are instruction only; no data caches.

(Table can be found on page 134 in the textbook.)



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

+ The 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:



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



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

There are two problems to contend with:



More than one device may have access to main memory



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



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

Table 4.4 Intel Cache Evolution

| Problem | Solution | Processor on Which Feature First Appears |
|---|--|--|
| External memory slower than the system bus. | Add external cache using faster memory technology. | 386 |
| Increased processor speed results in external bus becoming a bottleneck for cache access. | Move external cache on-chip, operating at the same speed as the processor. | 486 |
| Internal cache is rather small, due to limited space on chip. | Add external L2 cache using faster technology than main memory. | 486 |
| Contention occurs when both the Instruction Prefetcher and the Execution Unit simultaneously require access to the cache. In that case, the Prefetcher is stalled while the Execution Unit's data access takes place. | Create separate data and instruction caches. | Pentium |
| Increased processor speed results in external bus becoming a bottleneck for L2 cache access. | Create separate back-side bus that runs at higher speed than the main (front-side) external bus. The BSB is dedicated to the L2 cache. | Pentium Pro |
| | Move L2 cache on to the processor chip. | Pentium II |
| Some applications deal with massive databases and must have rapid access to large amounts of data. The on-chip caches are too small. | Add external L3 cache. | Pentium III |
| | Move L3 cache on-chip. | Pentium 4 |

+ Summary

Chapter 4

Cache Memory

- **Computer memory system overview**

- Characteristics of Memory Systems
- Memory Hierarchy

- **Cache memory principles**

- **Elements of cache design**

- Cache size
- Replacement algorithms
- Write policy
- Number of caches