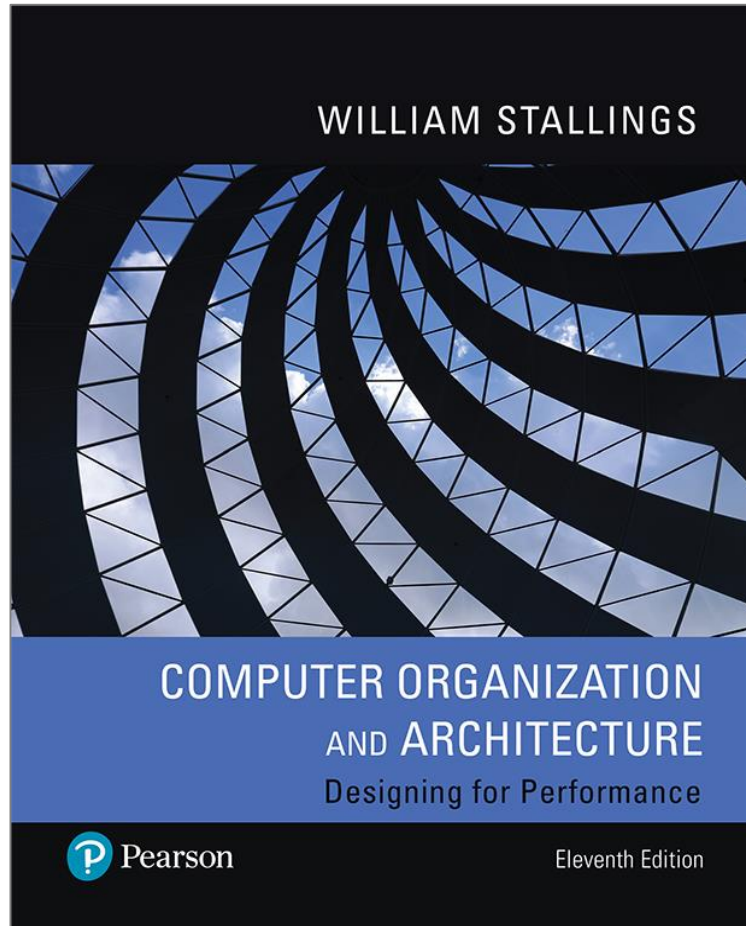


Computer Organization and Architecture

Designing for Performance

11th Edition



Chapter 4

The Memory Hierarchy:
Locality and Performance

Principle of Locality (1 of 2)

- Also referred to as the *locality of reference*
- Reflects the observation that during the course of execution of a program, memory references by the processor tend to **cluster**
- Think about memory access when:
 - Searching in an array – elements are sequentially stored in memory
 - Using loops – repeated small set of instructions
 - Executing a subroutine

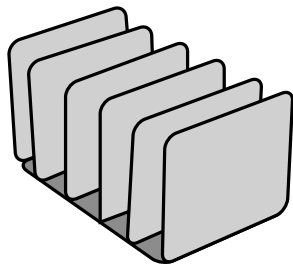
Principle of Locality (1 of 2)

- Locality is based on three assertions:
 - During any interval of time, a program references memory location **non-uniformly**, meaning some units of memory are **more likely to be accessed** than others.
 - As a function of time, the probability that a given unit of memory is referenced tends to change slowly. The **probability distribution** of memory references across the entire memory space tends to **change slowly over time**.
 - The **correlation** between **immediate past and immediate future memory reference patterns is high** and is weakened as the time interval increases

Principle of Locality (2 of 2)

- Two forms of locality
 - Temporal locality
 - Refers to the tendency of a program to reference in the near future those units of memory referenced in the recent past
 - Constants, temporary variables, and working stacks are also constructs that lead to this principle
 - Spatial locality
 - Spatial locality
 - Refers to the tendency of a program to reference units of memory whose addresses are near one another
 - Also reflects the tendency of a program to access data locations sequentially, such as when processing a table of data

Figure 4.1



**File organizer
on Bob's desk**



Filing cabinets in next room

Figure 4.1 Moving File Folders Between Smaller, Faster-Access Storage and Larger, Slower-Access Storage

Analogy

- Bob spends most of his time dealing with file folders.
- The folder are stored in the cabinets next room.
- For convenience, Bob has a file organizer on his desk which has slots that can hold few folders at one time.

Analogy – Bob Exploits Temporal Locality

- When Bob is working on a file and temporarily is finished, it may be likely that he will need to read or write one of the documents in that file in the near future, so he keeps it in his desk organizer.
- This is an example of temporal locality

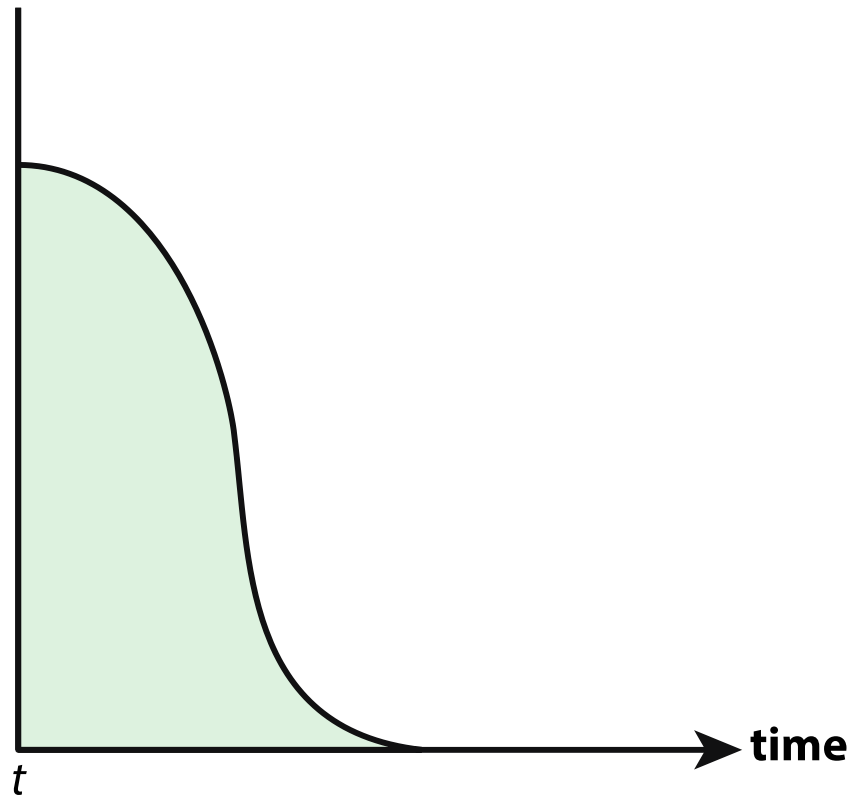
Analogy – Bob Exploits Spatial Locality

- Bob also observes that when he retrieves a folder from the filing cabinets, it is likely that in the near future he will need access to some of the nearby folders as well, so he retrieves the folder he needs plus a few folders on either side at the same time.
- This is an example of spatial locality.

Analogy – Returning Folders

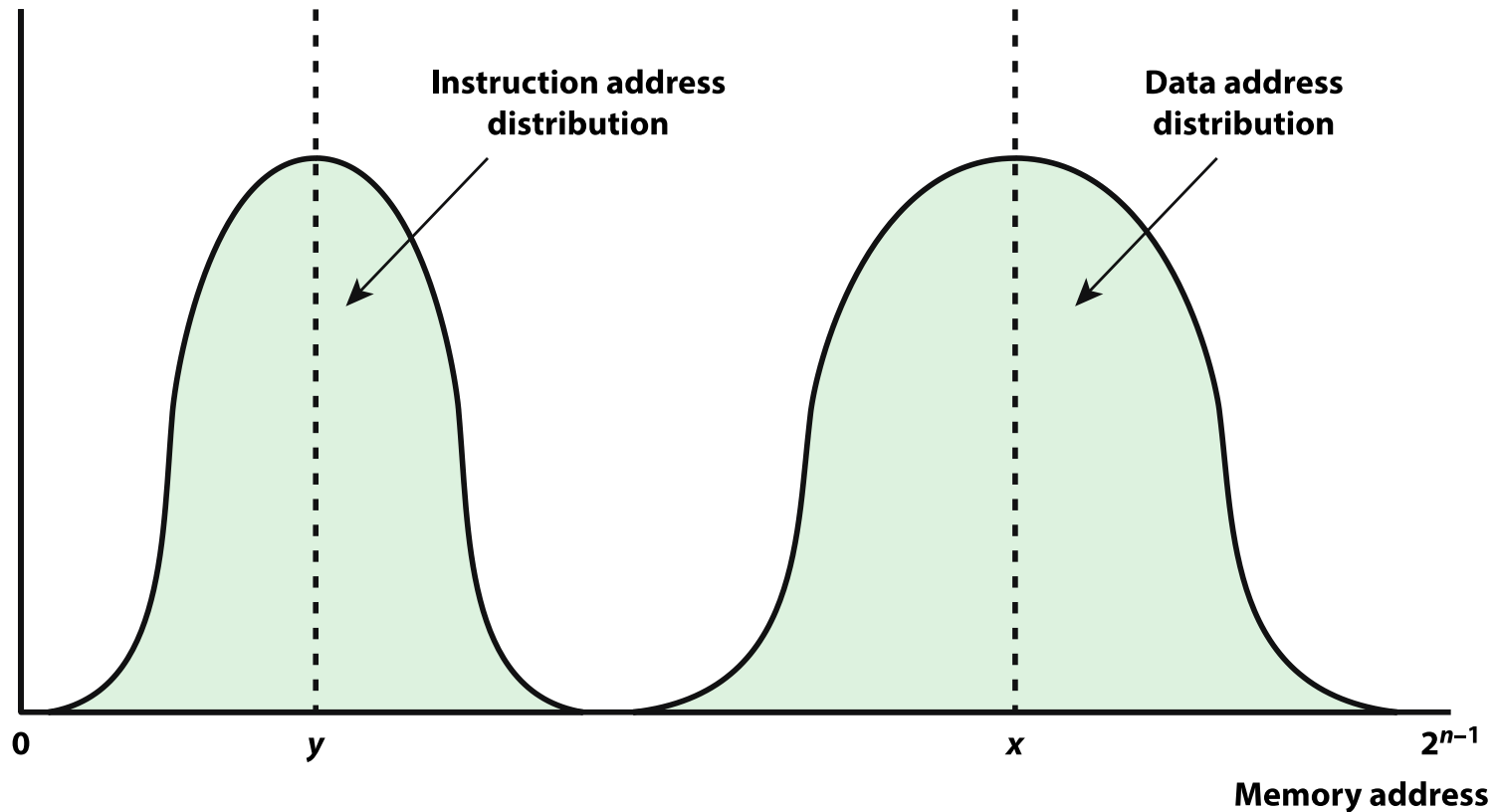
- Soon Bob's desk organizer is going to be filled up.
- Next time he goes to retrieve a folder from the cabinets, he must return some/all of the folders in his organizer.
- What policy should Bob follow in such case?
 - Replace the folder that has been in his desk organizer the longest or the least folder that has been used recently – Temporal Locality
 - Replace all the folders in his desk organizer with a new batch of contiguous folders including the one he needs plus nearby folders – Spatial Locality
- Neither policy is optimal. How about using both to a certain degree? E.g., batches of 10-20% of his desk organizer cap.

Figure 4.2



**Figure 4.2 Idealized Temporal Locality Behavior:
Probability Distribution for Time of Next Memory Access
to Memory Unit Accessed at Time t**

Figure 4.3



**Figure 4.3 Idealized Spatial Locality Behavior:
Probability Distribution for Next Memory Access
(most recent data memory access at location x ;
most recent instruction fetch at location y)**

Table 4.1

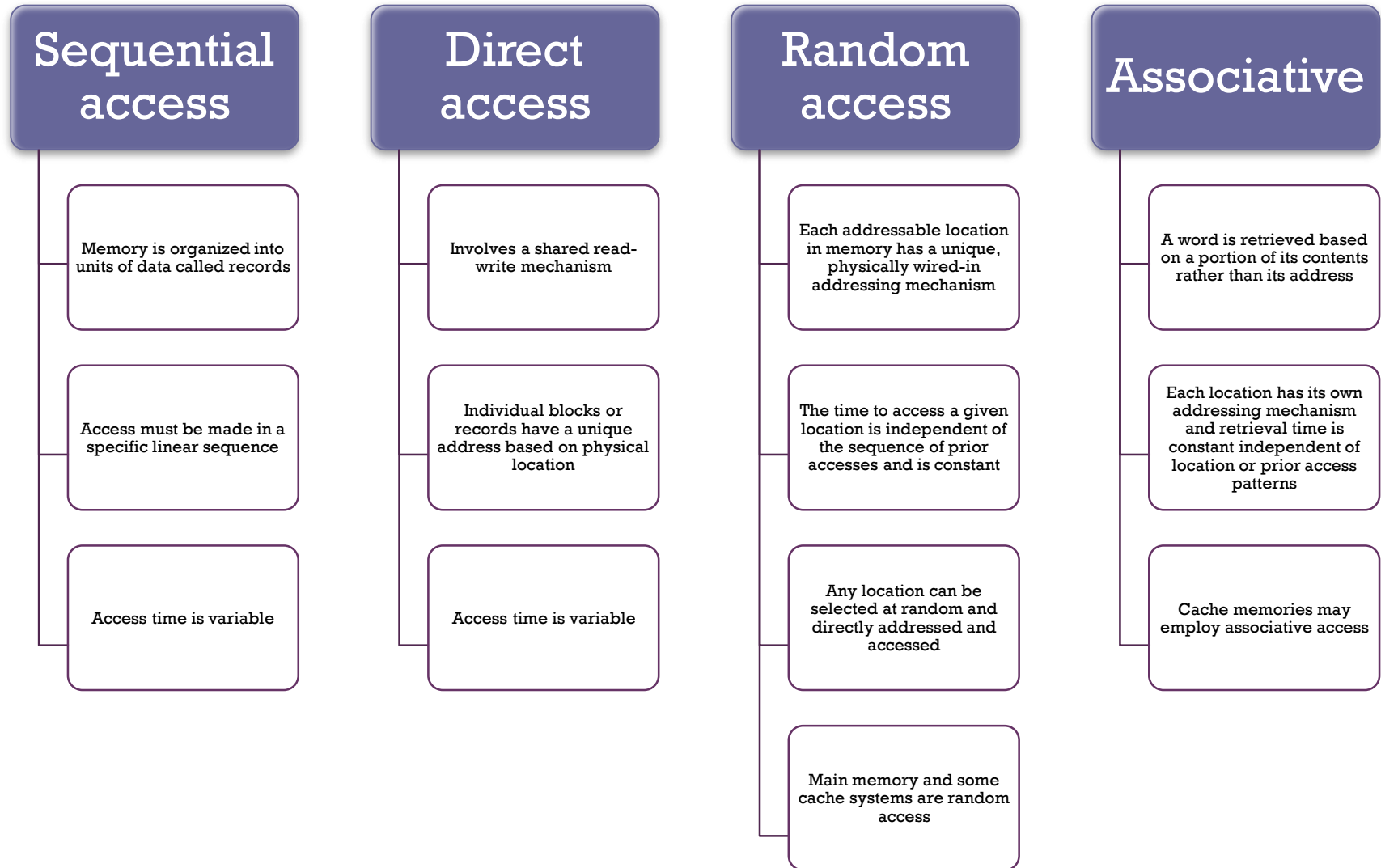
Key Characteristics of Computer Memory Systems

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

Characteristics of Memory Systems

- Location
 - Refers to whether memory is internal and 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 internal memory the unit of transfer is equal to the number of electrical lines into and out of the memory module

Method of Accessing Units of Data



Capacity and Performance:

The two most important characteristics of memory

Three performance parameters are used:

Access time (latency)

- For random-access memory it is the time it takes to perform a read or write operation
- For non-random-access memory it is the time it takes to position the read-write mechanism at the desired location

Memory cycle time

- Access time plus any additional time required before second access can commence
- Additional time may be required for transients to die out on signal lines or to regenerate data if they are read destructively
- Concerned with the system bus, not the processor

Transfer rate

- The rate at which data can be transferred into or out of a memory unit
- For random-access memory it is equal to $1/(\text{cycle time})$

Memory

- The most common forms are:
 - Semiconductor memory
 - Magnetic surface memory
 - 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

Copyright



This work is protected by United States copyright laws and is provided solely for the use of instructions in teaching their courses and assessing student learning. dissemination or sale of any part of this work (including on the World Wide Web) will destroy the integrity of the work and is not permitted. The work and materials from it should never be made available to students except by instructors using the accompanying text in their classes. All recipients of this work are expected to abide by these restrictions and to honor the intended pedagogical purposes and the needs of other instructors who rely on these materials.