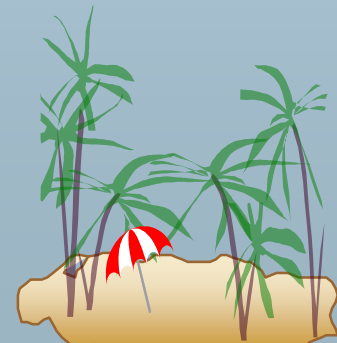# CS203 DB Principals

# Chapter 7:  Storage and File Structure
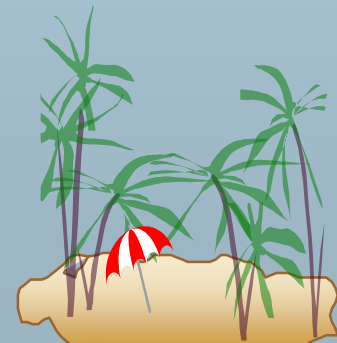
## *Asst.Prof : Asaad Alhijaj*

Reference:

**"Database System Concepts Fourth Edition" by Abraham Silberschatz Henry F. Korth S. Sudarshan , McGraw-Hill ISBN 0-07-255481-9**

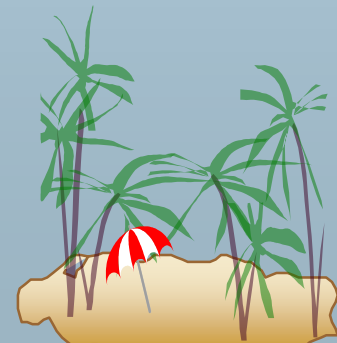# Chapter 7: Storage and File Structure

- Overview of Physical Storage Media
- Magnetic Disks
- RAID
- Tertiary Storage
- Storage Access
- File Organization
- Organization of Records in Files
- Data-Dictionary Storage
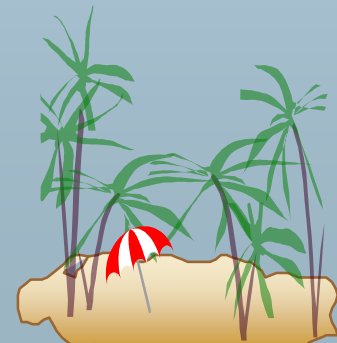
# Classification of Physical Storage Media

- Speed with which data can be accessed

- Cost per unit of data

- Reliability
    - data loss on power failure or system crash
    - physical failure of the storage device

- Can differentiate storage into:
    - **volatile storage**: loses contents when power is switched off
    - **non-volatile storage**:
        - Contents persist even when power is switched off.
        - Includes secondary and tertiary storage, as well as batter-backed up main-memory.

# Physical Storage Media

- **Cache** – fastest and most costly form of storage; volatile; managed by the computer system hardware.

- **Main memory**:
  - fast access (10s to 100s of nanoseconds; 1 nanosecond = $10^{-9}$ seconds)
  - generally too small (or too expensive) to store the entire database
    - capacities of up to a few Gigabytes widely used currently
    - Capacities have gone up and per-byte costs have decreased steadily and rapidly  (roughly factor of 2 every 2 to 3 years)
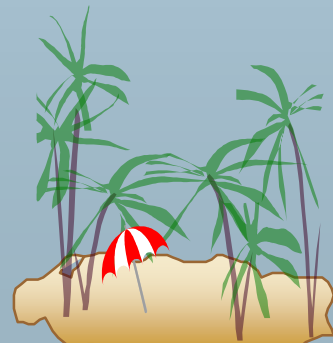  - **Volatile** — contents of main memory are usually lost if a power failure or system crash occurs.

# Physical Storage Media (Cont.)

- **Flash memory**
  - Data survives power failure
  - Data can be written at a location only once, but location can be erased and written to again
    - Can support only a limited number of write/erase cycles.
    - Erasing of memory has to be done to an entire bank of memory
  - Reads are roughly as fast as main memory
  - But writes are slow (few microseconds), erase is slower
  - Cost per unit of storage roughly similar to main memory
  - Widely used in embedded devices such as digital cameras
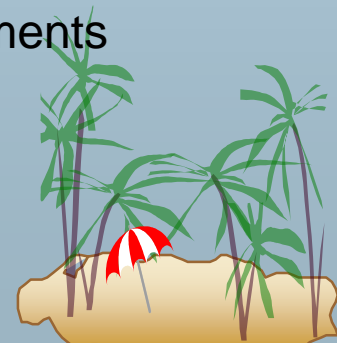  - also known as EEPROM (Electrically Erasable Programmable Read-Only Memory)

# Physical Storage Media (Cont.)

- **Magnetic-disk**
  - Data is stored on spinning disk, and read/written magnetically
  - Primary medium for the long-term storage of data; typically stores entire database.
  - Data must be moved from disk to main memory for access, and written back for storage
    - Much slower access than main memory (more on this later)
  - **direct-access** – possible to read data on disk in any order, unlike magnetic tape
  - Hard disks vs floppy disks
  - Capacities range up to roughly 100 GB currently
    - Much larger capacity and cost/byte than main memory/flash memory
    - Growing constantly and rapidly with technology improvements (factor of 2 to 3 every 2 years)
  - Survives power failures and system crashes
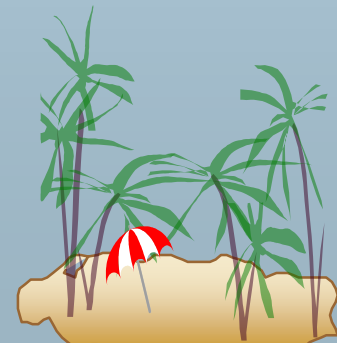    - disk failure can destroy data, but is very rare

# Physical Storage Media (Cont.)

- **Optical storage**
  - non-volatile, data is read optically from a spinning disk using a laser
  - CD-ROM (640 MB) and DVD (4.7 to 17 GB) most popular forms
  - Write-one, read-many (WORM) optical disks used for archival storage (CD-R and DVD-R)
  - Multiple write versions also available (CD-RW, DVD-RW, and DVD-RAM)
  - Reads and writes are slower than with magnetic disk
  - **Juke-box** systems, with large numbers of removable disks, a few drives, and a mechanism for automatic loading/unloading of disks available for storing large volumes of data
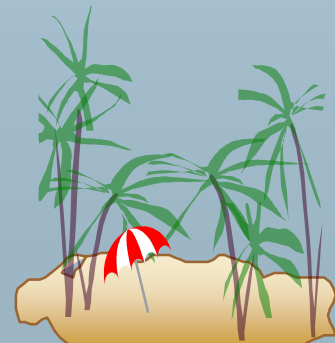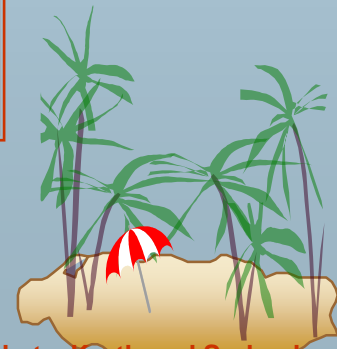
# Physical Storage Media (Cont.)

- **Tape storage**
  - non-volatile, used primarily for backup (to recover from disk failure), and for archival data
  - **sequential-access** – much slower than disk
  - very high capacity (40 to 300 GB tapes available)
  - tape can be removed from drive $\Rightarrow$ storage costs much cheaper than disk, but drives are expensive
  - Tape jukeboxes available for storing massive amounts of data
    - hundreds of terabytes (1 terabyte = $10^9$ bytes) to even a petabyte (1 petabyte = $10^{12}$ bytes)

# Storage Hierarchy

# Storage Hierarchy (Cont.)

- **primary storage:** Fastest media but volatile (cache, main memory).

- **secondary storage**: next level in hierarchy, non-volatile, moderately fast access time
  - ☞ also called **on-line storage**
  - ☞ E.g. flash memory, magnetic disks

- **tertiary storage**: lowest level in hierarchy, non-volatile, slow access time
  - ☞ also called **off-line storage**
  - ☞ E.g. magnetic tape, optical storage

# Magnetic Hard Disk Mechanism



**NOTE: Diagram is schematic, and simplifies the structure of actual disk drives**

# Magnetic Disks

- **Read-write head**
    - Positioned very close to the platter surface (almost touching it)
    - Reads or writes magnetically encoded information.
- Surface of platter divided into circular **tracks**
    - Over 16,000 tracks per platter on typical hard disks
- Each track is divided into **sectors**.
    - A sector is the smallest unit of data that can be read or written.
    - Sector size typically 512 bytes
    - Typical sectors per track: 200 (on inner tracks) to 400 (on outer tracks)
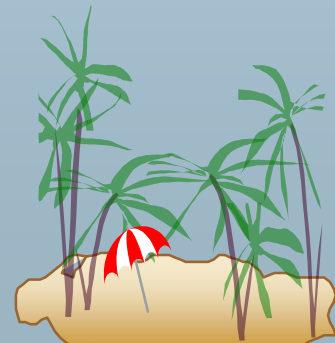- To read/write a sector
    - disk arm swings to position head on right track
    - platter spins continually; data is read/written as sector passes under head
- Head-disk assemblies
    - multiple disk platters on a single spindle (typically 2 to 4)
    - one head per platter, mounted on a common arm.
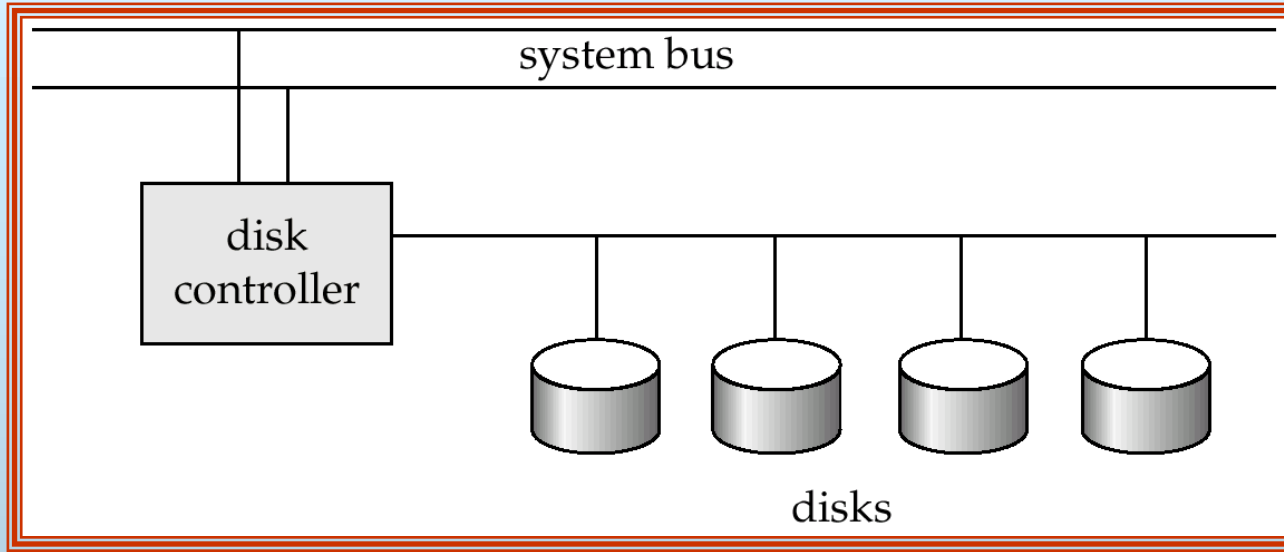- **Cylinder** $i$ consists of $i^{th}$ track of all the platters

# Magnetic Disks (Cont.)

- Earlier generation disks were susceptible to head-crashes
  - Surface of earlier generation disks had metal-oxide coatings which would disintegrate on head crash and damage all data on disk
  - Current generation disks are less susceptible to such disastrous failures, although individual sectors may get corrupted

- **Disk controller** – interfaces between the computer system and the disk drive hardware.
  - accepts high-level commands to read or write a sector
  - initiates actions such as moving the disk arm to the right track and actually reading or writing the data
  - Computes and attaches **checksums** to each sector to verify that data is read back correctly
    - If data is corrupted, with very high probability stored checksum won't match recomputed checksum
  - Ensures successful writing by reading back sector after writing it
  - Performs remapping of bad sectors

# Disk Subsystem



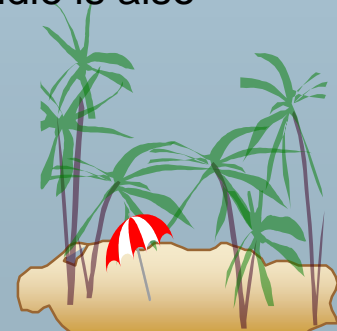- Multiple disks connected to a computer system through a controller
  - Controllers functionality (checksum, bad sector remapping) often carried out by individual disks; reduces load on controller
- Disk interface standards families
  - ATA (AT adaptor) range of standards
  - SCSI (Small Computer System Interconnect) range of standards
  - Several variants of each standard (different speeds and capabilities)
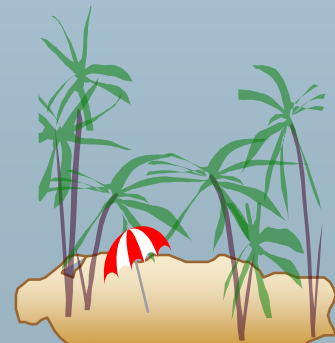
# Performance Measures of Disks

- **Access time** – the time it takes from when a read or write request is issued to when data transfer begins.  Consists of:
  - **Seek time** – time it takes to reposition the arm over the correct track.
    - Average seek time is 1/2 the worst case seek time.
      - Would be 1/3 if all tracks had the same number of sectors, and we ignore the time to start and stop arm movement
    - 4 to 10 milliseconds on typical disks
  - **Rotational latency** – time it takes for the sector to be accessed to appear under the head.
    - Average latency is 1/2 of the worst case latency.
    - 4 to 11 milliseconds on typical disks (5400 to 15000 r.p.m.)
- **Data-transfer rate** – the rate at which data can be retrieved from or stored to the disk.
  - 4 to 8 MB per second is typical
  - Multiple disks may share a controller, so rate that controller can handle is also important
    - E.g. ATA-5: 66 MB/second,  SCSI-3: 40 MB/s
    - Fiber Channel: 256 MB/s

# Performance Measures (Cont.)

- **Mean time to failure (MTTF)** – the average time the disk is expected to run continuously without any failure.

  - Typically 3 to 5 years

  - Probability of failure of new disks is quite low, corresponding to a "theoretical MTTF" of 30,000 to 1,200,000 hours for a new disk

    - E.g., an MTTF of 1,200,000 hours for a new disk means that given 1000 relatively new disks, on an average one will fail every 1200 hours

  - MTTF decreases as disk ages

# RAID

- **RAID: Redundant Arrays of Independent Disks**
    - disk organization techniques that manage a large numbers of disks, providing a view of a single disk of
        - high capacity and high speed by using multiple disks in parallel, and
        - high reliability by storing data redundantly, so that data can be recovered even if a disk fails

- The chance that some disk out of a set of $N$ disks will fail is much higher than the chance that a specific single disk will fail.
    - E.g., a system with 100 disks, each with MTTF of 100,000 hours (approx. 11 years), will have a system MTTF of 1000 hours (approx. 41 days)
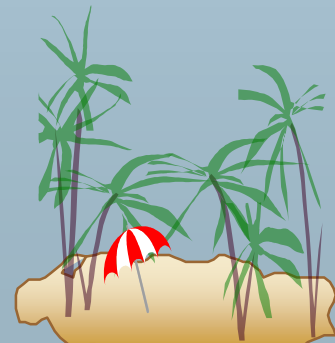    - Techniques for using redundancy to avoid data loss are critical with large numbers of disks

- Originally a cost-effective alternative to large, expensive disks
    - I in RAID originally stood for ``inexpensive''
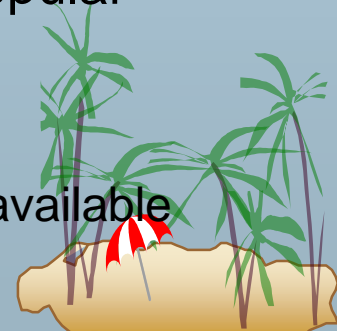    - Today RAIDs are used for their higher reliability and bandwidth.
        - The "I" is interpreted as independent

# Optical Disks

- Compact disk-read only memory (CD-ROM)
    - Disks can be loaded into or removed from a drive
    - High storage capacity (640 MB per disk)
    - High seek times or about 100 msec (optical read head is heavier and slower)
    - Higher latency (3000 RPM) and lower data-transfer rates (3-6 MB/s) compared to magnetic disks
- Digital Video Disk (DVD)
    - DVD-5  holds 4.7 GB , and DVD-9 holds 8.5 GB
    - DVD-10 and DVD-18 are double sided formats with capacities of 9.4 GB and 17 GB
    - Other characteristics similar to CD-ROM
- Record once versions (CD-R and DVD-R) are becoming popular
    - data can only be written once, and cannot be erased.
    - high capacity and long lifetime; used for archival storage
    - Multi-write versions (CD-RW, DVD-RW and DVD-RAM) also available

# Magnetic Tapes

- Hold large volumes of data and provide high transfer rates
  - Few GB for DAT (Digital Audio Tape) format, 10-40 GB with DLT (Digital Linear Tape) format, 100 GB+ with Ultrium format, and 330 GB with Ampex helical scan format
  - Transfer rates from few to 10s of MB/s
- Currently the cheapest storage medium
  - Tapes are cheap, but cost of drives is very high
- Very slow access time in comparison to magnetic disks and optical disks
  - limited to sequential access.
  - Some formats (Accelis) provide faster seek (10s of seconds) at cost of lower capacity
- Used mainly for backup, for storage of infrequently used information, and as an off-line medium for transferring information from one system to another.
- Tape jukeboxes used for very large capacity storage
  - (terabyte ($10^{12}$ bytes) to petabye ($10^{15}$ bytes)

# Storage Access

- A database file is partitioned into fixed-length storage units called **blocks**.  Blocks are units of both storage allocation and data transfer.

- Database system seeks to minimize the number of block transfers between the disk and memory.  We can reduce the number of disk accesses by keeping as many blocks as possible in main memory.

- **Buffer** – portion of main memory available to store copies of disk blocks.

- **Buffer manager** – subsystem responsible for allocating buffer space in main memory.

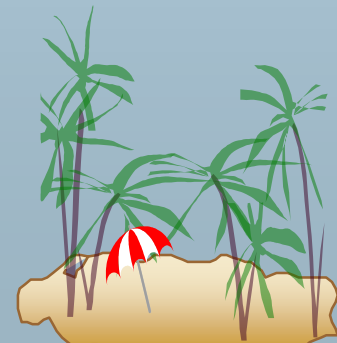# File Organization

- The database is stored as a collection of *files*.  Each file is a sequence of *records.*  A record is a sequence of fields.

- One approach:

    - assume record size is fixed

    - each file has records of one particular type only

    - different files are used for different relations

    This case is easiest to implement; will consider variable length records later.

# Fixed-Length Records

- Simple approach:
  - Store record $i$ starting from byte $n * (i - 1)$, where $n$ is the size of each record.
  - Record access is simple but records may cross blocks
    - Modification: do not allow records to cross block boundaries

- Deletion of record $l$: alternatives:
  - move records $i + 1, \ldots, n$ to $i, \ldots, n - 1$
  - move record $n$ to $i$
  - do not move records, but link all free records on a *free list*

| | | | |
|---|---|---|---|
| record 0 | A-102 | Perryridge | 400 |
| record 1 | A-305 | Round Hill | 350 |
| record 2 | A-215 | Mianus | 700 |
| record 3 | A-101 | Downtown | 500 |
| record 4 | A-222 | Redwood | 700 |
| record 5 | A-201 | Perryridge | 900 |
| record 6 | A-217 | Brighton | 750 |
| record 7 | A-110 | Downtown | 600 |
| record 8 | A-218 | Perryridge | 700 |

# Free Lists

- Store the address of the first deleted record in the file header.

- Use this first record to store the address of the second deleted record, and so on

- Can think of these stored addresses as pointers since they "point" to the location of a record.

- More space efficient representation: reuse space for normal attributes of free records to store pointers. (No pointers stored in in-use records.)

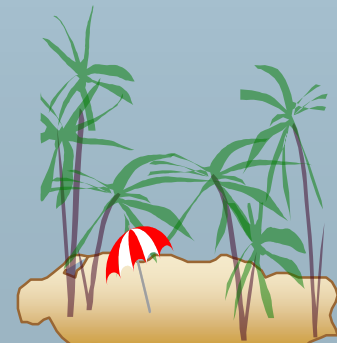| | | | |
|---|---|---|---|
| header | | | |
| record 0 | A-102 | Perryridge | 400 |
| record 1 | | | |
| record 2 | A-215 | Mianus | 700 |
| record 3 | A-101 | Downtown | 500 |
| record 4 | | | |
| record 5 | A-201 | Perryridge | 900 |
| record 6 | | | |
| record 7 | A-110 | Downtown | 600 |
| record 8 | A-218 | Perryridge | 700 |

# Variable-Length Records

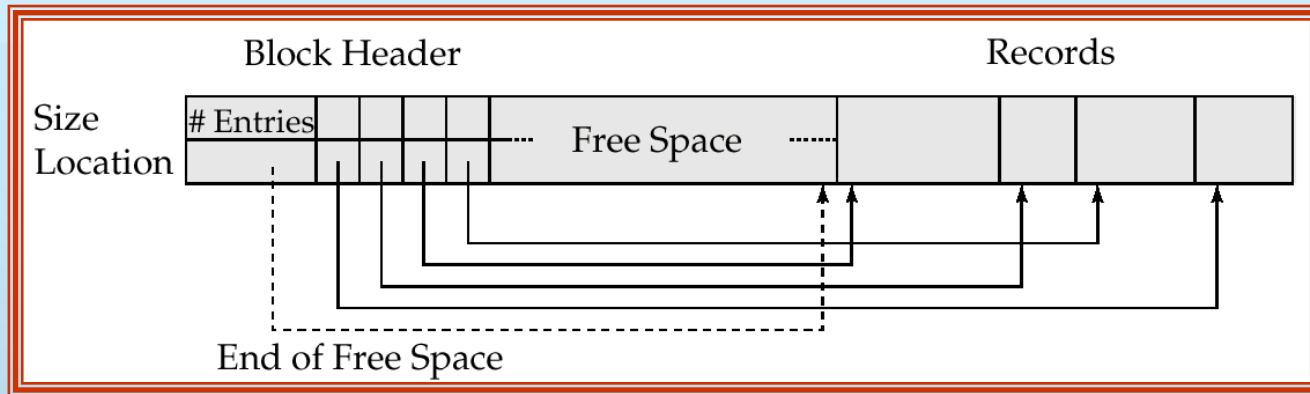- Variable-length records arise in database systems in several ways:
  - Storage of multiple record types in a file.
  - Record types that allow variable lengths for one or more fields.
  - Record types that allow repeating fields (used in some older data models).

- Byte string representation
  - Attach an *end-of-record* (⊥) control character to the end of each record
  - Difficulty with deletion
  - Difficulty with growth

# Variable-Length Records: Slotted Page Structure



- Slotted page header contains:
  - number of record entries
  - end of free space in the block
  - location and size of each record
- Records can be moved around within a page to keep them contiguous with no empty space between them; entry in the header must be updated.
- Pointers should not point directly to record — instead they should point to the entry for the record in header.

- Fixed-length representation:
  - ☞ reserved space
  - ☞ pointers

- Reserved space – can use fixed-length records of a known maximum length; unused space in shorter records filled with a null or end-of-record symbol.

| 0 | Perryridge | A-102 | 400 | A-201 | 900 | A-218 | 700 |
|---|------------|-------|-----|-------|-----|-------|-----|
| 1 | Round Hill | A-305 | 350 | ⊥ | ⊥ | ⊥ | ⊥ |
| 2 | Mianus | A-215 | 700 | ⊥ | ⊥ | ⊥ | ⊥ |
| 3 | Downtown | A-101 | 500 | A-110 | 600 | ⊥ | ⊥ |
| 4 | Redwood | A-222 | 700 | ⊥ | ⊥ | ⊥ | ⊥ |
| 5 | Brighton | A-217 | 750 | ⊥ | ⊥ | ⊥ | ⊥ |

# Organization of Records in Files

- **Heap** – a record can be placed anywhere in the file where there is space

- **Sequential** – store records in sequential order, based on the value of the search key of each record

- **Hashing** – a hash function computed on some attribute of each record; the result specifies in which block of the file the record should be placed

- Records of each relation may be stored in a separate file. In a **clustering file organization** records of several different relations can be stored in the same file

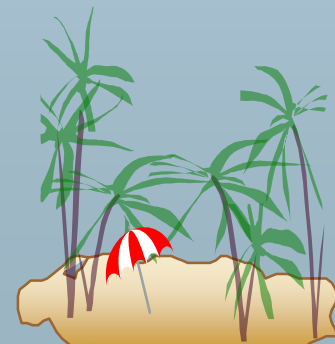  - Motivation: store related records on the same block to minimize I/O

# Sequential File Organization

■ Suitable for applications that require sequential processing of the entire file

■ The records in the file are ordered by a search-key

| A-217 | Brighton | 750 | |
|-------|----------|-----|---|
| A-101 | Downtown | 500 | |
| A-110 | Downtown | 600 | |
| A-215 | Mianus | 700 | |
| A-102 | Perryridge | 400 | |
| A-201 | Perryridge | 900 | |
| A-218 | Perryridge | 700 | |
| A-222 | Redwood | 700 | |
| A-305 | Round Hill | 350 | |

# Sequential File Organization (Cont.)

- Deletion – use pointer chains

- Insertion –locate the position where the record is to be inserted

  - ☞ if there is free space insert there

  - ☞ if no free space, insert the record in an overflow block

  - ☞ In either case, pointer chain must be updated

- Need to reorganize the file from time to time to restore sequential order

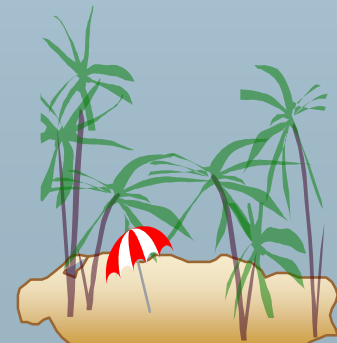| A-217 | Brighton | 750 | |
|-------|----------|-----|---|
| A-101 | Downtown | 500 | |
| A-110 | Downtown | 600 | |
| A-215 | Mianus | 700 | |
| A-102 | Perryridge | 400 | |
| A-201 | Perryridge | 900 | |
| A-218 | Perryridge | 700 | |
| A-222 | Redwood | 700 | |
| A-305 | Round Hill | 350 | |
| | | | |
| A-888 | North Town | 800 | |

# Clustering File Organization

- Simple file structure stores each relation in a separate file

- Can instead store several relations in one file using a **clustering** file organization

- E.g., clustering organization of *customer* and *depositor:*

| | | |
|---|---|---|
| Hayes | Main | Brooklyn |
| Hayes | A-102 | |
| Hayes | A-220 | |
| Hayes | A-503 | |
| Turner | Putnam | Stamford |
| Turner | A-305 | |

- ☞ good for queries involving depositor ⋈ customer, and for queries involving one single customer and his accounts
- ☞ bad for queries involving only customer
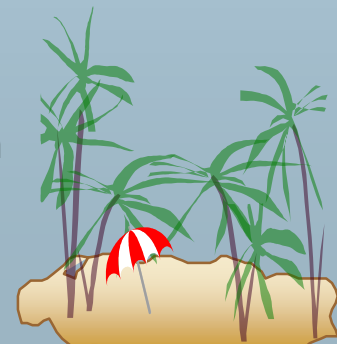- ☞ results in variable size records

# Data Dictionary Storage

Data dictionary (also called system catalog) stores metadata: that is, data about data, such as

- Information about relations
  - names of relations
  - names and types of attributes of each relation
  - names and definitions of views
  - integrity constraints
- User and accounting information, including passwords
- Statistical and descriptive data
  - number of tuples in each relation
- Physical file organization information
  - How relation is stored (sequential/hash/…)
  - Physical location of relation
    - operating system file name or
    - disk addresses of blocks containing records of the relation
- Information about indices (Chapter 12)

# Data Dictionary Storage (Cont.)

- Catalog structure:  can use either

    - specialized data structures designed for efficient access

    - a set of relations, with existing system features used to ensure efficient access

    The latter alternative is usually preferred

- A possible catalog representation:

*Relation-metadata = (<u>relation-name</u>, number-of-attributes,*
*storage-organization, location)*
*Attribute-metadata = (<u>attribute-name, relation-name</u>, domain-type,*
*position, length)*
*User-metadata = (<u>user-name</u>, encrypted-password, group)*
*Index-metadata = (<u>index-name, relation-name</u>, index-type,*
*index-attributes)*
*View-metadata = (<u>view-name</u>, definition)*