

# Fundamentals of DATABASE SYSTEMS

FOURTH EDITION

ELMASRI  NAVATHE

## Chapter 13

**Disk Storage, Basic File Structures, and  
Hashing.**



Copyright © 2004 Pearson Education, Inc.

## Chapter Outline

- Disk Storage Devices
- Files of Records
- Operations on Files
- Unordered Files
- Ordered Files
- Hashed Files
  - Dynamic and Extendible Hashing Techniques
- RAID Technology

## Disk Storage Devices (cont.)

- Preferred secondary storage device for high storage capacity and low cost.
- Data stored as magnetized areas on magnetic disk surfaces.
- A *disk pack* contains several magnetic disks connected to a rotating spindle.
- Disks are divided into concentric circular *tracks* on each disk *surface*. Track capacities vary typically from 4 to 50 Kbytes.

## Disk Storage Devices (cont.)

Because a track usually contains a large amount of information, it is divided into smaller blocks or sectors.

- The division of a track into *sectors* is hard-coded on the disk surface and cannot be changed. One type of sector organization calls a portion of a track that subtends a fixed angle at the center as a sector.
- A track is divided into *blocks*. The block size  $B$  is fixed for each system. Typical block sizes range from  $B=512$  bytes to  $B=4096$  bytes. Whole blocks are transferred between disk and main memory for processing.

## Disk Storage Devices (cont.)

INSERT FIGURE 13.2 here

## Disk Storage Devices (cont.)

- A *read-write* head moves to the track that contains the block to be transferred. Disk rotation moves the block under the read-write head for reading or writing.
- A physical disk block address consists of a surface number, track number (within surface), and block number (within track).
- Reading or writing a disk block is time consuming because of the seek time *s* and rotational delay (latency) **rd**.
- Double buffering can be used to speed up the transfer of contiguous disk blocks.

## Disk Storage Devices (cont.)

INSERT FIGURE 13.1 here

## Files of Records

- A file is a *sequence* of records, where each record is a collection of data values (or data items).
- A *file descriptor* (or *file header* ) includes information that describes the file, such as the *field names* and their *data types*, and the addresses of the file blocks on disk.
- Records are stored on disk blocks. The *blocking factor* BFR for a file is the (average) number of file records stored in a disk block.
- A file can have *fixed-length* records or *variable-length* records.

## Files of Records (cont.)

- File records can be *unspanned* (no record can span two blocks) or *spanned* (a record can be stored in more than one block).
- The physical disk blocks that are allocated to hold the records of a file can be *contiguous*, *linked*, or *indexed*.
- In a file of fixed-length records, all records have the same format. Usually, unspanned blocking is used with such files.
- Files of variable-length records require additional information to be stored in each record, such as *separator characters* and *field types*. Usually spanned blocking is used with such files.

## Operation on Files

Typical file operations include:

- **OPEN:** Readies the file for access, and associates a pointer that will refer to a *current* file record at each point in time.
- **FIND:** Searches for the first file record that satisfies a certain condition, and makes it the current file record.
- **FINDNEXT:** Searches for the next file record (from the current record) that satisfies a certain condition, and makes it the current file record.
- **READ:** Reads the current file record into a program variable.
- **INSERT:** Inserts a new record into the file, and makes it the current file record.

## Operation on Files (cont.)

- **DELETE:** Removes the current file record from the file, usually by marking the record to indicate that it is no longer valid.
- **MODIFY:** Changes the values of some fields of the current file record.
- **CLOSE:** Terminates access to the file.
- **REORGANIZE:** Reorganizes the file records. For example, the records marked deleted are physically removed from the file or a new organization of the file records is created.
- **READ\_ORDERED:** Read the file blocks in order of a specific field of the file.

## Unordered Files

- Also called a *heap* or a *pile* file.
- New records are inserted at the end of the file.
- To search for a record, a *linear search* through the file records is necessary. This requires reading and searching half the file blocks on the average, and is hence quite expensive.
- Record insertion is quite efficient.
- Reading the records in order of a particular field requires sorting the file records.

## Ordered Files

- Also called a *sequential file*.
- File records are kept sorted by the values of an *ordering field*.
- Insertion is expensive: records must be inserted in the *correct order*. It is common to keep a separate unordered *overflow* (or *transaction*) file for new records to improve insertion efficiency; this is periodically merged with the main ordered file.
- A *binary search* can be used to search for a record on its *ordering field value*. This requires reading and searching  $\log_2$  of the file blocks on the average, an improvement over linear search.
- Reading the records in order of the ordering field is quite efficient.

## Ordered Files (cont.)

**INSERT FIGURE 13.7**

## Hashed Files

- The file blocks are divided into  $M$  equal-sized *buckets*, numbered  $\text{bucket}_0, \text{bucket}_1, \dots, \text{bucket}_{M-1}$ . Typically, a bucket corresponds to one (or a fixed number of) disk block.
- One of the file fields is designated to be the hash key of the file.
- The record with hash key value  $K$  is stored in bucket  $i$ , where  $i=h(K)$ , and  $h$  is the *hashing function*.
- Search is very efficient on the hash key.
- Collisions occur when a new record hashes to a bucket that is already full. An overflow file is kept for storing such records. Overflow records that hash to each bucket can be linked together.



## Hashed Files (cont.)

There are numerous methods for collision resolution, including the following:

- **Open addressing:** Proceeding from the occupied position specified by the hash address, the program checks the subsequent positions in order until an unused (empty) position is found.
- **Chaining:** For this method, various overflow locations are kept, usually by extending the array with a number of overflow positions. In addition, a pointer field is added to each record location. A collision is resolved by placing the new record in an unused overflow location and setting the pointer of the occupied hash address location to the address of that overflow location.
- **Multiple hashing:** The program applies a second hash function if the first results in a collision. If another collision results, the program uses open addressing or applies a third hash function and then uses open addressing if necessary.

## Hashed Files (cont.)

**Insert figure 13.9**

## Hashed Files (cont.)

- To reduce overflow records, a hash file is typically kept 70-80% full.
- The hash function  $h$  should distribute the records uniformly among the buckets; otherwise, search time will be increased because many overflow records will exist.
- Main disadvantages of static external hashing:
  - Fixed number of buckets  $M$  is a problem if the number of records in the file grows or shrinks.
  - Ordered access on the hash key is quite inefficient (requires sorting the records).

## Hashed Files (cont.)

**INSERT FIGURE 13.10**

## Dynamic And Extendible Hashed Files

### Dynamic and Extendible Hashing Techniques

- Hashing techniques are adapted to allow the dynamic growth and shrinking of the number of file records.
- These techniques include the following: *dynamic hashing*, *extendible hashing*, and *linear hashing*.
- Both dynamic and extendible hashing use the *binary representation* of the hash value  $h(K)$  in order to access a *directory*. In dynamic hashing the directory is a binary tree. In extendible hashing the directory is an array of size  $2^d$  where  $d$  is called the *global depth*.

## Dynamic And Extendible Hashing (cont.)

- The directories can be stored on disk, and they expand or shrink dynamically. Directory entries point to the disk blocks that contain the stored records.
- An insertion in a disk block that is full causes the block to split into two blocks and the records are redistributed among the two blocks. The directory is updated appropriately.
- Dynamic and extendible hashing do not require an overflow area.
- Linear hashing does require an overflow area but does not use a directory. Blocks are split in *linear order* as the file expands.

## Extendible Hashing

**INSERT FIGURE 13.11**

## Parallelizing Disk Access using RAID Technology.

- Secondary storage technology must take steps to keep up in performance and reliability with processor technology.
- A major advance in secondary storage technology is represented by the development of **RAID**, which originally stood for **Redundant Arrays of Inexpensive Disks**.
- The main goal of raid is to even out the widely different rates of performance improvement of disks against those in memory and microprocessors.

## RAID Technology (cont.)

- A natural solution is a large array of small independent disks acting as a single higher-performance logical disk. A concept called **data striping** is used, which utilizes *parallelism* to improve disk performance.
- Data striping distributes data transparently over multiple disks to make them appear as a single large, fast disk.

**Insert Figure 13.12**

## RAID Technology (cont.)

Different raid organizations were defined based on different combinations of the two factors of granularity of data interleaving (striping) and pattern used to compute redundant information.

- Raid level 0 has no redundant data and hence has the best write performance.
- Raid level 1 uses mirrored disks.
- Raid level 2 uses memory-style redundancy by using Hamming codes, which contain parity bits for distinct overlapping subsets of components. Level 2 includes both error detection and correction.
- Raid level 3 uses a single parity disk relying on the disk controller to figure out which disk has failed.
- Raid Levels 4 and 5 use block-level data striping, with level 5 distributing data and parity information across all disks.
- Raid level 6 applies the so-called  $P + Q$  redundancy scheme using Reed-Soloman codes to protect against up to two disk failures by using just two redundant disks.

## Use of RAID Technology (cont.)

Different raid organizations are being used under different situations

- Raid level 1 (mirrored disks) is the easiest for rebuild of a disk from other disks
  - It is used for critical applications like logs
- Raid level 2 uses memory-style redundancy by using Hamming codes, which contain parity bits for distinct overlapping subsets of components. Level 2 includes both error detection and correction.
- Raid level 3 (single parity disks relying on the disk controller to figure out which disk has failed) and level 5 (block-level data striping) are preferred for Large volume storage, with level 3 giving higher transfer rates.
- Most popular uses of the RAID technology currently are: Level 0 (with striping), Level 1 (with mirroring) and Level 5 with an extra drive for parity.
- Design Decisions for RAID include – level of RAID, number of disks, choice of parity schemes, and grouping of disks for block-level striping.