

## Record Storage and Primary File Organization

The collection of data that makes up a computerized database must be stored physically on some computer **storage medium**. The DBMS software can then retrieve, update, and process this data as needed. Computer storage media form a *storage hierarchy* that includes two main categories:

**Primary storage.** This category includes storage media that can be operated on directly by the computer's *central processing unit* (CPU), such as the computer's main memory and smaller but faster cache memories. Primary storage usually provides fast access to data but is of limited storage capacity. Although main memory capacities have been growing rapidly in recent years, they are still more expensive and have less storage capacity than secondary and tertiary storage devices.

**Secondary and tertiary storage.** This category includes magnetic disks, optical disks (CD-ROMs, DVDs, and other similar storage media), and tapes. Hard-disk drives are classified as secondary storage, whereas removable media such as optical disks and tapes are considered tertiary storage. These devices usually have a larger capacity, cost less, and provide slower access to data than do primary storage devices. Data in secondary or tertiary storage cannot be processed directly by the CPU; first it must be copied into primary storage and then processed by the CPU.

## Memory Hierarchies and Storage Devices

In a modern computer system, data resides and is transported throughout a hierarchy of storage media. The highest-speed memory is the most expensive and is therefore available with the least capacity. The lowest-speed memory is offline tape storage, which is essentially available in indefinite storage capacity.

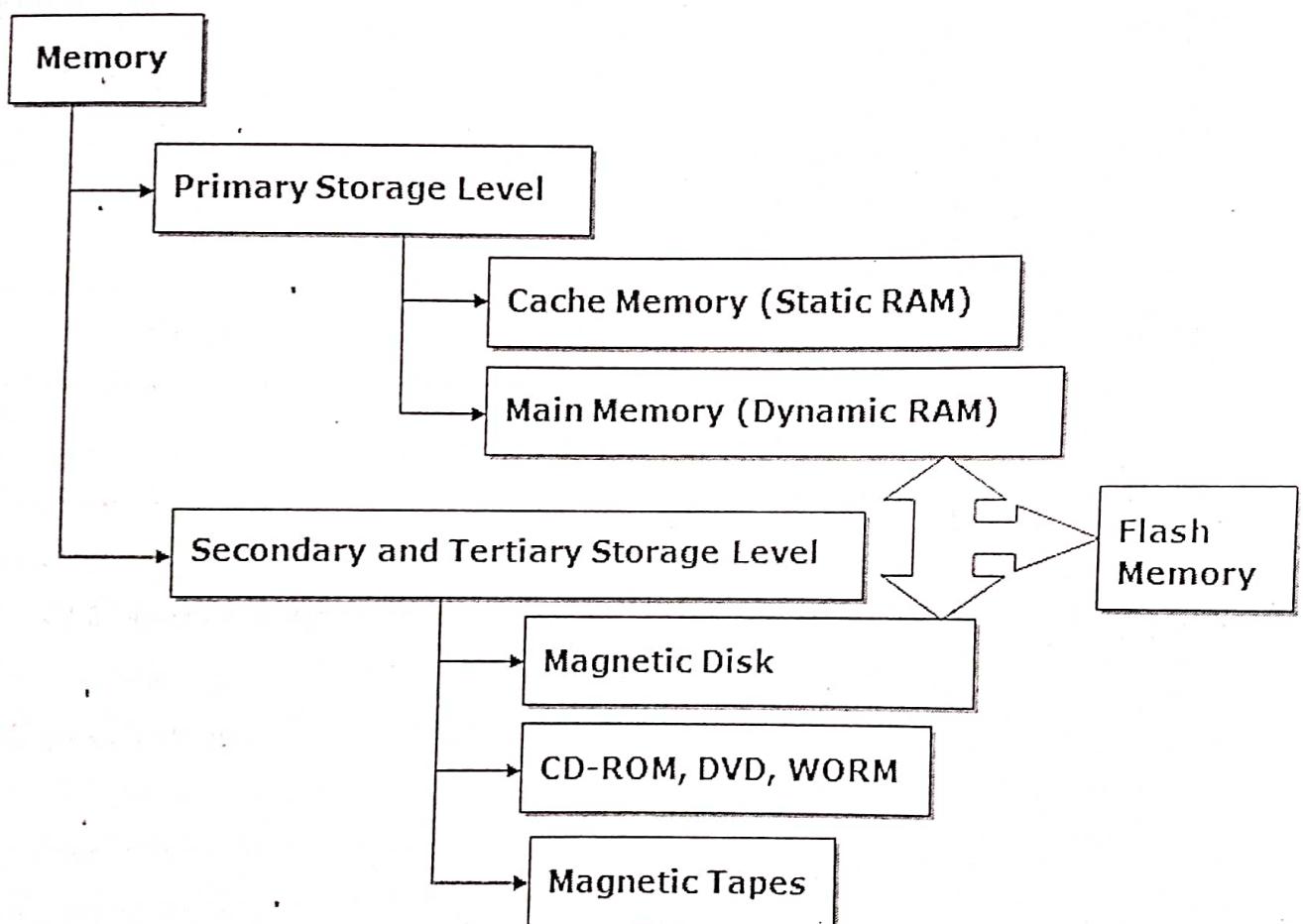


Fig: Memory Hierarchy

At the *primary storage level*, the memory hierarchy includes at the most expensive end, **cache memory**, which is a static RAM (Random Access Memory). Cache memory is typically used by the CPU to speed up execution of program instructions using techniques such as prefetching and pipelining.

The next level of primary storage is DRAM (Dynamic RAM), which provides the main work area for the CPU for keeping program instructions and data. It is popularly called **main memory**. The advantage of DRAM is its low cost, which continues to decrease; the drawback is its volatility<sup>1</sup> and lower speed compared with static RAM.

At the *secondary and tertiary storage level*, the hierarchy includes magnetic disks, as well as **mass storage** in the form of CD-ROM (Compact Disk-Read-Only Memory) and DVD (Digital Video Disk or Digital Versatile Disk) devices, and finally tapes at the least expensive end of the hierarchy. The **storage capacity** is measured in kilobytes (Kbyte or 1000 bytes), megabytes (MB or 1 million bytes), gigabytes (GB or 1 billion bytes), and even terabytes (1000 GB). The word petabyte (1000 terabytes or  $10^{15}$  bytes) is now becoming relevant in the context of very large repositories of data in physics, astronomy, earth sciences, and other scientific applications.

Programs reside and execute in DRAM. Generally, large permanent databases reside on secondary storage, (magnetic disks), and portions of the database are read into and written from buffers in main memory as needed. Nowadays, personal computers and workstations have large main memories of hundreds of megabytes of RAM and DRAM, so it is becoming possible to load a large part of the database into main memory.

Between DRAM and magnetic disk storage, another form of memory, **flash memory**, is becoming common, particularly because it is nonvolatile. Flash memories are high-density, high-performance memories using EEPROM (Electrically Erasable Programmable Read-Only Memory) technology.

The advantage of flash memory is the fast access speed; the disadvantage is that an entire block must be erased and written over simultaneously. Flash memory cards are appearing as the data storage medium in appliances with capacities ranging from a few megabytes to a few gigabytes.

These are appearing in cameras, MP3 players, cell phones, PDAs, and so on. USB (Universal Serial Bus) flash drives have become the most portable medium for carrying data between personal computers; they have a flash memory storage device integrated with a USB interface.

**CD-ROM** (Compact Disk – Read Only Memory) disks store data optically and are read by a laser. CD-ROMs contain prerecorded data that cannot be overwritten.

**WORM** (Write-Once-Read-Many) disks are a form of optical storage used for archiving data; they allow data to be written once and read any number of times without the possibility of erasing. They hold about half a gigabyte of data per disk and last much longer than magnetic disks.

The **DVD** is another standard for optical disks allowing 4.5 to 15 GB of storage per disk. Most personal computer disk drives now read CDROM and DVD disks.

Typically, drives are CD-R (Compact Disk Recordable) that can create CD-ROMs and audio CDs (Compact Disks), as well as record on DVDs.

Finally, **magnetic tapes** are used for archiving and backup storage of data. **Tape jukeboxes**—which contain a bank of tapes that are catalogued and can be automatically loaded onto tape drives—are becoming popular as **tertiary storage** to hold terabytes of data. For example, NASA's EOS (Earth Observation Satellite) system stores archived databases. Many large organizations are already finding it normal to have terabyte-sized databases.

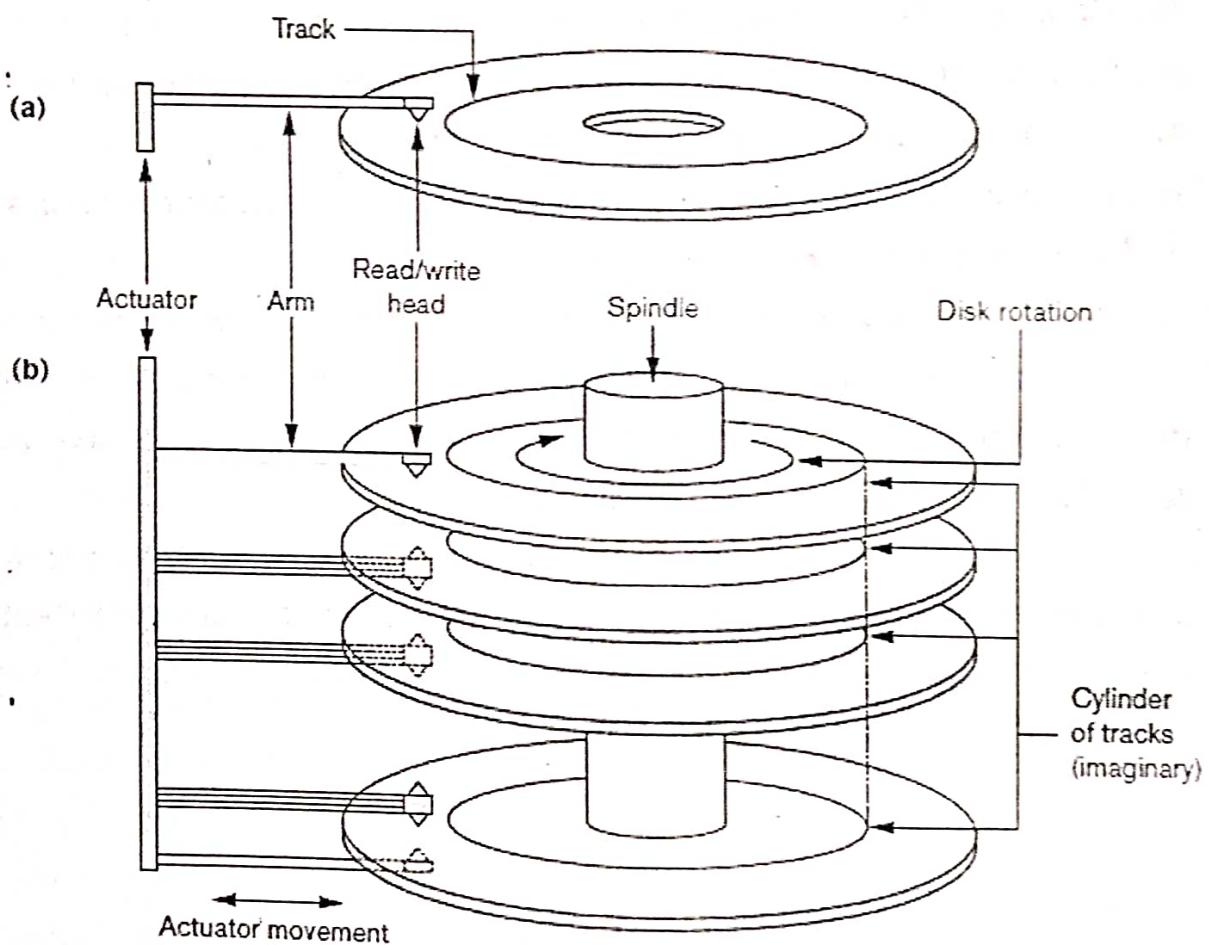
The term **very large database** can no longer be precisely defined because disk storage capacities are on the rise and costs are declining. Very soon the term may be reserved for databases containing tens of terabytes.

## Secondary Storage Devices

Magnetic disks are used for storing large amounts of data. The most basic unit of data on the disk is a single **bit** of information. Bits are grouped into **bytes** (or **characters**). Byte sizes are typically 4 to 8 bits, depending on the computer and the device. We assume that one character is stored in a single byte, and we use the terms *byte* and *character* interchangeably. The **capacity** of a disk is the number of bytes it can store, Hard disks for personal computers typically hold from several hundred MB up to tens of GB; and large disk packs used with servers and mainframes have capacities of hundreds of GB. Disk capacities continue to grow as technology improves.

**Figure**

- (a) A single-sided disk with read/write hardware.
- (b) A disk pack with read/write hardware.



Whatever their capacity, all disks are made of magnetic material shaped as a thin circular disk, as shown in Figure(a), and protected by a plastic or acrylic cover. A disk is **single-sided** if it stores information on one of its surfaces only and **doublesided** if both surfaces are used. To increase storage capacity, disks are assembled into a **disk pack**, as shown in Figure (b), which may include many disks and therefore many surfaces. Information is stored on a disk surface in concentric circles of *small width*, each having a distinct diameter. Each circle is called a **track**. In disk packs, tracks with the same diameter on the various surfaces are called a **cylinder** because of the shape they would form if connected in space. The concept of a cylinder is important because data stored on one cylinder can be retrieved much faster than if it were distributed among different cylinders.

The number of tracks on a disk ranges from a few hundred to a few thousand, and the capacity of each track typically ranges from tens of Kbytes to 150 Kbytes. 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.

The division of a track into equal-sized **disk blocks** (or **pages**) is set by the operating system during disk **formatting** (or **initialization**). Block size is fixed during initialization and cannot be changed dynamically. Typical disk block sizes range from 512 to 8192 bytes. A disk with hard-coded sectors often has the sectors subdivided into blocks during initialization. Blocks are separated by fixed-size **interblock gaps**, which include specially coded control information written during disk initialization. This information is used to determine which block on the track follows each interblock gap.

A disk is a *random access* addressable device. Transfer of data between main memory and disk takes place in units of disk blocks. The **hardware address** of a block—a combination of a cylinder number, track number (surface number within the cylinder on which the track is located), and block number (within the track) is supplied to the disk I/O (input/output) hardware.

The address of a **buffer**—a contiguous reserved area in main storage that holds one disk block—is also provided. For a **read** command, the disk block is copied into the buffer; whereas for a **write** command, the contents of the buffer are copied into the disk block. Sometimes several contiguous blocks, called a **cluster**, may be transferred as a unit. In this case, the buffer size is adjusted to match the number of bytes in the cluster.

The actual hardware mechanism that reads or writes a block is the disk **read/write head**, which is part of a system called a **disk drive**. A disk or disk pack is mounted in the disk drive, which includes a motor that rotates the disks. A read/write head includes an electronic component attached to a **mechanical arm**. Disk packs with multiple surfaces are controlled by several read/write heads—one for each surface, as shown in Figure (b). All arms are connected to an **actuator** attached to another electrical motor, which moves the read/write heads in unison and positions them precisely over the cylinder of tracks specified in a block address.

Disk drives for hard disks rotate the disk pack continuously at a constant speed (typically ranging between 5,400 and 15,000 rpm (revolutions per minute)). Once the read/write head is positioned on the right track and the block specified in the block address moves under the read/write head, the electronic component of the read/write head is activated to transfer the data. Some disk units have fixed read/write heads, with as many heads as there are tracks. These are called **fixed-head disks**, whereas disk units with an actuator are called **movable-head disks**.

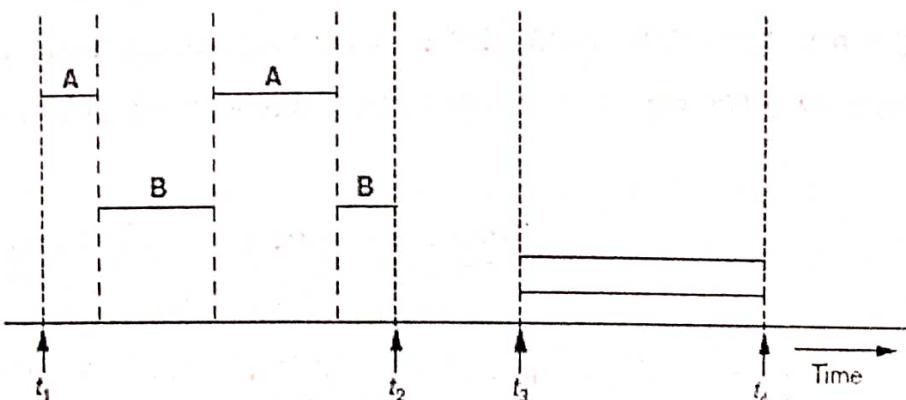
For fixed-head disks, a track or cylinder is selected by electronically switching to the appropriate read/write head rather than by actual mechanical movement; consequently, it is much faster. However, the cost of the additional read/write heads is quite high, so fixed-head disks are not commonly used.

A **disk controller**, typically embedded in the disk drive, controls the disk drive and interfaces it to the computer system. The controller accepts high-level I/O commands and takes appropriate action to position the arm and causes the read/write action to take place. To transfer a disk block, given its address, the disk controller must first mechanically position the read/write head on the correct track. The time required to do this is called the **seek time**. Typical seek times are 5 to 10 msec on desktops and 3 to 8 msecs on servers. There is another delay—called the **rotational delay or latency**—while the beginning of the desired block rotates into position under the read/write head. It depends on the rpm of the disk. Some additional time is needed to transfer the data; this is called the **block transfer time**. Hence, the total time needed to locate and transfer an arbitrary block, given its address, is the sum of the seek time, rotational delay, and block transfer time.

### **Buffering of Blocks**

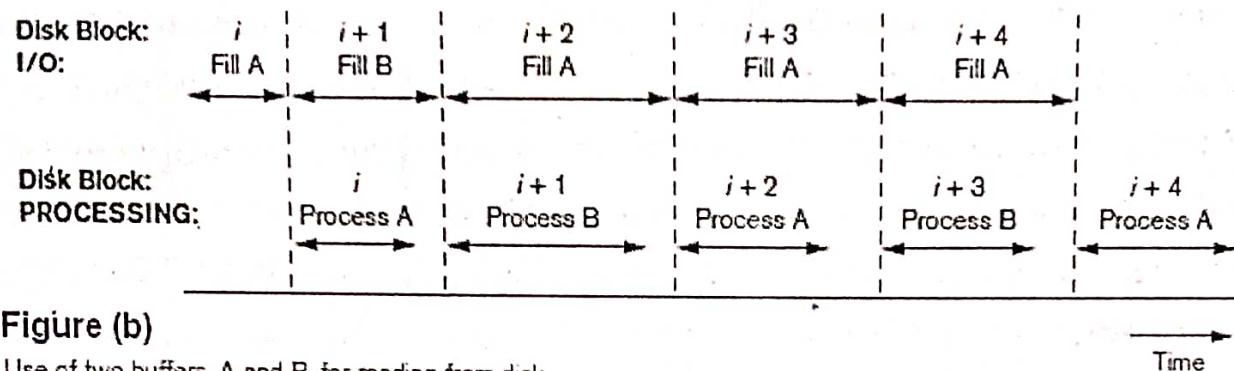
When several blocks need to be transferred from disk to main memory and all the block addresses are known, several buffers can be reserved in main memory to speed up the transfer. While one buffer is being read or written, the CPU can process data in the other buffer because an independent disk I/O processor (controller) exists that, once started, can proceed to transfer a data block between memory and disk independent of and in parallel to CPU processing.

Interleaved concurrency  
of operations A and B      Parallel execution of  
operations C and D



**Figure (a)**

Interleaved concurrency  
versus parallel execution.



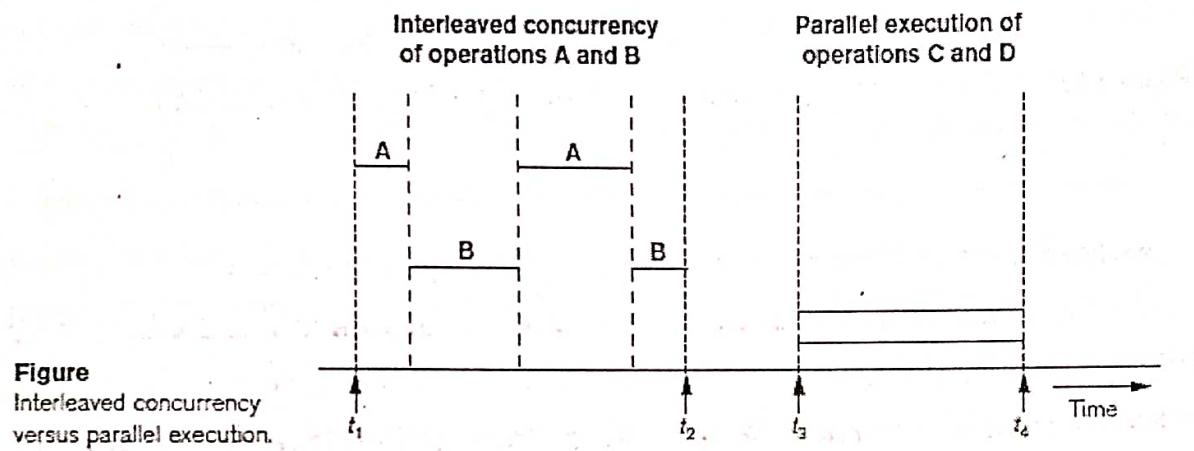
**Figure (b)**

Use of two buffers, A and B, for reading from disk.

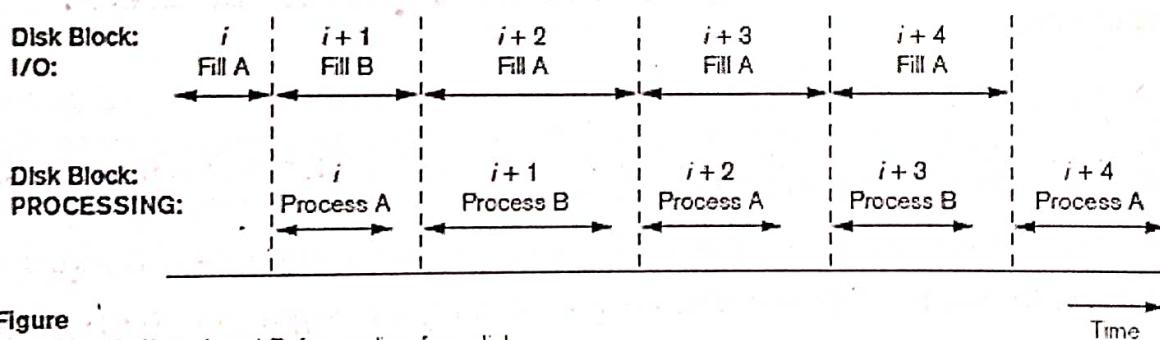
Figure (a) illustrates how two processes can proceed in parallel. Processes A and B are running **concurrently** in an **interleaved** fashion, whereas processes C and D are running **concurrently** in a **parallel** fashion. When a single CPU controls multiple processes, parallel execution is not possible. However, the processes can still run concurrently in an interleaved way. Buffering is most useful when processes can run concurrently in a parallel fashion, either because a separate disk I/O processor is available or because multiple CPU processors exist.

Figure (b) illustrates how reading and processing can proceed in parallel when the time required to process a disk block in memory is less than the time required to read the next block and fill a buffer. The CPU can start processing a block once its transfer to main memory is completed; at the same time, the disk I/O processor can be reading and transferring the next block into a different buffer. This technique is called **double buffering** and can also be used to read a continuous

stream of blocks from disk to memory. Double buffering permits continuous reading or writing of data on consecutive disk blocks, which eliminates the seek time and rotational delay for all but the first block transfer. Moreover, data is kept ready for processing, thus reducing the waiting time in the programs.



**Figure**  
Interleaved concurrency  
versus parallel execution.



**Figure**  
Use of two buffers, A and B, for reading from disk.

## Placing file Records on Disk

Data is usually stored in the form of **records**. Each record consists of a collection of related data **values** or **items**, where each value is formed of one or more bytes and corresponds to a particular **field** of the record. Records usually describe entities and their attributes.

For example, an **EMPLOYEE** record represents an employee entity, and each field value in the record specifies some attribute of that employee, such as EmployeeID, Name, DOB, Department, Salary.

A collection of field names and their corresponding data types constitutes a **record type** or **record format** definition.

A **data type**, associated with each field, specifies the types of values a field can take.

The data type of a field is usually include numeric (integer, long integer, or floating point), string of characters (fixed-length or varying), Boolean (having 0 and 1 or TRUE and FALSE values only), and sometimes specially coded **date** and **time** data types.

### **Files, Fixed-Length Records, and Variable-Length Records**

A **file** is a *sequence* of records. In many cases, all records in a file are of the same record type. If every record in the file has exactly the same size (in bytes), the file is said to be made up of **fixed-length records**. If different records in the file have different sizes, the file is said to be made up of **variable-length records**.

### **Allocating File Blocks on Disk**

There are several standard techniques for allocating the blocks of a file on disk.

In **contiguous allocation**, the file blocks are allocated to consecutive disk blocks. This makes reading the whole file very fast using double buffering, but it makes expanding the file difficult.

In **linked allocation**, each file block contains a pointer to the next file block. This makes it easy to expand the file but makes it slow to read the whole file. A combination of the two allocates **clusters** of consecutive disk blocks, and the clusters are linked. Clusters are sometimes called **file segments** or **extents**.

In **indexed allocation**, where one or more **index blocks** contain pointers to the actual file blocks. It is also common to use combinations of these techniques.

## **File Headers**

A **file header** or **file descriptor** contains information about a file that is needed by the system programs that access the file records. The header includes information to determine the disk addresses of the file blocks as well as to record format descriptions, which may include field lengths and the order of fields within a record for fixed-length unspanned records and field type codes, separator characters, and record type codes for variable-length records.

## **Operations on Files**

Operations on files are usually grouped into **retrieval operations** and **update operations**.

The **retrieval operations** do not change any data in the file, but only locate certain records so that their field values can be examined and processed. The latter change the file by insertion or deletion of records or by modification of field values. In either case, we may have to **select** one or more records for retrieval, deletion, or modification based on a **selection condition** (or **filtering condition**), which specifies criteria that the desired record or records must satisfy.

Actual operations for locating and accessing file records vary from system to system. DBMS software programs, access records by using these commands.

**Open.** Prepares the file for reading or writing. Allocates appropriate buffers (typically at least two) to hold file blocks from disk, and retrieves the file header. Sets the file pointer to the beginning of the file.

**Reset.** Sets the file pointer of an open file to the beginning of the file.

**Find (or Locate).** Searches for the first record that satisfies a search condition. Transfers the block containing that record into a main memory buffer (if it is not already there). The file pointer points to the record in the buffer and it becomes the *current record*.

**Read (or Get).** Copies the current record from the buffer to a program variable in the user program. This command may also advance the current record pointer to the next record in the file, which may necessitate reading the next file block from disk.

**FindNext.** Searches for the next record in the file that satisfies the search condition. Transfers the block containing that record into a main memory buffer.

**Delete.** Deletes the current record and (eventually) updates the file on disk to reflect the deletion.

**Modify.** Modifies some field values for the current record and (eventually) updates the file on disk to reflect the modification.

**Insert:** Inserts a new record in the file by locating the block where the record is to be inserted; transferring that block into a main memory buffer (if it is not already there), writing the record into the buffer, and (eventually) writing the buffer to disk to reflect the insertion.

**Close.** Completes the file access by releasing the buffers and performing any other needed cleanup operations.

The preceding (except for Open and Close) are called **record-at-a-time** operations

because each operation applies to a single record. It is possible to streamline the operations Find, FindNext, and Read into a single operation, Scan, whose description is as follows:

**Scan.** If the file has just been opened or reset, Scan returns the first record; otherwise it returns the next record. If a condition is specified with the operation, the returned record is the first or next record satisfying the condition.

**FindAll.** Locates *all* the records in the file that satisfy a search condition.

**Find (or Locate) *n*.** Searches for the first record that satisfies a search condition and then continues to locate the next  $n - 1$  records satisfying the same condition. Transfers the blocks containing the *n* records to the main memory

buffer (if not already there).

**FindOrdered.** Retrieves all the records in the file in some specified order.

**Reorganize.** Starts the reorganization process. As we shall see, some file organizations require periodic reorganization. An example is to reorder the file records by sorting them on a specified field.

### **Files of Unordered Records (Heap Files)**

In this simplest and most basic type of organization, records are placed in the file in the order in which they are inserted, so new records are inserted at the end of the file. Such an organization is called a **heap** or **pile file**.

**Inserting a new record** is very efficient. The last disk block of the file is copied into a buffer, the new record is added, and the block is then **rewritten** back to disk. The address of the last file block is kept in the file header.

**Searching for a record** using any search condition involves a **linear search** through the file block by block—an expensive procedure. If only one record satisfies the search condition, then, on the average, a program will read into memory and search half the file blocks before it finds the record. For a file of  $b$  blocks, this requires searching  $(b/2)$  blocks, on average. If no records or several records satisfy the search condition, the program must read and search all  $b$  blocks in the file.

To **delete a record**, a program must first find its block, copy the block into a buffer, delete the record from the buffer, and finally **rewrite the block** back to the disk. This leaves unused space in the disk block. Deleting a large number of records in this way results in wasted storage space. Another technique used for record deletion is to have an extra byte or bit, called a **deletion marker**, stored with each record. A record is deleted by setting the deletion marker to a certain value.

## Files of Ordered Records (Sorted Files)

We can physically order the records of a file on disk based on the values of one of their fields—called the **ordering field**. This leads to an **ordered** or **sequential** file. If the ordering field is also a **key field** of the file—a field guaranteed to have a unique value in each record—then the field is called the **ordering key** for the file.

A **binary search** for disk files can be done on the blocks rather than on the records. Suppose that the file has  $b$  blocks numbered 1, 2, ...,  $b$ ; the records are ordered by ascending value of their ordering key field; and we are searching for a record whose ordering key field value is  $K$ . Assuming that disk addresses of the file blocks are available in the file header, the binary search can be described by Algorithm shown below.

Binary Search on an Ordering Key of a Disk File

$L \leftarrow 1$ ;  $U \leftarrow b$ ; (\*  $b$  is the number of file blocks \*)

while ( $U \geq L$ ) do

**begin**  $i \leftarrow (L + U) \text{ div } 2$ ;

        read block  $i$  of the file into the buffer;

        if  $K <$  (ordering key field value of the *first* record in block  $i$ )

            then  $U \leftarrow i - 1$

        else if  $K >$  (ordering key field value of the *last* record in block  $i$ )

            then  $L \leftarrow i + 1$

        else if the record with ordering key field value =  $K$  is in the buffer

            then goto found

        else goto notfound;

**end**;

    goto notfound;

A binary search usually accesses  $\log_2(b)$  blocks, whether the record is found or not—an improvement over linear searches, where, on the average,  $(b/2)$  blocks

are accessed when the record is found and  $b$  blocks are accessed when the record is not found.

**Inserting and deleting records** are expensive operations for an ordered file because the records must remain physically ordered. To insert a record, we must find its correct position in the file, based on its ordering field value, and then make space in the file to insert the record in that position. For a large file this can be very time consuming because, on the average, half the records of the file must be moved to make space for the new record. This means that half the file blocks must be read and rewritten after records are moved among them. For record deletion, the problem is less severe if deletion markers and periodic reorganization are used.

**Modifying** a field value of a record depends on two factors: the search condition to locate the record and the field to be modified. If the search condition involves the ordering key field, we can locate the record using a binary search; otherwise we must do a linear search. A nonordering field can be modified by changing the record and rewriting it in the same physical location on disk—assuming fixed-length records. Modifying the ordering field means that the record can change its position in the file. This requires deletion of the old record followed by insertion of the modified record.

## Hashing Techniques

Hashing is a method which provides very fast access to records under certain search conditions. This organization is usually called a **hash file**. The search condition must be an equality condition on a single field, called the **hash field**. In most cases, the hash field is also a key field of the file, in which case it is called the **hash key**. The idea behind hashing is to provide a function  $h$ , called a **hash function** or **randomizing function**, which is applied to the hash field value of a record and yields the **address** of the disk block in which the record is stored. A search for the record within the block can be carried out in a main memory buffer. For most records, we need only a single-block access to retrieve that record.

Hashing is also used as an internal search structure within a program whenever a group of records is accessed exclusively by using the value of one field.

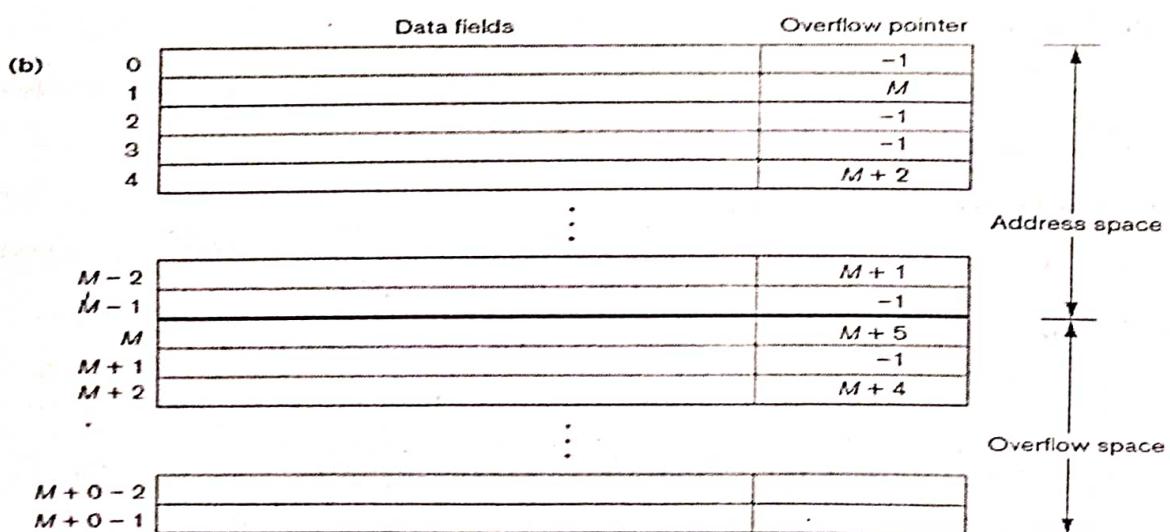
Hashing can be used for internal files. It can be modified to store external files on disk. It can also be extended to dynamically growing files.

### Internal Hashing

Hashing is typically implemented as a **hash table** through the use of an array of records. Let the array index range is from 0 to  $M-1$ , as shown in Figure(a); then we have  $M$  **slots** whose addresses correspond to the array indexes. We choose a hash function that transforms the hash field value into an integer between 0 and  $M-1$ . One common hash function is the  $h(K) = K \bmod M$  function, which returns the remainder of an integer hash field value  $K$  after division by  $M$ ; this value is then used for the record address.

(a)	Name	Ssn	Job	Salary
0				
1				
2				
3				
$M-2$				
$M-1$				

**Figure**  
Internal hashing data structures.  
(a) Array of  $M$  positions for use in internal hashing.  
(b) Collision resolution by chaining records.



\* null pointer = -1

\* overflow pointer refers to position of next record in linked list

Noninteger hash field values can be transformed into integers before the mod function is applied. For character strings, the numeric (ASCII) codes associated with characters can be used in the transformation—for example, by multiplying those code values. For a hash field whose data type is a string of 20 characters, Algorithm (a) can be used to calculate the hash address.

**Algorithm** Two simple hashing algorithms: (a) Applying the mod hash function to a character string K. (b) Collision resolution by open addressing.

(a)  $\text{temp} \leftarrow 1;$

```
for  $i \leftarrow 1$  to 20 do  $\text{temp} \leftarrow \text{temp} * \text{code}(K[i]) \text{ mod } M;$   
                 $\text{hash\_address} \leftarrow \text{temp} \text{ mod } M;$ 
```

(b)  $i \leftarrow \text{hash\_address}(K); a \leftarrow i;$

```
if location  $i$  is occupied  
    then begin  $i \leftarrow (i + 1) \text{ mod } M;$   
            while ( $i \neq a$ ) and location  $i$  is occupied  
                do  $i \leftarrow (i + 1) \text{ mod } M;$   
            if ( $i = a$ ) then all positions are full  
            else  $\text{new\_hash\_address} \leftarrow i;$   
            end;
```

Other hashing functions can be used. One technique, called **folding**, involves applying an arithmetic function such as *addition* or a logical function such as *exclusive or* to different portions of the hash field value to calculate the hash address.

A **collision** occurs when the hash field value of a record that is being inserted hashes to an address that already contains a different record. In this situation, we must insert the new record in some other position. The process of finding another position is called **collision resolution**.

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

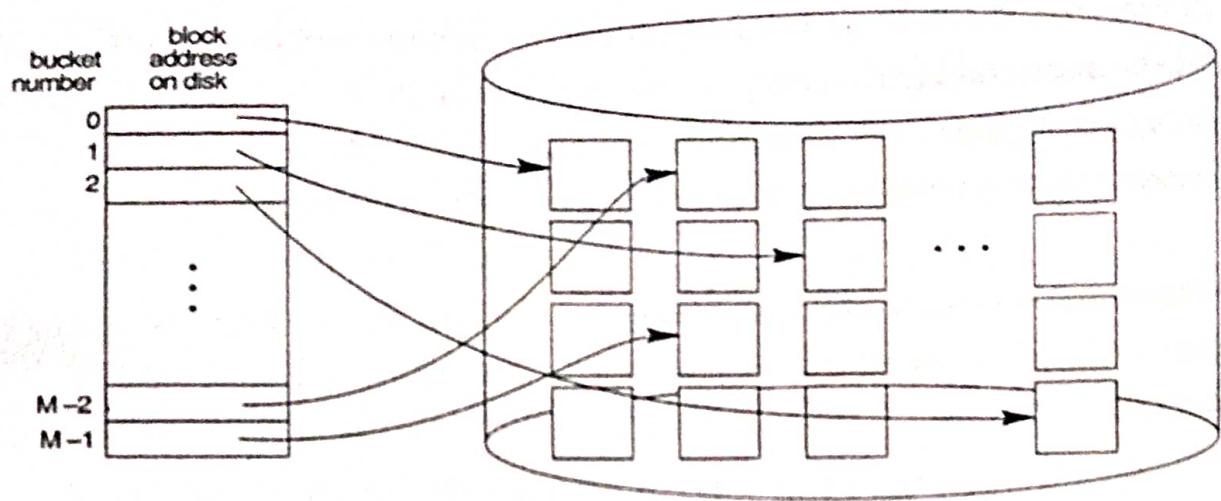
**Open addressing.** Once a position specified by the hash address is found to be occupied, the program checks the subsequent positions in order until an unused (empty) position is found. Algorithm (b) may be used for this purpose.

**Chaining.** For this method, various overflow locations are kept, usually by extending the array with a number of overflow positions. Additionally, 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. A linked list of overflow records for each hash address is thus maintained, as shown in Figure (b).

**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.

### **External Hashing for Disk Files**

Hashing for disk files is called **external hashing**. the target address space is in external hashing is made of **buckets**, A bucket is either one disk block or a cluster of contiguous blocks. The hashing function maps the indexing field's value into a relative bucket number. A table maintained in the file header converts the bucket number into the corresponding disk block address as shown in below Figure.



The hashing scheme is called **static hashing** if a fixed number of buckets is allocated.

A major drawback of static hashing is that the number of buckets must be chosen large enough that can handle large files. That is, it is difficult to expand or shrink the file dynamically.

### **Dynamic Hashing**

Two hashing techniques are

Extendible hashing and Liner hashing.

### **Extendible hashing**

Basic Idea:

- Maintain a directory of bucket addresses instead of just hashing to buckets directly. (indirection)
- The directory can grow, but its size is always a power of 2.
- At any time, the directory consists of  $d$  levels, and a directory of depth  $d$  has  $2^d$  bucket pointers.
  - However, not every directory entry (bucket pointer) has to point to a unique bucket. More than one directory entry can point to the same one.
  - Each bucket has a local depth  $d'$  that indicates how many of the  $d$  bits of

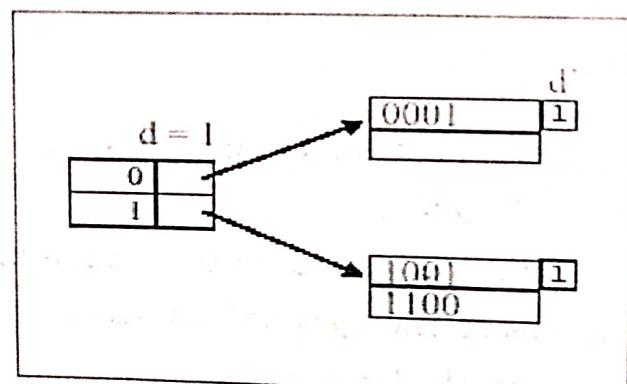
the hash value are actually used to indicate membership in the bucket.

- The depth  $d$  of the directory is based on the # of bits we use from each hash value.
  - The hash function produces an output integer which can be treated as a sequence of  $k$  bits. We use the first  $d$  bits in the hash value produced to look up an entry in the directory. The directory entry points us to the block that contains records whose keys hash to that value.

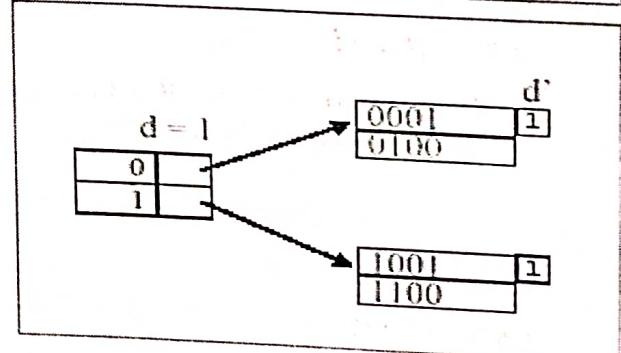
Example:

- Assume each hashed key is a sequence of four binary digits.

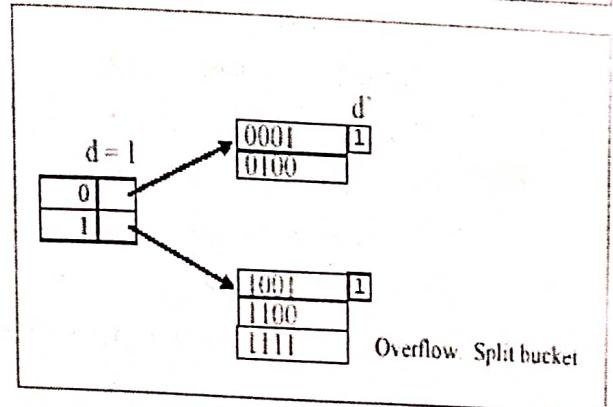
- Store values 0001, 1001, 1100.



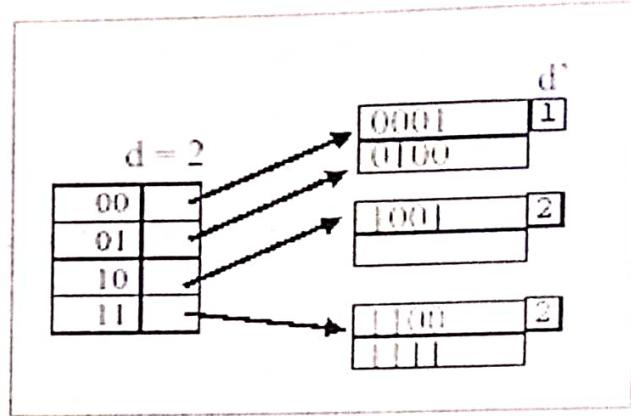
Insert 0100



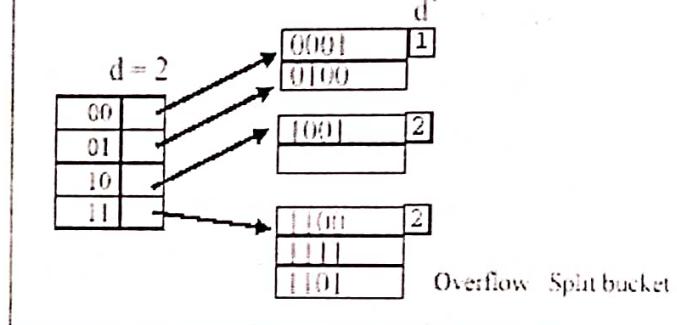
Insert 1111



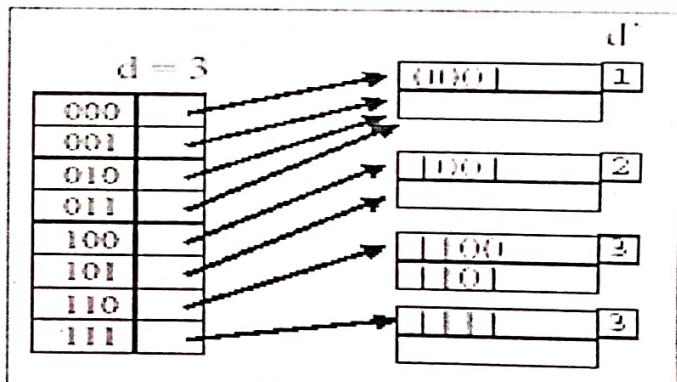
Directory grows one level.



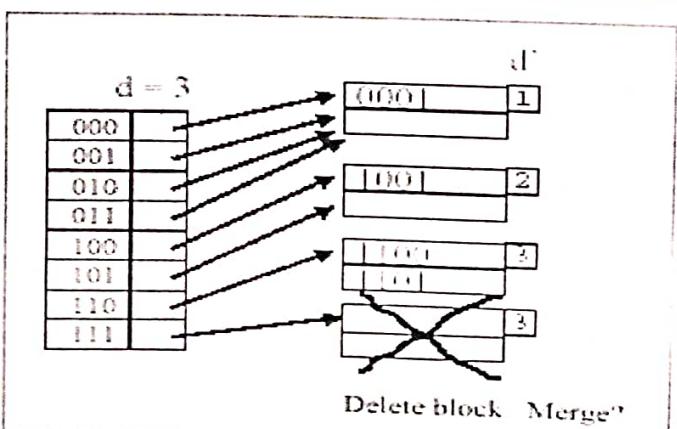
Insert 1101



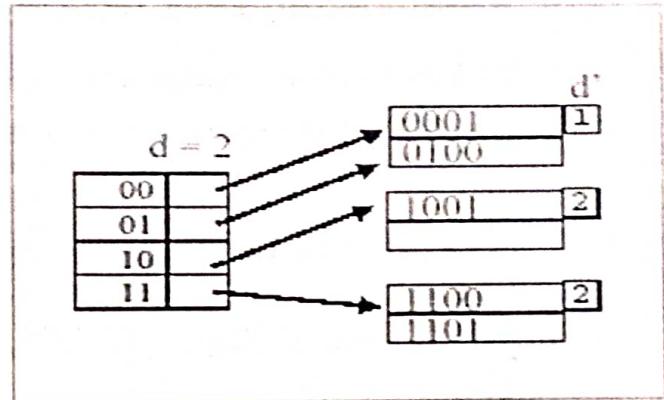
Directory grows one level.



Delete 1111



## Merge Blocks and Shrink Directory



## Linear hashing

- Linear hashing allows a hash file to expand and shrink dynamically without the need of a directory.
  - Thus, directory issues like extendible hashing are not present.
- A linear hash table starts with  $2d$  buckets where  $d$  is the # of bits used from the hash value to determine bucket membership.
  - The size of the table will grow gradually, but not double in size.

The growth of the hash table can either be triggered:

- 1) Every time there is a bucket overflow.
- 2) When the load factor of the hash table reaches a given point.

We will examine the second growth method.

Since overflows may not always trigger growth, note that each bucket may use overflow blocks.

## Parallelizing Disk Access Using RAID Technology

**RAID** (originally **redundant array of inexpensive disks**; now commonly **redundant array of independent disks**) is a data storage that combines multiple disk drive components into a logical unit for the purposes of data redundancy or performance improvement.

Disk striping is the process of dividing a body of data into blocks and spreading the data blocks across several partitions on several hard disks.

Striping can be done at the byte level, or in blocks.

**Byte-level striping** means that the file is broken into "byte-sized pieces". The first byte of the file is sent to the first drive, then the second to the second drive, and so on. Sometimes byte-level striping is done as a sector of 512 bytes.

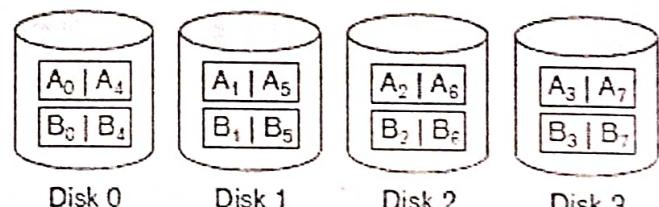
**Block-level striping** means that each file is split into blocks of a certain size and those are distributed to the various drives. The size of the blocks used is also called the stripe size (or block size, or several other names), and can be selected from a variety of choices when the array is set up.

**Figure**

Striping of data across multiple disks.  
(a) Bit-level striping across four disks.  
(b) Block-level striping across four disks.

A <sub>0</sub>   A <sub>1</sub>   A <sub>2</sub>   A <sub>3</sub>   A <sub>4</sub>   A <sub>5</sub>   A <sub>6</sub>   A <sub>7</sub>
B <sub>0</sub>   B <sub>1</sub>   B <sub>2</sub>   B <sub>3</sub>   B <sub>4</sub>   B <sub>5</sub>   B <sub>6</sub>   B <sub>7</sub>

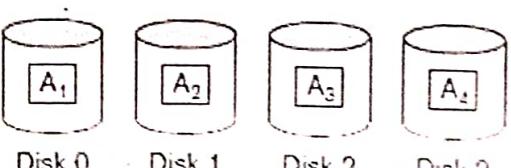
(a) Data



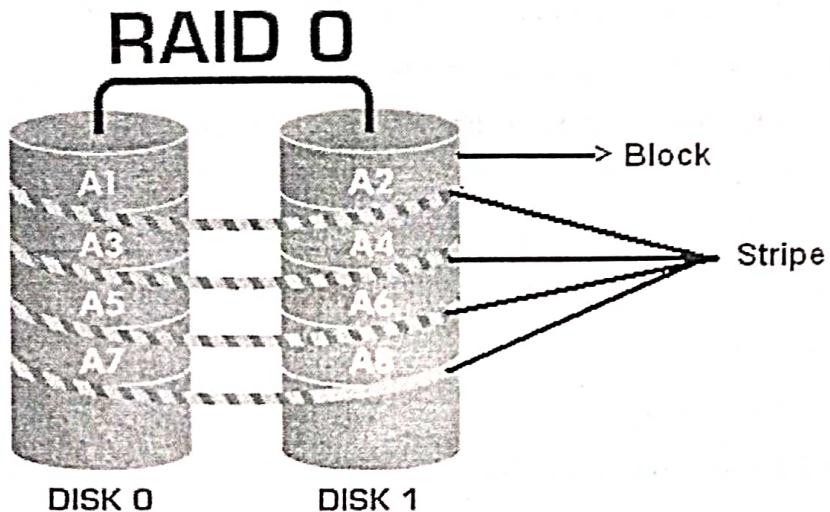
File A:

Block A <sub>1</sub>	Block A <sub>2</sub>	Block A <sub>3</sub>	Block A <sub>4</sub>
----------------------	----------------------	----------------------	----------------------

(b)



## RAID 0



RAID 0 performs what is called “Block Striping” across multiple drives. The data is fragmented, or broken up, into blocks and striped among the drives. This level increases the data transfer rate and data storage since the controller can access the hard disk drives simultaneously. However, this level has no redundancy. If single drive fails, the entire array becomes inaccessible. The more drivers in the array the higher the data transfer but higher risk of a failure. RAID level 0 requires 2 drives to implement.

### Advantages

- I/O performance is greatly improved by spreading the I/O load across many channels and drives
- Best performance is achieved when data is striped across multiple controllers with only one drive per controller
- No parity calculation overhead is involved
- Very simple design
- Easy to implement

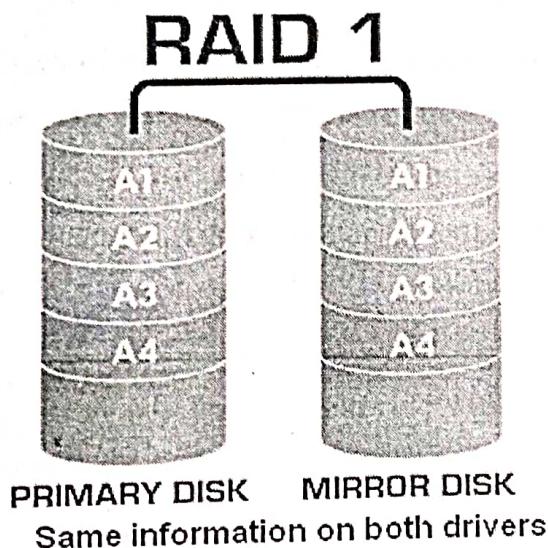
## **Disadvantages**

- Not a "True" RAID because it is NOT fault-tolerant
- The failure of just one drive will result in all data in an array being lost
- Should never be used in mission critical environments

## **Recommended Applications**

- Video Production and Editing
- Image Editing
- Pre-Press Applications
- Any application requiring high bandwidth

## **RAID 1**



**RAID 1** is used to create a complete mirrored set of disks. Usually employing two disks, RAID 1 writes data as a normal hard drive would but makes an exact copy of the primary drive with the second drive in the array. This mirrored copy is constantly updated as data on the primary drive changes. By keeping this mirrored backup, the array decreases the chance of failure from 5% over three years to 0.25%. Should a drive fail, the failed drive should be replaced as soon as possible.

to keep the mirror updated. RAID Level 1 requires a minimum of 2 drives to implement.

### **Advantages**

- Twice the Read transaction rate of single disks, same Write transaction rate as single disks
- 100% redundancy of data means no rebuild is necessary in case of a disk failure, just a copy to the replacement disk
- Transfer rate per block is equal to that of a single disk
- Under certain circumstances, RAID 1 can sustain multiple simultaneous drive failures
- Simplest RAID storage subsystem design

### **Disadvantages**

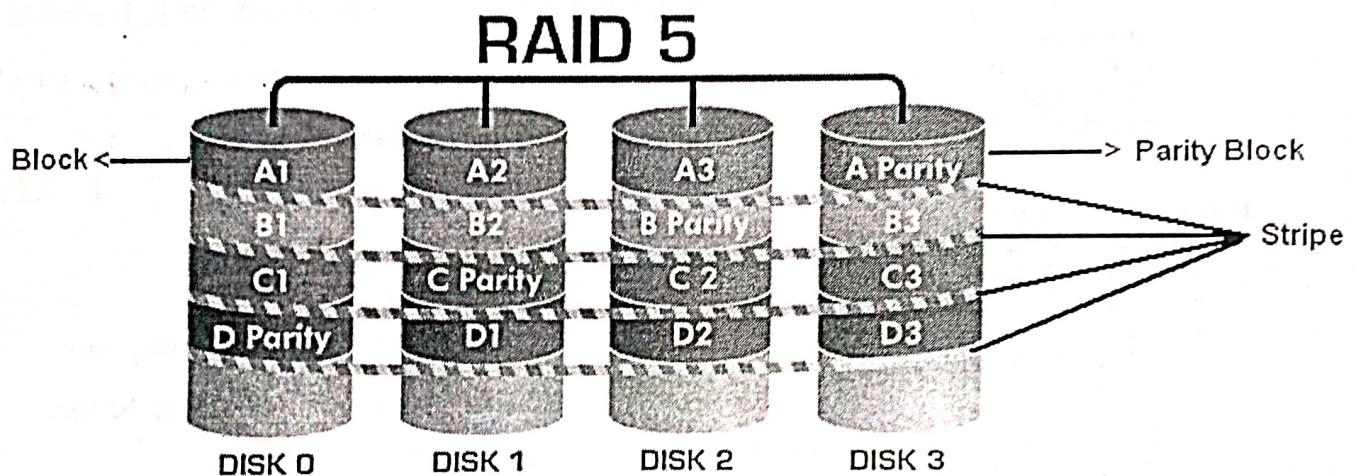
- Highest disk overhead of all RAID types (100%) - inefficient
- Typically the RAID function is done by system software, loading the CPU/Server and possibly degrading throughput at high activity levels. Hardware implementation is strongly recommended
- May not support hot swap of failed disk when implemented in "software"

### **Recommended Applications**

- Accounting
- Payroll
- Financial
- Any application requiring very high availability

## RAID 5

**RAID 5** is the most stable of the more advanced RAID levels and offers redundancy, speed, and the ability to rebuild a failed drive. RAID 5 uses the same block level striping as other RAID levels but adds another level of data protection by creating a "parity block". These blocks are stored alongside the other blocks in the array in a staggered pattern and are used to check that the data has been written correctly in the drive.



when an error or failure occur, the parity block will be used to locate the information stored on the other member disks and parity blocks in the array to rebuild the data correctly. This can be done on-the-fly without interruptions to applications or other programs running on the computer. The computer will notify the user of the failed drive, but will continue to operate normally. This state of operation is known as Interim Data Recovery Mode. Performance may suffer while the disk is being rebuilt, but operation should continue. The failed drive should be replaced as soon as possible. RAID Level 5 requires a minimum of 3 drives to implement.

### **Advantages**

- Highest Read data transaction rate
- Medium Write data transaction rate
- Low ratio of Error Correction Code (Parity) disks to data disks means high efficiency
- Good aggregate transfer rate

### **Disadvantages**

- Disk failure has a medium impact on throughput
- Most complex controller design
- Difficult to rebuild in the event of a disk failure (as compared to RAID 1)
- Individual block data transfer rate same as single disk

### **Recommended Applications**

- File and Application servers
- Database servers
- Web, E-mail, and News servers
- Intranet servers
- Most versatile RAID level

## 7. Indexing

### Contents:

- Single-Level Ordered Indexes
- Multi-Level Indexes
- B<sup>+</sup> Tree based Indexes
- Index Definition in SQL

### Basic Concepts

- Indexing mechanisms are used to optimize certain accesses to data (records) managed in files. For example, the author catalog in a library is a type of index.
- *Search Key (definition)*: attribute or combination of attributes used to look up records in a file.
- An *Index File* consists of records (called index entries) of the form

search key value	pointer to block in data file
------------------	-------------------------------

- Index files are typically much smaller than the original file because only the values for search key and pointer are stored.
- There are two basic types of indexes:
  - *Ordered indexes*: Search keys are stored in a sorted order (main focus here in class).
  - *Hash indexes*: Search keys are distributed uniformly across “buckets” using a hash function.

## Index Evaluation Criteria

Indexing techniques are evaluated on the basis of:

- Access types that are efficiently supported; for example,
  - search for records with specified values for an attribute  
**(select \* from EMP where EmpNo = 4711;)**
  - search for records with an attribute value in a specified range  
**(select \* from EMP where DeptNo between 20 and 50;)**
- Access time (index entry → record)
- Insertion time (record → index entry)
- Deletion time (record → index entry)
- Space and time overhead (for maintaining index)

## Types of Single Level Ordered Indexes

- In an *ordered index file*, index entries are stored sorted by the search key value.
  - most versatile kind of index: supports lookup by search key value or by range of search key values
- *Primary Index*: in a sequentially ordered file (e.g., for a relation), the index whose search key specifies the sequential order of the file. For a relation, there can be at most one primary index. ( $\leadsto$  index-sequential file)
- *Secondary Index*: an index whose search key is different from the sequential order of the file (i.e., records in the file are not ordered according to secondary index).
- If search key does not correspond to primary key (of a relation), then multiple records can have the same search key value
- *Dense Index Files*: index entry appears for every search key value in the record file.
- *Sparse Index Files*: only index entries for some search key values are recorded.
  - To locate a record with search key value  $K$ , first find index entry with largest search key value  $< K$ , then search file sequentially starting at the record the index entry points to
  - Less space and maintenance overhead for insertions and deletions
  - Generally slower than dense index for directly locating records

## Secondary Indexes

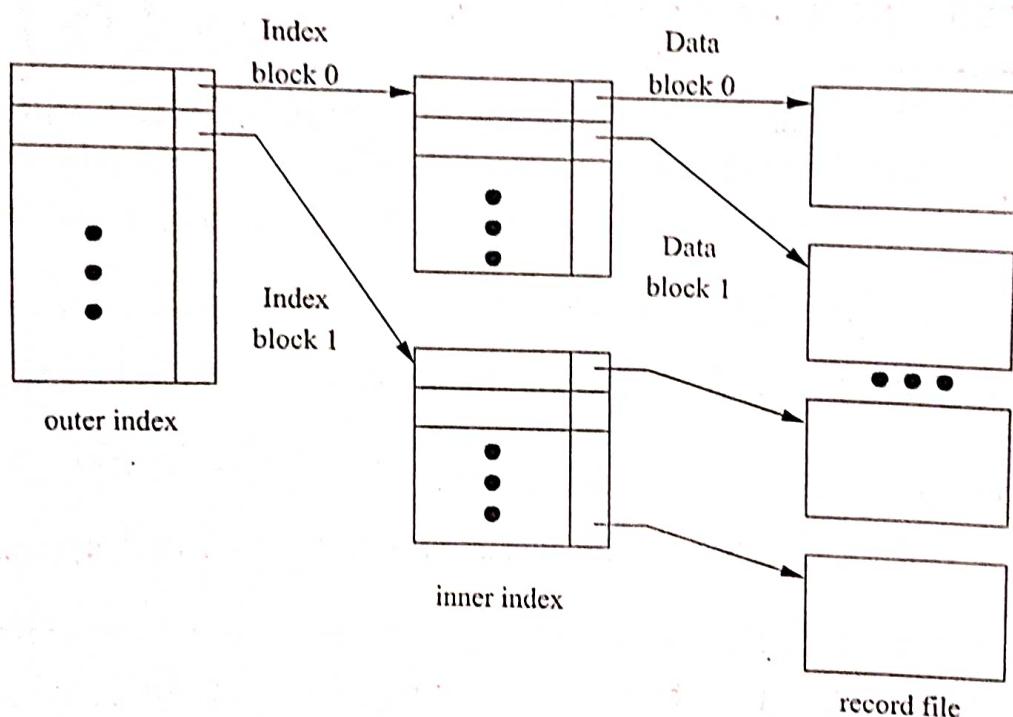
- Often one wants to find all records whose values in a certain field (which is not the search key of the primary index) satisfy some condition
  - Example 1: In the EMPLOYEE database, records are stored sequentially by EmpNo, we want to find employees working in a particular department.
  - Example 2: as above, but we want to find all employees with a specified salary or range of salary
- One can specify a secondary index with an index entry for each search key value; index entry points to a *bucket*, which contains pointers to all the actual records with that particular search key.

## Primary Indexes vs. Secondary Indexes

- Secondary indexes have to be dense
- Indexes offer substantial benefits when searching for records
- When a record file is modified (e.g., a relation), every index on that file must be updated. Updating indexes imposes overhead on database performance.
- Sequential scan using primary index is efficient, but a sequential scan using a secondary index is expensive (each record access may fetch a new block from disk)

## Multi-Level Index

- If primary index does not fit in memory, access to records becomes expensive
- To reduce number of disk accesses to index entries, treat primary index on disk as sequential file and construct a sparse index on it.
  - outer index → a sparse index of primary index
  - inner index → the primary index file
- Multilevel Index structure



- If even outer index is too large to fit in main memory, yet another level of index can be created, and so on.
- Note that indexes at all levels must be updated on insertions or deletions of records from a file.

## **Dynamic Multi-Level Indexes using B<sup>+</sup>-Trees**

B<sup>+</sup>-Tree indexes are an alternative to index sequential files.

- Disadvantage of index-sequential files: performance degrades as sequential file grows, because many overflow blocks are created. Periodic reorganization of entire file is required.
- Advantage of B<sup>+</sup>-Tree index file: automatically reorganizes itself with small, local changes in the case of insertions and deletions. Reorganization of entire file is not required to maintain performance.
- Disadvantage of B<sup>+</sup>-Trees: extra insertions and deletion overhead, space overhead.
- Advantages of B<sup>+</sup>-Trees outweigh disadvantages, and B<sup>+</sup>-Trees are used extensively in all DBMS.

A  $B^+$ -Tree is a rooted tree satisfying the following properties:

- All paths from the root to leaf have the same length ( $\implies$  a  $B^+$  tree is a balanced tree).
- Each node that is not the root or a leaf node has between  $\lceil n/2 \rceil$  and  $n$  children (where  $n$  is fixed for a particular tree).
- A leaf node has between  $\lceil (n - 1)/2 \rceil$  and  $n - 1$  values.
- Special case: if the root is not a leaf, it has at least 2 children. If the root is a leaf, it can have between 0 and  $n - 1$  values.
- Typical structure of a node:

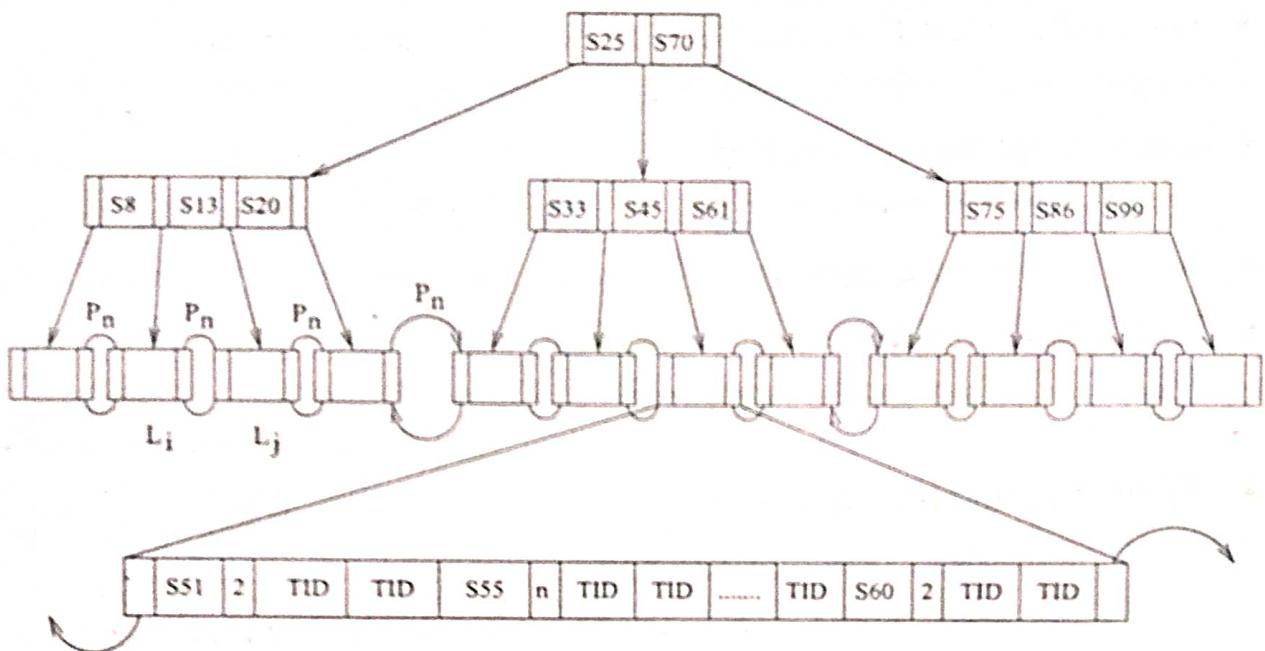
$P_1$	$K_1$	$P_2$	$\dots$	$P_{n-1}$	$K_{n-1}$	$P_n$
-------	-------	-------	---------	-----------	-----------	-------

- $K_i$  are the search key values
- $P_i$  are pointers to children (for non-leaf nodes) or pointers to records or buckets of records (for leaf nodes)

- The search keys in a node are ordered, i.e,

$$K_1 < K_2 < K_3 \dots < K_{n-1}$$

## Example of a $B^+$ -Tree



## Leaf Nodes in a $B^+$ -Tree

- For  $i = 1, 2, \dots, n - 1$ , pointer  $P_i$  either points to a file record with search key value  $K_i$  (using the tuple identifier, tid), or to a bucket of pointers to file records, each record having search key value  $K_i$ .  
Note that one only needs bucket structure if search key does not correspond to primary key of relation the index is associated with.
- If  $L_i, L_j$  are leaf nodes and  $i < j$ ,  $L_i$ 's search key values are less than  $L_j$ 's search key values.
- $P_n$  points to next leaf node in search key order.

## Non-Leaf Nodes in a $B^+$ -Tree

- All the search keys in the subtree to which  $P_1$  points are less than  $K_1$ ; all search keys in the subtree to which  $P_m$  points are greater than or equal to  $K_{m-1}$ .

## Observations about $B^+$ -Trees

- Since the inter-node connections are done by pointers, there is no assumption that in the  $B^+$ -Tree logically close blocks are also "physically" close.
- The non-leaf levels of the  $B^+$ -Tree form a hierarchy of sparse indices.
- The  $B^+$ -Tree contains a relatively small number of levels (logarithmic in size of the main file), thus searches can be done efficiently.
- Insertions and deletions to the main file can be handled efficiently, as the index can be restructured in logarithmic time. (☞ ECS 110)

## Queries on $B^+$ -Trees

Find all records with a search key value of  $k$

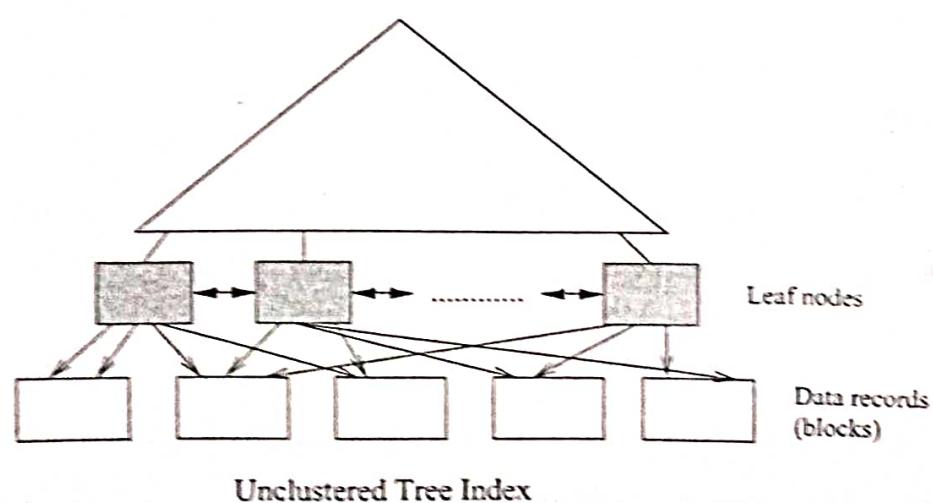
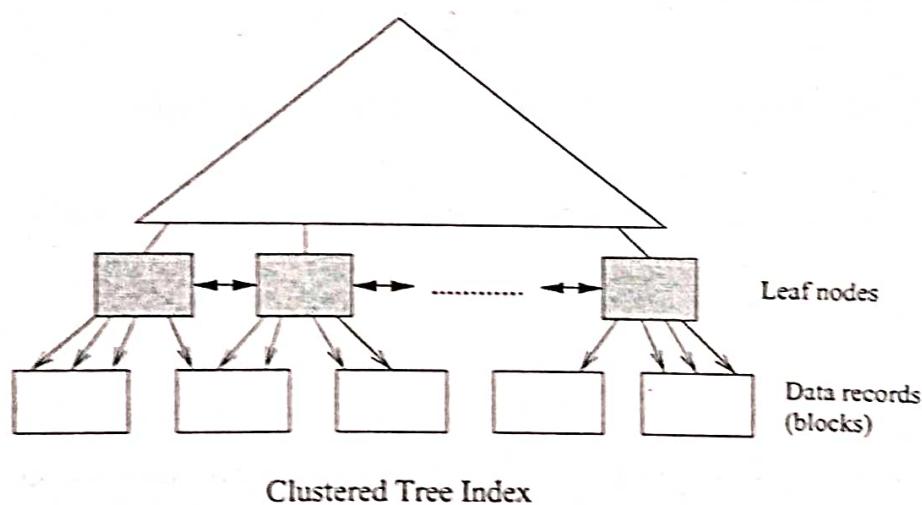
- Start with the root node
  - Examine the node for the smallest search key value  $> k$ .
  - If such a value exists, assume it is  $K_i$ . Then follow  $P_i$  to the child node.
  - Otherwise,  $k \geq K_{m-1}$ , where are  $m$  pointers in the node. Then follow  $P_m$  to the child node.
- If the node is reached by following the pointer above is not a leaf node, repeat the above procedure on the node, and follow the corresponding pointer.
- Eventually reach a leaf node. Scan entries  $K_i$  in the leaf node. If  $K_i = k$ , follow pointer  $P_i$  to the desired record or bucket. Otherwise no record with search key value  $k$  exists.
- Further comments:
  - If there are  $V$  search key values in the file, the path from the root to a leaf node is no longer than  $\lceil \log_{\lceil n/2 \rceil}(V) \rceil$ .
  - In general a node has the same size as a disk block, typically 4KB, and  $n \approx 100$  (40 bytes per index entry).
  - With 1,000,000 search key values and  $n = 100$ , at most  $\log_{50}(1,000,000) = 4$  nodes are accessed in the lookup!

## B<sup>+</sup>-Tree File Organization

- Index file degradation problem is solved by using B<sup>+</sup>-Tree indices. Data file degradation problem is solved by using a B<sup>+</sup>-Tree file organization.
- Leaf nodes in a B<sup>+</sup>-Tree file organization can store records instead of just pointers.

### Clustered vs. Unclustered Indices

Clustered: Order of data records is the same as order of index entries.



## **Index Definition in PostgreSQL**

- Indexes are not part of SQL standard, but nearly all DBMS's support them via a syntax like the one below.
- PostgreSQL syntax:

**create [unique] index <index name> on <relation name> (<list of attributes>);**

**drop index <index name>;**

- Many more options available, including clauses to specify sort order, partial indexes, fill factor, tablespace, concurrent construction, index method, . . .
- By default, indexes are created in ascending order.
- With primary key in a relation, an index is associated.

- Information about indexes is stored in the system catalogs.  
Relevant tables are pg\_index and pg\_class.

The system catalog table pg\_index:

Column	Description
indexrelid	The OID of the pg_class entry for this index
indrelid	The OID of the pg_class entry for the table this index is for
indnatts	The number of columns in the index (duplicates pg_class.relnatts)
indisunique	If true, this is a unique index
indisprimary	If true, this index represents the primary key of the table
...	

- Example:

```
create index city_name_idx on CITY(name);
```