

UNIT 3 FILES

Structure

- 3.0 **Introduction**
- 3.1 **Objectives**
- 3.2 **Terminology**
- 3.3 **File Organisation**
- 3.4 **Sequencial Files**
- 3.5 **Direct File Organisation**
- 3.6 **Indexed Sequencial File Organisation**
- 3.7 **Summary**
- 3.8 **Model Answer**

3.0 INTRODUCTION

Many tasks in an information oriented society require people to deal with voluminous data and to use computers to deal with this data efficiently and speedily. For example in an airlines reservation, office data regarding flights, routes, seats available etc. is required to facilitate booking of seats. A university might like to store data related to all students-the courses they sign up for etc. all this implies the following:

- Data will be stored on external storage devices like magnetic tapes, floppy disks etc.
- Data will be accessed by many people and software programs
- Users of the data will expect that
 - it is always reliably available for processing
 - it is secure
 - it is stored in a manner flexible enough to allow the users to add new types of data as per his changing needs

Files deal with storage and retrieval of data in a computer. Programming languages incorporate statements and structures which allow users to write programs to access and use the data in the file.

3.1 OBJECTIVES

After a study of this unit, you will be able to:

- * define the various terms related to files;
- * discuss the types of files

It will be useful for the student to recapitulate the units of Block 1 to refresh their knowledge of syntax related to file handling in PASCAL & C and also Unit 4 of Block 4 on Tree structures.

3.2 TERMINOLOGY

We Will now define the term's of the hierarchical structure of computer stored data collections.

1. **Field:** It is an elementary data item characterized by its size, length and type.

For example:

Name	:	a character type Of size 10
Age	:	a numeric type

2. **Record:** It is a collection of related fields that can be treated as a unit from an applications point of view.

For example:

A university could use a student record with the , fields, university enrolment no., Name Major subjects

3. **File:** Data is organized for storage in files. A file is a collection of similar, related records. it has an identifying name.

For example. "STUDENT" could be a file Consisting of student records for all the pupils in a university.

4. **Index:** An index file corresponds to a data file. It's records contain a key field and a Pointer to that record of the data file which has the same value of the key field.

Indexing will be discussed in detail later in the unit.

The data stored in files is accessed by software which can be divided into the following two categories:

1. **User Programs:** These are usually written by a Programmer to manipulate retrieved data in the manner required by the application.
2. **File Operations:** These deal with the physical movement of data in and out of files. User programs effectively use file operations through appropriate programming language syntax. The File Management system manages the independent files and acts as the software interface between the user programs and the file operations.

File operations can be categorized as

1. CREATION of the file
2. INSERTION of records in the file
3. UPDATION of previously inserted records
4. RETRIEVAL of previously inserted records
5. DELETION of records
6. DELETION of the file.

3.3 FILE ORGANISATION

File organization can most simply be defined as the method of sorting storing Data record in file in a file and the subsequent implications on the way these records can be accessed. The factors involved in selecting a particular file organization for use are:

- Ease of retrieval

- Convenience of updates
- Economy of storage
- Reliability
- Security
- Integrity

Different file organizations accord the above factors differing weightages. The choice must be made depending upon the individual needs of the particular application in question.

We now introduce in brief the various commonly encountered file organizations.

1. Sequential Files

Data records are stored in some specific sequence e.g. order of arrival value of key field etc. Records of a sequential file cannot be accessed at random i.e. to access the nth record, one must traverse the preceding (n-1) records. Sequential files will be dealt with at length in the next section.

2 Relative Files

Each data record has a fixed place in a relative file. Each record must have associated with it an integer key value that will help identify this slot. This key, therefore, will be used for insertion and retrieval of the record. Random as well as sequential access is possible. Relative files can exist only on random access devices like disks.

3. Direct Files

These are similar to relative files, except that the key value need not be an integer. The user can specify keys which make sense to his application.

4. Indexed Sequential Files

An index is added to the sequential file to provide random access. An overflow area needs to be maintained to permit insertion in sequence.

5. Indexed Files

In this file organization, no sequence is imposed on the storage of records in the data file, therefore, no overflow area is needed. The index, however, is maintained in strict sequence. Multiple indexes are allowed on a file to improve access.

3.4 SEQUENTIAL FILES

We will now discuss in detail the Sequential file organization as defined in 3.3. Sequential files have data records stored in a specific sequence.

A sequentially organized file may be stored on either a serial-access or a direct-access storage medium.

3.4.1 Structure

To provide the 'sequence' required a 'key' must be defined for the data records. Usually a field whose values can uniquely identify data records is selected as the key. If a single field cannot fulfil this criterion, then a combination of fields can serve as the key. For example in a file which keeps student records, a key could be student no.

3.4.2 Operations

1. **Insertion:** Records must be inserted at the place dictated by the sequence of the keys. As is obvious, direct insertions into the main data file would lead to frequent rebuilding of the file. This problem could be mitigated by reserving 'overflow areas' in the file for insertions. But this leads to wastage of space and also the overflow areas may also be filled.

The common method is to use transaction logging. This works as follows:

- collect records for insertion in a transaction file in their order of arrival
- when population of the transactions file has ceased, sort the transaction file in the order of the key of the primary data file
- merge the two files on the basis of the key to get a new copy of the primary sequential file.

Such insertions are usually done in a batch mode when the activity/program which populates the transaction file have ceased. The structure of the transaction files records will be identical to that of the primary file.

2. **Deletion:** Deletion is the reverse process of insertion. The space occupied by the record should be freed for use. Usually deletion (like-insertion) is not done immediately. The concerned record (al 9 with a marker or 'tombstone' to indicate deletion) is written to a transaction file. At the time of merging the corresponding data record will be dropped from the primary data file.
3. **Updation :** Updation is a combination of insertion and deletions. The record with the new values is inserted and the earlier version deleted. This is also done using transaction files.
4. **Retrieval:** User programs will often retrieve data for viewing prior to making decisions, therefore, it is vital that this data reflects the latest state of the data if the merging activity has not yet taken Place.

Retrieval is usually done for a particular value of the key field. Before return in to the user, the data record should be merged with the transaction record (if any) for that key value.

The other two operations 'creation' and 'deletion' of files are achieved by simple programming language statements.

3.4.3 Disadvantages

Following are some of the disadvantages of sequential file organization:

- * Updates are not easily accommodated
- * By definition, random access is not possible
- * All records must be structurally identical. If a new field has to be added, then every record must be rewritten to provide space for the new field.

- * Continuous areas may not be possible because both the primary data file and the transaction file must be looked during merging.

3.4.4 Areas of Use

Sequential files are most frequently used in commercial batch oriented data processing where there is the concept of a master file to which details are added periodically. Ex. payroll applications.

Check Your Progress 1

1. Describe the record structure to be used for the lending section of a library.
2. Write a program in Pascal to insert the following records into a file 'PERSONAL'.

- Adam Bede, 47, , Engineer
- Silas Marner, 50, Doctor

Use a name field of size 20, age field of size 2 and profession field of size 20.

3. Merge the following, sequenced on NO:

Transactions		Master	
No.	Name	No:	Name
6	Beta	1	Delta
4	Alpha	2	Lamba
7	Gamma	8	Phi

3.5 DIRECT FILE ORGANIZATION

It offers an effective way to organize data when there is a need to access individual records directly.

To access a record directly (or random access) a relationship is used to translate the key value into a physical address. This is called the mapping function R R.(key value) -- Address

Direct files are stored on DASD (Direct Access Storage Device)

A calculation is performed on the key value to get an address. This address calculation technique is often termed as hashing. The calculation applied is called a hash function.

Here we discuss a very commonly used hash function called Division - Remainder.

Division-Remainder Hashing:

According to this method, key value is divided by an appropriate number, generally a prime number, and the division of remainder is used as the address for the record.

The choice of appropriate divisor may not be so simple. If it is known that the file is to contain n records, then we must, assuming that only one record can be stored at a given address, have divisor n.

Also we may have a very large key space as compared to the address space. Key space refers to all the possible key values. Although only a part of this space will be used as key values in the file, the size of key space, therefore a one to one mapping may not be there. That is calculated address may not be unique. It is called Collision, i.e.

$$R(K1) = R(K2) \text{ but } K1 \neq K2$$

Two unequal keys have been calculated to have the same address. The keys are called synonyms.

There are various approaches to handle the problem of collisions. One of these is to hash to buckets. A bucket is a space that can accommodate multiple records. A discussion on buckets and other such methods to handle collisions is out of the scope of this Unit. However the student is advised to read some text on Bucket Addressing and related topics.

3.6 INDEXED SEQUENTIAL FILE ORGANIZATION

When there is need to access records sequentially by some key value and also to access records directly by the same key value, the collection of records may be organized in an effective manner called Indexes Sequential Organization.

You must be familiar with search process for a word in a language dictionary. The data in the dictionary is stored in sequential manner. However an index is provided in terms of thumb tabs. To search for a word we do not search sequentially. We access the index that is the appropriate thumb tab, locate an approximate location for the word and then proceed to find the word sequentially.

To implement the concept of indexed sequential file organizations, we consider an approach in which the index part and data part reside on a separate file. The index file has a tree structure and data file has a sequential structure. Since the data file is sequenced, it is not necessary for the index to have an entry for each record. Following figure shows a sequential file with a two-level index.

Level 1 of the index holds an entry for each three-record section of the main file. The level 2 indexes level 1 in the same way.

When new records are inserted in the data file, the sequence of records need to be preserved and also the index is accordingly updated.

Two approaches used to implement indexes are static indexes and dynamic indexes.

As the main data file changes due to insertions and deletions, the static index contents may change but the structure does not change. In case of dynamic indexing approach, insertions and deletions in the main data file may lead to changes in the index structure. Recall the change in height of B-Tree as records are inserted and deleted.

Both dynamic and static indexing techniques are useful depending on the type of application.

3.7 SUMMARY

This Unit dealt with the methods of physically storing data in the files. The terms fields, records and files were defined. The organization types were introduced.

The various file organization were discussed. Sequential File Organization finds in use in application areas where batch processing is more common. Sequential Files are simple to use and can be stored on inexpensive media. They are suitable for applications that require direct access to only particular records of the collection. They do not provide adequate support for interactive applications.

In Direct file organization there exists a predictable relationship between the key used and by program to identify a Particular record and or programmer that record's location on secondary storage. A direct file must be stored on a direct access device. Direct files are used extensively in application areas where interactive processing is used.

An Indexed Sequential file supports both sequential access by key value and direct access to a particular record given its key value. It is implemented by building an index on top of a sequential data file that resides on a direct access storage device.

3.8 MODEL ANSWERS

Check Your Progress

1. The following record structure could take care of the general requirements of a lending library.

Member No., Member Name, Book Classification, i.e. Book Name, Author, Issue Date, Due Date.
2. {Personal, Per} program Insert (Personal);

```
type
    person =several
        name [20]: char;
        age : integer
        profession [20]: char
    end,
    filep = file of person;
    var
        personal: filep;
        persons [2] : person;
        x :integer
begin
    persons [1].name = "Adam Bede";
    persons [1]. age = 47
    persons [1]. profession = "Engineer";
    persons [2]. name = "Silai Mainu"
    persons [2]. age - 50
    Persons [2]. profession = 'Doctor"
    Rewrite (personal);
    For X:= 1 to 2 do
        write (Personal, persons [x]);
    End.
```
3.
(1) Sort Transaction file

No.	Name
4	Alpha
6	Beta
7	Gamma

(2) Merge to get

No.	Name
1	Delta
2	Lambda
4	Alpha
6	Beta
7	Gamma
8	Phi

If a sequential file on a disc is to occupy the least possible space its records must be stored contiguously i.e. with no unused space between them.

In case of addition or deletion of a record, the file must be rewritten to maintain its sequential order without spaces