# Experiment-04

1) Spark frame work is from apache vendor which is open source
2) Spark is only meant for big data processing (using In-memory cluster computing primitive) --- however it's not meant for storage.
3) Spark does not have any storage component --- and it can read the input data from local file system (LFS), HDFS, RDBMS, NoSQL, cloud vendors (Azure/AwS).

→ one of the processing is done from spark, it can write the o/p data to any of the above storage systems.
→ When we install spark we will see (spark core, sparksol, spark streaming, sparkmllib, spark Graph-x)
   5 modules
→ Spark is a unified stack
→ spark processing can be achieved through
1. Scala language
2. Python cpyspark
3. R language
4. Java language

⎫
⎬  four compatable
⎭

## Advantage:-

The main advantage of spark is because of its In-memory cluster computing a processing which is not present in our hadoop hadoop is a desk based-processing system.

→ Hadoop follows <u>top to bottom</u> approach where as spark follows <u>bottom to top</u> approach.
→ If we want to see the different methods in spark, we need to type    scala> sc. (And press tab)
   It will show the different methods in spark,
→ In spark we have 2 different RDD operations like
   ① Transformation   ② Action

1

1) Transformation is a method of an RDD creates a new RDD by performing a computation on the source RDD
   1) Map                  5) Union
   2) Filter                6) Intersection
   3) flat map            7) substract
   4) Map partitions       8) Sort by

2) Actions are RDD methods that return a value to a driver.
① collect - method returns elements in source RDD as an array.
② count - The method returns a count of elements in source RDD.
→ To run spark processing in a standalone cluster manage node
   1) Using scala → Spark-shell
   2) Using python → Pyspark
   3) Using R → SparkR
   4) Using java → No support

I scala> var file =sc.textFile("file:/// home/vboxuser/prac/spark/
                             Input.log)

scala> var mapfile - file.map(x=> (x,x.length))
scala > mapfile. collect. foreach (println)
output:- (Apache spark is an open source, ,32)
(It is wide range data processing engine, 38)
(There is a common belief that Apache spark ,42)
(Is an extension of Hadoop ,28)
(It uses Hadoop for storage purpose only., 41)

II scala> Var data =sc. paralleliz (List ("spark", "scala", "c",
    "cobol", "vc++", "Hadoop", "Ruby", "information", "blockchain",
   " generative AI", "chat Gpt", "java", "python", "c++"))
scala > var mapdata = data.map (x => (x,x.Length))
scala > Mapdata. collect. foreach (println)

output:- (spark,5)
          (scala,5)
          (c,1)
          (cobol,5)
          (vc++,4)
          (hadoop,6)
          (ruby,4)
          (information,11)
          (blockchain,10)
          (generativeAI,12)
          (chatGPT,7)
          (java,4)
          (python,6)
          (c++,3)

→ Here collect gathers all the elements of RDD or dataframe.
→ for each is an action that implies a function to each element of collected data.
→ println is simply printing out each element
→ It collects all elements from mapdata & then prints each tuple to the console.

scala > var Mapdata = data.map (x => (x, x.reverse, x.length))
scala > mapdata.collect . foreach (println)
output:- (spark, kraps, 5)
          (scala, alacs, 5)
          (c,c, 1)
          (cobol, loboc, 5)
          (vc++, ++cv., 4)
          (hadoop, poodah, 6)
          (Ruby, ybuR, 4)
          (information, noitamrofni, 11)
          (blockchain, niahckcolb, 10)
          (generativeAI, IA.evitareneg, 12)
          (chatGPT, TPGtahc, 7)
          (java, avaj, 4)
          (python, nohtyp, 6)     (c++, ++c, 3)

filter transformation:- filter is used like where clause type

```scala
scala> var file = sc. textfile ("file: ///home/vboxuser/prac/spark/Test.log")
scala> var filfile = file. filter (x=>x.length >=10)
scala > file. cocent
O/P:- res7: hong = 19
scala > filfile. coount
      res8: hong = 17
scala> filfile. collect. foreach (println)
```

output:- Spark SQL is Apache Spark's module
         for working with structured data.
         The SQL syntax section describes
         the SQL syntax in detail along with usage
         examples when applicable.

Spark SQL>

     To work with structured & semi-structured data by
making use of tabular kind of view, we are going ahead
with spark SQL.

⇒ What is main entry point in spark core =>sc

=) What is main programming abstraction in spark core is __
   --> RDD (Resilient Distributed database)

=) What is main entry point in Spark SQL → spark

=) What is main programming abstraction in spark SQL-->
Dataframe (DF)

=) What is a Dataframe (DF)

Dataframe is a row-column object which is a schema based RDD
schema + RDD => Its a kind a table in RDBMS => Dataframe
    in Spark.

we create an RDD by making use of sc.textfile sc. parallelize

how to convert **RDD to Data frame.**

⇒ from the created RDD

syntax:- rdd - object . toDF

↓
Default method available & which is automatically

convert rdd to dataframe.

Eg:- scala > var data = sc. parallelize (1 to 40)
scala > var mapdata = data.map (x => (x, (x*x)))
scala > var dfobj = mapdata. to DF
scala > dfobj. columns
scala > dfobj. collect.
output:- res12: Array [org. apache. spark. sql. row] = Array ( [1,1],
[2,4] , [3,9], [4,16], [5,25], [6,36], [7,49], [8,64], [9,81], [10,100],
- - - - - [40,1600])

scala > dfobj . show()
output:-

```
+----+----+
| -1 | -2 |
+----+----+
| 1  | 1  |
| 2  | 4  |
| 3  | 9  |
| 4  | 16 |
| 5  | 25 |
| 6  | 36 |
| 7  | 49 |
| 8  | 64 |
| 9  | 81 |
| 10 | 100|
| 11 | 121|
| 12 | 144|
| 13 | 144|
| 14 | 109|
| 15 | 196|
| 16 | 225|
| 17 | 256|
| 18 | 289|
| 19 | 324|
| 20 | 361|
| 20 | 400|
+----+----+
```

only showing top 20 rows

scala > dfobj .show(40)

output:-

| -1 | -2 |
|---|---|
| 1 | 1 |
| 2 | 4 |
| 3 | 9 |
| 4 | 16 |
| 5 | 25 |
| 6 | 36 |
| 7 | 49 |
| 8 | 64 |
| 9 | 81 |
| 10 | 100 |
| 11 | 121 |
| 12 | 144 |
| 13 | 169 |
| 14 | 196 |
| 15 | 225 |
| 16 | 256 |
| 17 | 289 |
| 18 | 324 |
| 19 | 361 |
| 20 | 400 |
| 21 | 441 |
| 22 | 484 |
| 23 | 529 |
| 24 | 576 |
| 25 | 625 |
| 26 | 676 |
| 27 | 729 |
| 28 | 784 |
| 29 | 841 |
| 30 | 900 |
| 31 | 961 |
| 32 | 1024 |
| 33 | 1089 |
| 34 | 1156 |
| 35 | 1225 |
| 36 | 1296 |
| 37 | 1369 |
| 38 | 1444 |
| 39 | 1521 |
| 40 | 1600 |

→ If we use only show it will display only 20 records in given table, if we having more than 20 records, we should mention the number in [show()] → it is an action.

Eg: The number should be above 10 & the square should be less than 1500

```
scala > spark.sql ("show tables"). show
scala > dfobj.registerTempTable ("numbertable")
scala > spark.sql ("show tables"). show
```

output:

(without table name)

| namespace | tableName | isTemporary |
|-----------|-----------|-------------|

(with table name)

| namespace | tablename | isTemporary |
|-----------|-------------|-------------|
|           | numbertable | true        |

```
scala > spark.sql (" select * from numbertable where _1 > 10
                            AND    _2 < 1500 "). show
```

output:—

| _1 | _2 |
|----|-----|
| 11 | 121 |
| 12 | 144 |
| 13 | 169 |
| 14 | 196 |
| 15 | 225 |
| 16 | 256 |
| 17 | 289 |
| 18 | 324 |
| 19 | 361 |
| 20 | 400 |
| 21 | 441 |
| 22 | 484 |
| 23 | 529 |
| 24 | 576 |
| 25 | 625 |
| 26 | 676 |
| 27 | 729 |
| 28 | 784 |
| 29 | 841 |
| 30 | 900 |

```
scala> var data = sc. parallelize (List (("spark",50), ("c",40),
("Javascript",80), ("c++", 30), ("cobol", 10), ("vc++", 90),
("python",40), ("Groovy",30), ("DS",50), ("ratig",90), ("Blockchain",60),
("AI", 40)))
scala > var mapval data = data. mapvalues(x => (x+100))
scala > var mapvaldata. toDF ("technology", "rating")

scala > var df = mapvaldata.
scala > var df. registerTempTable ("techtable")
scala > spark. sql ("select * from techtable where rating >160")- show
```

| technology   | rating |
|--------------|--------|
| java script  | 180    |
| vc++         | 190    |
|              | 198    |
| blockchain   | 188    |