Big Data And Hadoop

Session 13 – Assignment 3

**Problem Statement:**

Write the code to Turn a collection into a RDD and perform map operation on it to cube every number and filter the number which are divided by two and three.

**Solution:**

We first start spark context using the command **spark-shell.** The scala shell appears as follows:



Step 1:
Let us consider the collection to be a List with the following elements: (1, 2, 3, 4, 5, 6, 7, 8, 9, 10).
We create an RDD representing this data using parallelize(). Parallelized collections are created by calling SparkContext's **parallelize** method on an existing collection in your driver program (a Scala Seq). The elements of the collection are copied to form a distributed dataset that can be operated on in parallel. In the our example, we use parallelize as follows:

**val data = sc.parallelize( List(1, 2, 3, 4, 5, 6, 7, 8, 9, 10) ).**

This is a transformation performed .



On performing **collect** action, we get the output as shown in the above snapshot. Here, the RDD got created with the specified data.

**Collect** being an action, its execution appears in the Spark UI as a job as follows:

This action was split into two tasks as mentioned in the above snapshot. The DAG representation of this job is as follows:



This action had only one node in the DAG. Detailed DAG is as follows:

Step2:

Now, we use **map** transformation, wherein we will generate a cube of every element of the RDD as follows:

**val cubes =data.map( x=> x*x*x )**

This can be seen in the below snapshot as follows:

```
                                    acadgild@localhost:~                              _ □ □
File  Edit  View  Search  Terminal  Help

Using Scala version 2.10.5 (Java HotSpot(TM) 64-Bit Server VM, Java 1.8.0_65)
Type in expressions to have them evaluated.
Type :help for more information.
17/09/12 22:38:10 WARN Utils: Your hostname, localhost.localdomain resolves to a loopback address: 127.0.0.1; using 192.168.
6.101 instead (on interface eth1)
17/09/12 22:38:10 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address
17/09/12 22:38:10 WARN SparkConf:
SPARK_WORKER_INSTANCES was detected (set to '2').
This is deprecated in Spark 1.0+.

Please instead use:
 - ./spark-submit with --num-executors to specify the number of executors
 - Or set SPARK_EXECUTOR_INSTANCES
 - spark.executor.instances to configure the number of instances in the spark config.

Spark context available as sc.
17/09/12 22:38:15 WARN Connection: BoneCP specified but not present in CLASSPATH (or one of dependencies)
17/09/12 22:38:16 WARN Connection: BoneCP specified but not present in CLASSPATH (or one of dependencies)
17/09/12 22:38:22 WARN ObjectStore: Version information not found in metastore. hive.metastore.schema.verification is not en
bled so recording the schema version 1.2.0
17/09/12 22:38:22 WARN ObjectStore: Failed to get database default, returning NoSuchObjectException
17/09/12 22:38:25 WARN Connection: BoneCP specified but not present in CLASSPATH (or one of dependencies)
17/09/12 22:38:26 WARN Connection: BoneCP specified but not present in CLASSPATH (or one of dependencies)
SQL context available as sqlContext.

scala> val data = sc.parallelize(List(1,2,3,4,5,6,7,8,9,10));
data: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[0] at parallelize at <console>:27

scala> data.collect
res0: Array[Int] = Array(1, 2, 3, 4, 5, 6, 7, 8, 9, 10)

scala> val cubes = data.map(x=>x*x*x)
cubes: org.apache.spark.rdd.RDD[Int] = MapPartitionsRDD[1] at map at <console>:29

scala> cubes.collect
res1: Array[Int] = Array(1, 8, 27, 64, 125, 216, 343, 512, 729, 1000)

scala>
```
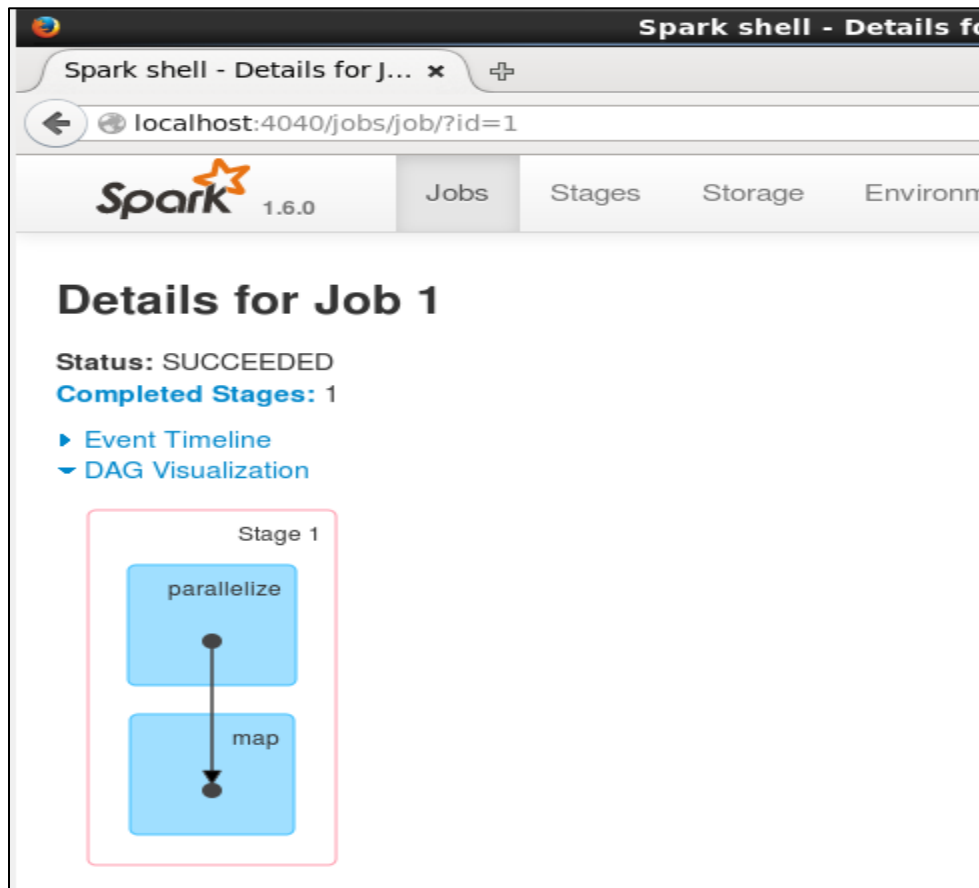
The above snapshot also shows the output on performing the **collect** action.
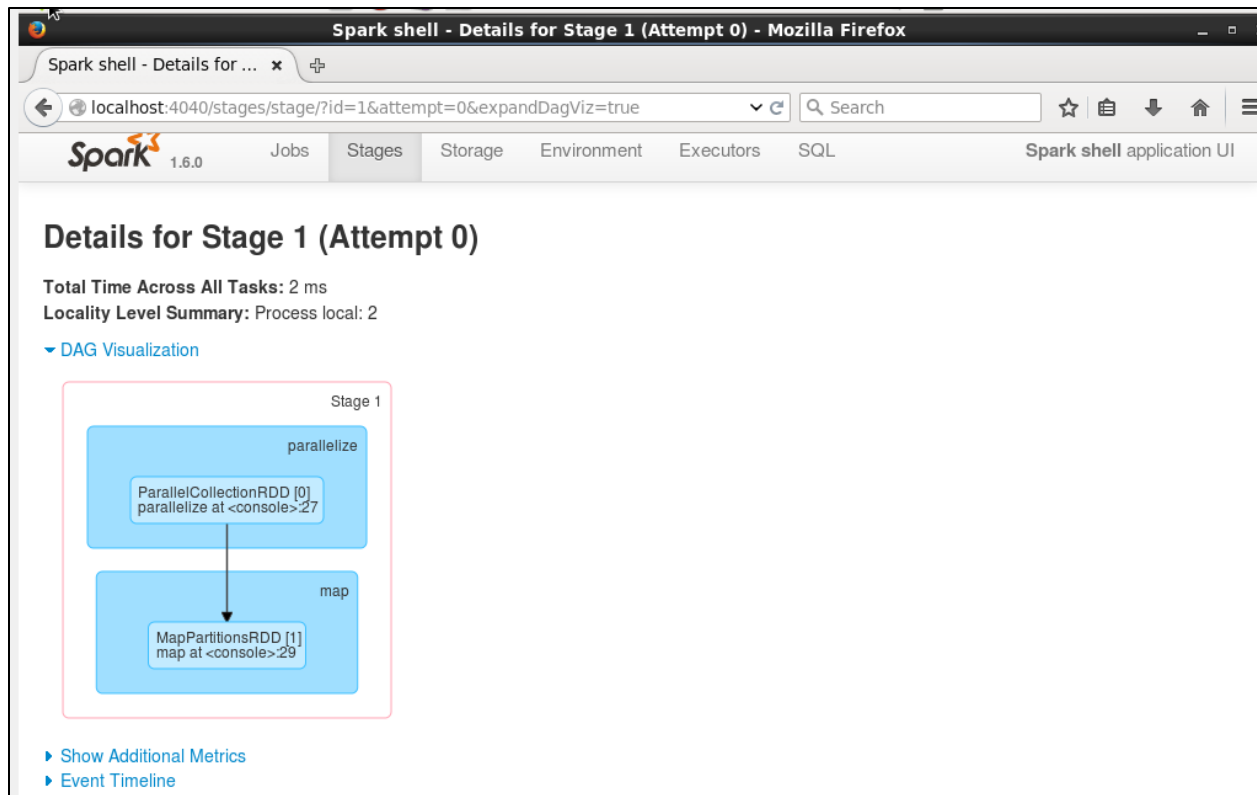
The execution of this action appeared in UI as follows.

Its DAG representation is as follows:



Detailed DAG representation is as follows:

Step3:

Next is to filter the RDD. In this, we need only those elements of RDD that are divisible by 2 and 3. Hence, made use of **modulo** as follows:

**val filteredData = cubes.filter( x=> ( x%2==0 && x%3==0 ) )**

Its execution is shown below as follows:

```
                                    acadgild@localhost:~                              _  □

 File  Edit  View  Search  Terminal  Help

17/09/12 22:38:10 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address
17/09/12 22:38:10 WARN SparkConf:
SPARK_WORKER_INSTANCES was detected (set to '2').
This is deprecated in Spark 1.0+.

Please instead use:
 - ./spark-submit with --num-executors to specify the number of executors
 - Or set SPARK_EXECUTOR_INSTANCES
 - spark.executor.instances to configure the number of instances in the spark config.

Spark context available as sc.
17/09/12 22:38:15 WARN Connection: BoneCP specified but not present in CLASSPATH (or one of dependencies)
17/09/12 22:38:16 WARN Connection: BoneCP specified but not present in CLASSPATH (or one of dependencies)
17/09/12 22:38:22 WARN ObjectStore: Version information not found in metastore. hive.metastore.schema.verification is not ena
bled so recording the schema version 1.2.0
17/09/12 22:38:22 WARN ObjectStore: Failed to get database default, returning NoSuchObjectException
17/09/12 22:38:25 WARN Connection: BoneCP specified but not present in CLASSPATH (or one of dependencies)
17/09/12 22:38:26 WARN Connection: BoneCP specified but not present in CLASSPATH (or one of dependencies)
SQL context available as sqlContext.

scala> val data = sc.parallelize(List(1,2,3,4,5,6,7,8,9,10));
data: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[0] at parallelize at <console>:27

scala> data.collect
res0: Array[Int] = Array(1, 2, 3, 4, 5, 6, 7, 8, 9, 10)

scala> val cubes = data.map(x=>x*x*x)
cubes: org.apache.spark.rdd.RDD[Int] = MapPartitionsRDD[1] at map at <console>:29

scala> cubes.collect
res1: Array[Int] = Array(1, 8, 27, 64, 125, 216, 343, 512, 729, 1000)

scala> val filteredData = cubes.filter(x=>(x%2==0  && x%3==0))
filteredData: org.apache.spark.rdd.RDD[Int] = MapPartitionsRDD[2] at filter at <console>:31

scala> filteredData.collect
res2: Array[Int] = Array(216)

scala>
```

Performing the action **collect** on filtered RDD gives the results as shown above.
The execution of this action appeared in UI as follows.

The DAG representation is as follows:



Detailed representation of DAG is as follows:

This is how we have performed required steps to solve the above mentioned problem statement.