ASSIGNMENT

NAME-TUSHAR RAJ VERMA

Q1- Setup Spark stand alone cluster with 1 master and 2 worker nodes (refer the github link provided) and demonstrate running an application (using spark-shell) on this stand alone cluster.

This application should run minimum 4 tasks in 1 stage in 1 job. It will be good to see if this application tasks utilize both worker nodes. Display this output via spark UI.

Ans-1
Spark stand alone cluster is set up with spark running

```
Terminal

File Edit View Search Terminal Help

tushar@tushar:spark $ jps
5298 Worker
5080 Master
5502 Jps
tushar@tushar:spark $
```

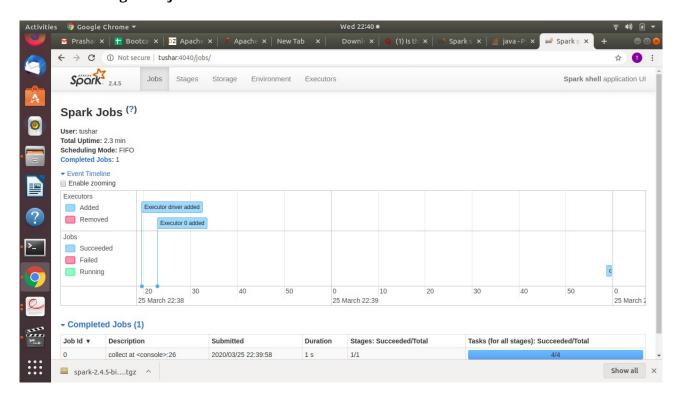
spark-shell running on both nodes

```
Terminal
 File Edit View Search Terminal Help
20/03/25 19:30:53 WARN NativeCodeLoader: Unable to load native-hadoop library for your pl
atform... using builtin-java classes where applicable
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLeve
31).
20/03/25 19:31:00 WARN MacAddressUtil: Failed to find a usable hardware address from the
network interfaces; using random bytes: 2f:cc:9f:47:33:b5:4f:67
Spark context Web ÚI available at http://tushar:4040
Spark context available as 'sc' (master = spark://tushar:7077, app id = app-2020032519310
1-0000).
 Spark session available as 'spark'.
Welcome to
Using Scala version 2.11.12 (OpenJDK 64-Bit Server VM, Java 11.0.6)
Type in expressions to have them evaluated.
Type :help for more information.
scala>
```

Application running

```
Terminal
File Edit View Search Terminal Help
Welcome to
                                      version 2.4.5
Using Scala version 2.11.12 (OpenJDK 64-Bit Server VM, Java 1.8.0 242)
Type in expressions to have them evaluated.
Type :help for more information.
scala> val sr=sc.parallelize(1 to 100,4)
sr: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[0] at parallelize at <
console>:24
scala > sr.map(x => x).collect
res0: Array[Int] = Array(1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56,
57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76,
77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96,
97, 98, 99, 100)
scala>
```

4 tasks in 1 stage in 1 job



Application tasks utilize both worker nodes



Q2. Using pyspark, demonstrate using a method of spark context that creates a series of RDDs

when working on a collection of numbers. Create a collection of numbers say list of 1 to 20. While using a method of spark context to create RDDs, print a list of even numbers. Demonstrate usage of cache and benefit in performance.

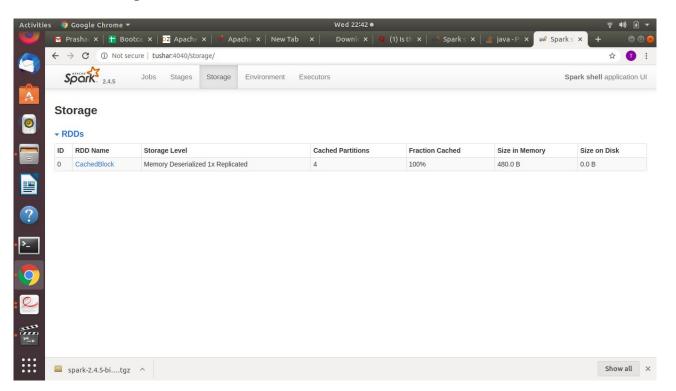
Ans-2 pyspark series(1 to 20) and printing even numbers using lambda function

```
Terminal
                                                                             File Edit View Search Terminal Help
 Type "help", "copyright", "credits" or "license" for more information.
 20/03/25 22:11:19 WARN Utils: Your hostname, tushar resolves to a loopback addre
ss: 127.0.1.1; using 172.20.10.3 instead (on interface wlp3s0)
20/03/25 22:11:19 WARN Utils: Set SPARK LOCAL IP if you need to bind to another
laddress
20/03/25 22:11:23 WARN NativeCodeLoader: Unable to load native-hadoop library fo
r your platform... using builtin-java classes where applicable
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLeve
l(newLevel).
Welcome to
Using Python version 3.6.9 (default, Nov 7 2019 10:44:02)
 SparkSession available as 'spark'.
 >>> r=sc.parallelize([1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20])
 >>> r.filter(lambda x: x % 2 == 0).collect()
 [2, <u>4</u>, 6, 8, 10, 12, 14, 16, 18, 20]
```

Usage of cache

```
Terminal
                                                                                        File Edit View Search Terminal Help
sr: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[0] at parallelize at <
console>:24
scala> sr.map(x => x).collect
res0: Array[Int] = Array(1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16,
17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56,
57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76,
77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96,
97, 98, 99, 100)
scala> sr.setName("CachedBlock").cache
res1: sr.type = CachedBlock ParallelCollectionRDD[0] at parallelize at <console>
:24
scala> sr.map(x => x).collect
res2: Array[Int] = Array(1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36,
37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56,
57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76,
77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96,
97, 98, 99, 100)
```

"CacheBlock" in gui



BENEFITS IN PERFORMANCE:-

Caching requires a consideration for number of nodes available, the relative priority of each RDD.

and the amount of memory available for caching, it often requires a good deal of trial-anderror.

And the more RDDs you have to consider, all the more complex this method becomes.

That's where automation plays a role. A tool that provides full-stack performance intelligence can

tell you if caching is a viable option in any particular case, and, if it would, how you could know

what to cache to get maximum performance. It can tell you the performance you can achieve if you

had five additional servers available, and whether adding those five servers to the mix makes sense.

There are generally three factors to consider in tuning memory usage: the amount of memory used

by your objects, the cost of accessing those objects, and the overhead of "garbage collection" (if

you have high turnover in terms of objects). Automating performance management across the stack

takes these factors into consideration, relieving you of the burden of manually tuning your system.

And:-

Due to Caching we get:-

- 1.) It Reduced latency
- 2.) It provides Content availability
- 3.) It Avoids network congestion
- 4.) Increases Response time:
- 5.) Increases System throughput
- 6.) It reduces the overhead from server resources.
- 7.) It increases the performance of the application by serving user with cached output.