

1. Pseudo-code for PageRank Program in Spark Scala

- The following is the Pseudo-code which heavily borrowed Scala syntax.

```
1  /* PAGERANK-SPARK-SCALA */
2
3
4  val k ← user set k value (number of linear chains)
5
6  val outlinkListBuffer ← List Buffer of tuple (Int, Int)
7  val rankListBuffer ← list buffer of tuple (Int, Double) containing dummy page (0, 0)
8
9  // generate outlink (p1 → p2) values
10 for (i ← 1 to pow(k, 2)) {
11   if (i % k == 0) {
12     outlinkListBuffer += ((i, 0))
13   } else {
14     outlinkListBuffer += ((i, i + 1))
15   }
16 }
17
18 // generate initial PageRank values
19 for (i ← 1 to pow(k, 2).toInt) {
20   rankListBuffer += ((i, 1 / pow(k, 2)))
21 }
22
23 // create outlink RDD
24 val outlinkRDD = sc.parallelize(outlinkListBuffer).partitionBy(new HashPartitioner(10))
25
26 // create inlink RDD and join with outlink RDD
27 val inlinkRDD = outlinkRDD
28   .map{case (from, to) => (to, from)}
29   .partitionBy(outlinkRDD.partitioner.get)
30   // outer-join with outlink graph to get inlink/outlink of each page
31   .fullOuterJoin(outlinkRDD)
32   // map in appropriate format to help join with PageRank graph later on
33   .map {
34     case (page, (inlink, outlink)) => {
35       if (page == 0) {
36         (inlink.getOrElse(0), page) // handle dummy page
37       } else {
38         (inlink.getOrElse(-1), page) // handle other every page
39       }
40     }
41   }
42   .persist()
43
44 // create PageRank RDD
45 var rankRDD = sc.parallelize(rankListBuffer).partitionBy(outlinkRDD.partitioner.get)
46
47 val alpha = 0.85
48
49 // iteration until convergence
50 val iteration ← user set number of PageRank iterations
51 for (i ← 1 to iteration) {
52
53   // join edge RDD with PageRank RDD and output according to the page type
54   val joinRDD = inlinkRDD
55     .leftOuterJoin(rankRDD)
56     .map {
57       case (outlink, (page, someRank)) => {
58
59         val pageRank = (1 - alpha) / pow(k, 2)
60
61         // if the page has no inlink, the value is set to 0 by default
62         val rank = someRank.getOrElse(0.toDouble)
63
64         if (rank == 0.toDouble) {
65           (page, pageRank) // handle page with no inlink (only outgoing link)
66         } else if (page == 0) {
67           (page, rank) // handle dangling page
68         } else {
69           (page, pageRank + alpha * rank) // handle page with outlink
70         }
71       }
72     }
73     .reduceByKey(_ + _) // aggregate to get dangling page PR mass
74
75   // read out total dangling PageRank mass in dummy page 0
76   val danglingPageRankSum = joinRDD.lookup(0).head
77
78   // distribute dangling PageRank mass evenly over all real pages
79   rankRDD = joinRDD
80     .map {
81       case (page, rank) => {
82         if (page == 0) {
83           (page, 0.toDouble) // if page is dummy, set its PageRank to 0
84         } else {
85           (page, rank + alpha * danglingPageRankSum / pow(k, 2)) // update PageRank
86         }
87       }
88     }
89
90 }
91
92 // for debugging purpose, check if PageRank values sum up to 1
93 print ← sum up all the PageRank values in RankRDD
94
95 // group to one output file, sort, and print first 20 records (page 0 ~ 19)
96 print ← rankRDD.repartition(1).sortByKey().take(20)
```

2. Source code for PageRank Program in Spark Scala

- <https://github.com/CS6240/hw-4-spark-jill666666/blob/master/src/main/scala/pr/PageRank.scala>

3. PageRank Program Output for k=100 and iterations=10

```
(0,0.0)
(1,1.5888125735610845E-5)
(2,2.9391220523344354E-5)
(3,4.086705694684288E-5)
(4,5.061974152456419E-5)
(5,5.890776462131797E-5)
(6,6.595084287305438E-5)
(7,7.193573524791716E-5)
(8,7.702118669812162E-5)
(9,8.134213026403417E-5)
(10,8.501325886262295E-5)
(11,1.047006992966952E-4)
(12,1.047006992966952E-4)
(13,1.047006992966952E-4)
(14,1.047006992966952E-4)
(15,1.047006992966952E-4)
(16,1.047006992966952E-4)
(17,1.047006992966952E-4)
(18,1.047006992966952E-4)
(19,1.047006992966952E-4)
```

- For each tuple (p1, p2), p1 is a page, and p2 is its corresponding PageRank.

4. PageRank RDD Lineage Report

- Iteration 1

```
2021-11-05 22:58:38,698 INFO root: (4) MapPartitionsRDD[12] at map at PageRank.scala:103 []
| ShuffledRDD[11] at reduceByKey at PageRank.scala:96 []
+- (4) MapPartitionsRDD[10] at map at PageRank.scala:79 []
|   MapPartitionsRDD[9] at leftOuterJoin at PageRank.scala:78 []
|   MapPartitionsRDD[8] at leftOuterJoin at PageRank.scala:78 []
|   CoGroupedRDD[7] at leftOuterJoin at PageRank.scala:78 []
+- (4) MapPartitionsRDD[5] at map at PageRank.scala:55 []
|   MapPartitionsRDD[4] at fullOuterJoin at PageRank.scala:53 []
|   MapPartitionsRDD[3] at fullOuterJoin at PageRank.scala:53 []
|   CoGroupedRDD[2] at fullOuterJoin at PageRank.scala:53 []
+- (4) MapPartitionsRDD[1] at map at PageRank.scala:51 []
|   | ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 []
+- (4) ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 []
+- (4) ParallelCollectionRDD[6] at parallelize at PageRank.scala:66 []
```

- Iteration 2

```
2021-11-05 22:57:42,685 INFO root: (4) MapPartitionsRDD[18] at map at PageRank.scala:103 []
| ShuffledRDD[17] at reduceByKey at PageRank.scala:96 []
+- (4) MapPartitionsRDD[16] at map at PageRank.scala:79 []
|   MapPartitionsRDD[15] at leftOuterJoin at PageRank.scala:78 []
|   MapPartitionsRDD[14] at leftOuterJoin at PageRank.scala:78 []
|   CoGroupedRDD[13] at leftOuterJoin at PageRank.scala:78 []
+- (4) MapPartitionsRDD[5] at map at PageRank.scala:55 []
|   MapPartitionsRDD[4] at fullOuterJoin at PageRank.scala:53 []
|   MapPartitionsRDD[3] at fullOuterJoin at PageRank.scala:53 []
|   CoGroupedRDD[2] at fullOuterJoin at PageRank.scala:53 []
+- (4) MapPartitionsRDD[1] at map at PageRank.scala:51 []
|   | ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 []
+- (4) ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 []
+- (4) MapPartitionsRDD[12] at map at PageRank.scala:103 []
|   ShuffledRDD[11] at reduceByKey at PageRank.scala:96 []
+- (4) MapPartitionsRDD[10] at map at PageRank.scala:79 []
|   MapPartitionsRDD[9] at leftOuterJoin at PageRank.scala:78 []
|   MapPartitionsRDD[8] at leftOuterJoin at PageRank.scala:78 []
|   CoGroupedRDD[7] at leftOuterJoin at PageRank.scala:78 []
+- (4) MapPartitionsRDD[5] at map at PageRank.scala:55 []
|   MapPartitionsRDD[4] at fullOuterJoin at PageRank.scala:53 []
|   MapPartitionsRDD[3] at fullOuterJoin at PageRank.scala:53 []
|   CoGroupedRDD[2] at fullOuterJoin at PageRank.scala:53 []
+- (4) MapPartitionsRDD[1] at map at PageRank.scala:51 []
|   | ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 []
+- (4) ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 []
+- (4) ParallelCollectionRDD[6] at parallelize at PageRank.scala:66 []
```

- Iteration 3

```

2021-11-05 22:53:03,644 INFO root: (4) MapPartitionsRDD[24] at map at PageRank.scala:103 []
| ShuffledRDD[23] at reduceByKey at PageRank.scala:96 []
+- (4) MapPartitionsRDD[22] at map at PageRank.scala:79 []
| MapPartitionsRDD[21] at leftOuterJoin at PageRank.scala:78 []
| MapPartitionsRDD[20] at leftOuterJoin at PageRank.scala:78 []
| CoGroupedRDD[19] at leftOuterJoin at PageRank.scala:78 []
+- (4) MapPartitionsRDD[5] at map at PageRank.scala:55 []
| MapPartitionsRDD[4] at fullOuterJoin at PageRank.scala:53 []
| MapPartitionsRDD[3] at fullOuterJoin at PageRank.scala:53 []
| CoGroupedRDD[2] at fullOuterJoin at PageRank.scala:53 []
+- (4) MapPartitionsRDD[1] at map at PageRank.scala:51 []
| ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 []
+- (4) ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 []
+- (4) MapPartitionsRDD[18] at map at PageRank.scala:103 []
| ShuffledRDD[17] at reduceByKey at PageRank.scala:96 []
+- (4) MapPartitionsRDD[16] at map at PageRank.scala:79 []
| MapPartitionsRDD[15] at leftOuterJoin at PageRank.scala:78 []
| MapPartitionsRDD[14] at leftOuterJoin at PageRank.scala:78 []
| CoGroupedRDD[13] at leftOuterJoin at PageRank.scala:78 []
+- (4) MapPartitionsRDD[5] at map at PageRank.scala:55 []
| MapPartitionsRDD[4] at fullOuterJoin at PageRank.scala:53 []
| MapPartitionsRDD[3] at fullOuterJoin at PageRank.scala:53 []
| CoGroupedRDD[2] at fullOuterJoin at PageRank.scala:53 []
+- (4) MapPartitionsRDD[1] at map at PageRank.scala:51 []
| ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 []
+- (4) ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 []
+- (4) MapPartitionsRDD[12] at map at PageRank.scala:103 []
| ShuffledRDD[11] at reduceByKey at PageRank.scala:96 []
+- (4) MapPartitionsRDD[10] at map at PageRank.scala:79 []
| MapPartitionsRDD[9] at leftOuterJoin at PageRank.scala:78 []
| MapPartitionsRDD[8] at leftOuterJoin at PageRank.scala:78 []
| CoGroupedRDD[7] at leftOuterJoin at PageRank.scala:78 []
+- (4) MapPartitionsRDD[5] at map at PageRank.scala:55 []
| MapPartitionsRDD[4] at fullOuterJoin at PageRank.scala:53 []
| MapPartitionsRDD[3] at fullOuterJoin at PageRank.scala:53 []
| CoGroupedRDD[2] at fullOuterJoin at PageRank.scala:53 []
+- (4) MapPartitionsRDD[1] at map at PageRank.scala:51 []
| ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 []
+- (4) ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 []
+- (4) ParallelCollectionRDD[6] at parallelize at PageRank.scala:66 []

```

Details for Job 0

Status: SUCCEEDED
Completed Stages: 6

Event Timeline
DAG Visualization

Completed Stages (6)

Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
5	lookup at PageRank.scala:99	2021/11/06 03:06:55	0.2 s	1/1			19.7 KiB	
4	map at PageRank.scala:79	2021/11/06 03:06:53	2 s	2/2			133.6 KiB	39.5 KiB
3	map at PageRank.scala:55	2021/11/06 03:06:52	0.9 s	2/2			158.3 KiB	81.8 KiB
2	parallelize at PageRank.scala:66	2021/11/06 03:06:46	0.8 s	2/2				51.8 KiB
1	map at PageRank.scala:51	2021/11/06 03:06:46	1.0 s	2/2				79.1 KiB
0	parallelize at PageRank.scala:47	2021/11/06 03:06:46	0.9 s	2/2				79.2 KiB

- If we look at the details for Job 0, it parallelizes the Ranks RDD which is the n followed by map and lookup operations.

Details for Job 1

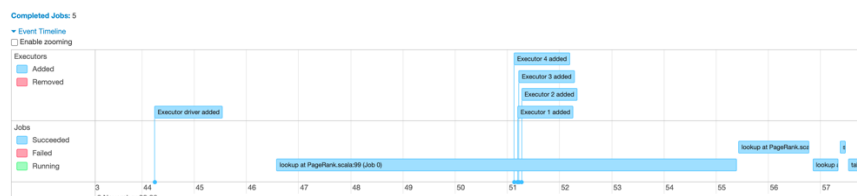
Status: SUCCEEDED
Completed Stages: 4
Skipped Stages: 5

Event Timeline
DAG Visualization

Completed Stages (4)

Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
14	lookup at PageRank.scala:99	2021/11/06 03:06:56	0.2 s	1/1			19.8 KiB	
13	map at PageRank.scala:79	2021/11/06 03:06:55	0.8 s	2/2			135.1 KiB	39.6 KiB
12	map at PageRank.scala:103	2021/11/06 03:06:55	0.1 s	2/2			39.5 KiB	53.3 KiB
11	map at PageRank.scala:55	2021/11/06 03:06:55	0.4 s	2/2			158.3 KiB	81.8 KiB

- Job 1 and Job 2 is almost identical in terms of the operations. We can also check skipped stages since they have already been computed in the previous round.



- Event timeline shows how each job is being triggered. The jobs include joining the Ranks RDD with the graph RDD, mapping to get the PageRank mass and calculate total dangling pages PageRank values.

Completed Jobs (5)

Page: 1

1 Pages. Jump to 1 . Show 100 items in a page. Go

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
4	take at PageRank.scala:126 take at PageRank.scala:126	2021/11/06 03:06:57	0.3 s	3/3 (11 skipped)	4/4 (22 skipped)
3	sum at PageRank.scala:120 sum at PageRank.scala:120	2021/11/06 03:06:57	0.1 s	1/1 (11 skipped)	2/2 (22 skipped)

- And we have Job 3 which sums up the PageRank values of all pages. This is for the debugging purpose to check whether the result sums up to 1. Finally, we have Job 4 that ‘take’ first 20 records (page 0 ~ 19) which is our final output.

5. Spark Cache and Reusability

- Below is the lineage of Ranks RDD after 3 iterations.

```
2021-11-05 21:59:58,537 INFO root: (4) MapPartitionsRDD[24] at map at PageRank.scala:103 []
| ShuffledRDD[23] at reduceByKey at PageRank.scala:96 []
+- (4) MapPartitionsRDD[22] at map at PageRank.scala:79 []
| MapPartitionsRDD[21] at leftOuterJoin at PageRank.scala:78 []
| MapPartitionsRDD[20] at leftOuterJoin at PageRank.scala:78 []
| CoGroupedRDD[19] at leftOuterJoin at PageRank.scala:78 []
+- (4) MapPartitionsRDD[5] at map at PageRank.scala:55 []
| MapPartitionsRDD[4] at fullOuterJoin at PageRank.scala:53 []
| MapPartitionsRDD[3] at fullOuterJoin at PageRank.scala:53 []
| CoGroupedRDD[2] at fullOuterJoin at PageRank.scala:53 []
+- (4) MapPartitionsRDD[1] at map at PageRank.scala:51 []
| ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 []
+- (4) ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 []
+- (4) MapPartitionsRDD[18] at map at PageRank.scala:103 []
| ShuffledRDD[17] at reduceByKey at PageRank.scala:96 []
+- (4) MapPartitionsRDD[16] at map at PageRank.scala:79 []
| MapPartitionsRDD[15] at leftOuterJoin at PageRank.scala:78 []
| MapPartitionsRDD[14] at leftOuterJoin at PageRank.scala:78 []
| CoGroupedRDD[13] at leftOuterJoin at PageRank.scala:78 []
+- (4) MapPartitionsRDD[5] at map at PageRank.scala:55 []
| MapPartitionsRDD[4] at fullOuterJoin at PageRank.scala:53 []
| MapPartitionsRDD[3] at fullOuterJoin at PageRank.scala:53 []
| CoGroupedRDD[2] at fullOuterJoin at PageRank.scala:53 []
+- (4) MapPartitionsRDD[1] at map at PageRank.scala:51 []
| ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 []
+- (4) ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 []
+- (4) MapPartitionsRDD[12] at map at PageRank.scala:103 []
| ShuffledRDD[11] at reduceByKey at PageRank.scala:96 []
+- (4) MapPartitionsRDD[10] at map at PageRank.scala:79 []
| MapPartitionsRDD[9] at leftOuterJoin at PageRank.scala:78 []
| MapPartitionsRDD[8] at leftOuterJoin at PageRank.scala:78 []
| CoGroupedRDD[7] at leftOuterJoin at PageRank.scala:78 []
+- (4) MapPartitionsRDD[5] at map at PageRank.scala:55 []
| MapPartitionsRDD[4] at fullOuterJoin at PageRank.scala:53 []
| MapPartitionsRDD[3] at fullOuterJoin at PageRank.scala:53 []
| CoGroupedRDD[2] at fullOuterJoin at PageRank.scala:53 []
+- (4) MapPartitionsRDD[1] at map at PageRank.scala:51 []
| ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 []
+- (4) ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 []
+- (4) ParallelCollectionRDD[6] at parallelize at PageRank.scala:66 []
```

Adding the print statement didn't change the lineage, and each of the iteration print statement has been printed only once (e.g., 1, 2, 3, ...). We can also find out that Spark is smart enough to re-use the pre-computed results.

- We can certainly see from the log how Spark manages the jobs efficiently by performing the lookup actions. Below is part of lookup logs copied from the log file.

```
- 2021-11-05 21:55:16,253 INFO spark.SparkContext: Starting job: lookup at PageRank.scala:99
- 2021-11-05 21:55:16,583 INFO scheduler.DAGScheduler: Got job 0 (lookup at PageRank.scala:99) with 1 output partitions
- 2021-11-05 21:55:16,583 INFO scheduler.DAGScheduler: Final stage: ResultStage 5 (lookup at PageRank.scala:99)
- 2021-11-05 21:55:18,001 INFO scheduler.DAGScheduler: ResultStage 5 (lookup at PageRank.scala:99) finished in 0.075 s
- 2021-11-05 21:55:18,006 INFO scheduler.DAGScheduler: Job 0 finished: lookup at PageRank.scala:99, took 1.752810 s
- 2021-11-05 21:55:18,024 INFO spark.SparkContext: Starting job: lookup at PageRank.scala:99
- 2021-11-05 21:55:18,029 INFO scheduler.DAGScheduler: Got job 1 (lookup at PageRank.scala:99) with 1 output partitions
- 2021-11-05 21:55:18,029 INFO scheduler.DAGScheduler: Final stage: ResultStage 14 (lookup at PageRank.scala:99)
- ...
- 2021-11-05 21:55:19,402 INFO scheduler.DAGScheduler: Job 8 finished: lookup at PageRank.scala:99, took 0.118823 s
- 2021-11-05 21:55:19,418 INFO scheduler.DAGScheduler: Got job 9 (lookup at PageRank.scala:99) with 1 output partitions
- 2021-11-05 21:55:19,418 INFO scheduler.DAGScheduler: Final stage: ResultStage 194 (lookup at PageRank.scala:99)
- 2021-11-05 21:55:19,525 INFO scheduler.DAGScheduler: ResultStage 194 (lookup at PageRank.scala:99) finished in 0.012 s
- 2021-11-05 21:55:19,525 INFO scheduler.DAGScheduler: Job 9 finished: lookup at PageRank.scala:99, took 0.110396 s
```

As we can see from the logs, even without Cache() or Persist() operations, Spark re-uses Ranks RDD result that has been pre-computed in the previous iteration.

- c. Refer to the link to the log files created after running two programs, one without Cache(), and the other one using Cache().
- Link: <https://github.com/CS6240/hw-4-spark-jill666666/tree/master/logs>

One most prominent difference between two cases is that the caching makes Spark program finds the blocks and store the blocks as values in memory to re-use them.

Below is some of the logs that only appears when cached.

- 2021-11-05 22:22:58,946 INFO storage.BlockManager: Found block rdd_66_2 locally
- 2021-11-05 22:22:58,946 INFO storage.BlockManager: Found block rdd_66_0 locally
- 2021-11-05 22:22:58,946 INFO storage.BlockManager: Found block rdd_66_3 locally
- 2021-11-05 22:22:58,946 INFO storage.BlockManager: Found block rdd_66_1 locally

Even though program without Cache() stores blocks such as broadcast_8 as values in memory, following are the logs that are only distinct to the program with Cache().

- 2021-11-05 22:22:57,028 INFO memory.MemoryStore: Block rdd_12_0 stored as values in memory (estimated size 87.9 KiB, free 366.2 MiB)
- 2021-11-05 22:22:57,029 INFO memory.MemoryStore: Block rdd_12_1 stored as values in memory (estimated size 87.9 KiB, free 366.1 MiB)

One another interesting thing is that we can see the CachedPartitions in the lineage of the Ranks RDD.

```
2021-11-05 22:22:57,884 INFO root: (4) MapPartitionsRDD[36] at map at PageRank.scala:103 [Memory Deserialized 1x Replicated]
| ShuffledRDD[35] at reduceByKey at PageRank.scala:96 [Memory Deserialized 1x Replicated]
+-(4) MapPartitionsRDD[34] at map at PageRank.scala:79 [Memory Deserialized 1x Replicated]
| MapPartitionsRDD[33] at leftOuterJoin at PageRank.scala:78 [Memory Deserialized 1x Replicated]
| MapPartitionsRDD[32] at leftOuterJoin at PageRank.scala:78 [Memory Deserialized 1x Replicated]
| CoGroupedRDD[31] at leftOuterJoin at PageRank.scala:78 [Memory Deserialized 1x Replicated]
+-(4) MapPartitionsRDD[5] at map at PageRank.scala:55 [Memory Deserialized 1x Replicated]
| MapPartitionsRDD[4] at fullOuterJoin at PageRank.scala:53 [Memory Deserialized 1x Replicated]
| MapPartitionsRDD[3] at fullOuterJoin at PageRank.scala:53 [Memory Deserialized 1x Replicated]
| CoGroupedRDD[2] at fullOuterJoin at PageRank.scala:53 [Memory Deserialized 1x Replicated]
+-(4) MapPartitionsRDD[1] at map at PageRank.scala:51 [Memory Deserialized 1x Replicated]
| ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 [Memory Deserialized 1x Replicated]
+-(4) ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 [Memory Deserialized 1x Replicated]
+-(4) MapPartitionsRDD[130] at map at PageRank.scala:103 [Memory Deserialized 1x Replicated]
| CachedPartitions: 4; MemorySize: 351.6 KiB; DiskSize: 0.0 B
| ShuffledRDD[29] at reduceByKey at PageRank.scala:96 [Memory Deserialized 1x Replicated]
+-(4) MapPartitionsRDD[28] at map at PageRank.scala:79 [Memory Deserialized 1x Replicated]
| MapPartitionsRDD[27] at leftOuterJoin at PageRank.scala:78 [Memory Deserialized 1x Replicated]
| MapPartitionsRDD[26] at leftOuterJoin at PageRank.scala:78 [Memory Deserialized 1x Replicated]
| CoGroupedRDD[25] at leftOuterJoin at PageRank.scala:78 [Memory Deserialized 1x Replicated]
+-(4) MapPartitionsRDD[5] at map at PageRank.scala:55 [Memory Deserialized 1x Replicated]
| MapPartitionsRDD[4] at fullOuterJoin at PageRank.scala:53 [Memory Deserialized 1x Replicated]
| MapPartitionsRDD[3] at fullOuterJoin at PageRank.scala:53 [Memory Deserialized 1x Replicated]
| CoGroupedRDD[2] at fullOuterJoin at PageRank.scala:53 [Memory Deserialized 1x Replicated]
+-(4) MapPartitionsRDD[1] at map at PageRank.scala:51 [Memory Deserialized 1x Replicated]
| ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 [Memory Deserialized 1x Replicated]
+-(4) ParallelCollectionRDD[0] at parallelize at PageRank.scala:47 [Memory Deserialized 1x Replicated]
+-(4) MapPartitionsRDD[24] at map at PageRank.scala:103 [Memory Deserialized 1x Replicated]
| CachedPartitions: 4; MemorySize: 351.6 KiB; DiskSize: 0.0 B
```

6. Spark Scala PageRank Program on EMR (1 Master & 5 Workers)

- k = 10,000 caused java heap space 'OutOfMemory' error. Setting k = 1,000 instead with 10 iterations resulted in the following output.
- Output & Log: <https://github.com/CS6240/hw-4-spark-jill666666/tree/master/aws-output>
- Running Time: Approximately 1 minute

PageRank Spark Cluster		j-ZVB8S8P6OC084	Terminated	2021-11-05 15:38 (UTC-4)	9 minutes	48
Summary		Master public DNS: ec2-50-19-15-88.compute-1.amazonaws.com				
Termination protection: Off		Tags: --				
Hardware		Master: Terminated 1 m5.xlarge Core: Terminated 5 m5.xlarge Task: --				
View cluster details		View monitoring details				
Steps		Name	Status	Start time (UTC-4)	Elapsed time	View all interactive jobs
		PageRank	Completed	2021-11-05 15:44 (UTC-4)	54 seconds	Bootstrap actions Name No bootstrap actions available
		Setup Hadoop Debugging	Completed	2021-11-05 15:44 (UTC-4)	8 seconds	

- Thoughts: There are some approaches I would like to take, if more time has been allowed, to address the memory problem. One thing can be generating the input in parallel and create input RDD that can save a way much more memory. Another solution can be storing the intermediate output every iteration in the file system, like what we did in the MapReduce PageRank program.

7. Pseudo-code for PageRank Program in MapReduce

```
/* PAGERANK-MAPREDUCE */

input <- each line of the text follows the format (page, [outlinks], PageRank)

PR_MASS_COUNTER <- global counter that sums up the PageRank mass values
TOTAL_NUM_VERTICES <- global counter that counts the total number of vertices

// takes in the input text file and calculate the total number of vertices
// and the dangling page PageRank mass sum.
DanglingPageMapper(text) {
  tokens <- split input text file by "\n"
  for (String token in tokens) {
    adjacencyList, pageRank <- unpack tokens
    if (adjacencyList is empty) {
      PR_MASS_COUNTER <- increment by PageRank value
    }
    TOTAL_NUM_VERTICES <- increment by 1
  }
}

// also takes in the input text file and emit pair (vertex ID, Vertex object).
// each Vertex object contains its ID, pageRank, adjacencyList, and boolean
// value flag (dangling page ? true : else false).
// emit twice, one with key vertex ID, and the another one with key inlink ID.
PageRankMapper(text) {
  tokens <- split input text file by "\n"
  for (String token in tokens) {
    vertexId, adjacencyList, pageRank <- unpack tokens
    if (adjacencyList is empty) {
      inlinkId <- get inlink vertex from adjacencyList
      emit(inlinkId, Vertex(vertexId, pageRank, adjacencyList, flag false))
    }
    emit(vertexId, Vertex(vertexId, pageRank, adjacencyList, flag true))
  }
}

PageRankReducer(vertexId, [vertices ...]) {

  // retrieve value from the global counters which will help calculate and
  // update the new PageRank for each page
  Setup() {
    danglingPageRankSum <- get value from global counter PR_MASS_COUNTER
    totalNumVertices <- get value from global counter TOTAL_NUM_VERTICES
  }

  alpha = 0.85
  adjacencyList = []
  newPageRank = ((1 - alpha) / totalNumVertices)
               + (alpha * danglingPageRankSum / totalNumVertices)

  // iterate through object vertices
  // for each vertex, if flag shows false (not a dangling page), calculate the
  // final PageRank.
  // get the page's true adjacency list.
  for (object vertex in vertices) {
    if (vertex.flag is false) {
      newPageRank += alpha * vertex.pageRank / vertex.adjacencyList.length
    } else {
      adjacencyList <- vertex.adjacencyList
    }
  }

  // here the reducer output has same format as the input, since this output
  // should become the input for the mapper in next iteration
  emit(vertexId, (adjacencyList, newPageRank))
}

// chain multiple jobs and iteratively run PageRank program
RunMapReduce() {
  numIterations <- number of iterations
  for (iter from 1 to numIterations) {
    job1 <- assign job DanglingPageMapper only
    job2 <- assign job PageRankMapper followed by PageRankReducer
    // takes in the input and run job1
    input -> run job1
    // run job2 and store the result in output
    output <- run job2
    // the output becomes the new input and will be read in the next iteration
    input <- output
  }
}
```

8. Source code for PageRank Program in MapReduce

- <https://github.com/CS6240/hw-4-mapreduce-jill666666/blob/master/src/main/java/pr/PageRank.java>

9. MapReduce PageRank Program on EMR (1 Master & 5 Workers)

- Logs: <https://github.com/CS6240/hw-4-mapreduce-jill666666/tree/master/aws-log>
- Output: <https://github.com/CS6240/hw-4-mapreduce-jill666666/tree/master/aws-output>
- Running Time: Approximately 50 minutes

Name	ID	Status	Creation time (UTC-4) ▾	Elapsed time	Normalized instance hours
<div><div></div><div>PageBark MacReduce Cluster</div></div>	j-2SIUS7DRG3X7	Terminated All steps completed	2021-11-05 16:11 (UTC-4)	55 minutes	48
Summary	Steps			View all interactive jobs	Bootstrap actions
Master public DNS: ec2-54-152-144-92.compute-1.amazonaws.com	Name	Status	Start time (UTC-4) ▾	Elapsed time	Name
Termination protection: Off	Custom JAR	Completed	2021-11-05 16:17 (UTC-4)	48 minutes	
Tags: --	Setup Hadoop Debugging	Completed	2021-11-05 16:17 (UTC-4)	4 seconds	
Hardware					No bootstrap actions available
Master: Terminated 1 m5.xlarge					
Core: Terminated 5 m5.xlarge					
Task: --					
View cluster details	View monitoring details				