

1. Pseudo-code for 5 Twitter-follower-count Programs

- RDD-G

FOLLOWER-COUNT-RDD-G

```
val textFile <- read input csv file
val counts =
  textFile
    .map(edge => edge.split(",")(1)) // split edge to get followedID
    .map(followedID => (followedID, 1)) // map each node to (node, 1)
    .groupByKey() // group by key (followedID)
    .map(pair => (pair._1, pair._2.sum)) // for each node, aggregate count
    .filter{case (userID, count) => count % 100 == 0} // get counts divisible by 100

output <- counts
```

- RDD-R

FOLLOWER-COUNT-RDD-R

```
val textFile <- read input csv file
val counts =
  textFile
    .map(edge => edge.split(",")(1)) // split edge to get followedID
    .map(followedID => (followedID, 1)) // map each node to (node, 1)
    .reduceByKey(_ + _) // reduce by key (followedID) which aggregates the count
    .filter{case (userID, count) => count % 100 == 0} // get counts divisible by 100

output <- counts
```

- RDD-F

FOLLOWER-COUNT-RDD-F

```
val textFile <- read input csv file
val counts =
  textFile
    .map(edge => edge.split(",")(1)) // split edge to get followedID
    .map(followedID => (followedID, 1)) // map each node to (node, 1)
    .foldByKey(0)(_ + _) // fold by key (followedID) which aggregates the count
    .filter{case (userID, count) => count % 100 == 0} // get counts divisible by 100

output <- counts
```

- RDD-A

FOLLOWER-COUNT-RDD-A

```
val textFile <- read input csv file

val initialCount = 0
val sumFollowerCounts = (c1: Int, c2: Int) => c1 + c2 // rule summing up int values
val sumPartitionCounts = (p1: Int, p2: Int) => p1 + p2 // rule summing up int values

val counts = textFile
  .map(edge => edge.split(",")(1)) // split edge to get followedID
  .map(followedID => (followedID, 1)) // map each node to (node, 1)
  .aggregateByKey(initialCount)(sumFollowerCounts, sumPartitionCounts)
  // sum up each of the follower counts & aggregate the partition counts
  .filter{case (userID, count) => count % 100 == 0}
  // get counts divisible by 100

output <- counts
```

- DSET

```
FOLLOWER-COUNT-DSET

val textFile <- read input csv file
val rdd = textFile
    .map(edge => edge.split(",")(1)) // split edge to get followedID
    .map(followedID => (followedID, 1)) // map each node to (node, 1)

val ds = spark.createDataset(rdd) // convert RDD to DataSet

val counts =
    ds
    .groupBy(ds.col("_1")) // group by the first column (followedID)
    .agg(sum(ds.col("_2"))) // aggregate the values of second column (count)
    .filter("sum(_2) % 100 == 0") // get counts divisible by 100

output <- counts
```

2. Source code for 5 Twitter-follower-count Programs

- RDD-G: <https://github.com/CS6240/hw-3-jill666666/blob/master/src/main/scala/fc/RDD-G.scala>
- RDD-R: <https://github.com/CS6240/hw-3-jill666666/blob/master/src/main/scala/fc/RDD-R.scala>
- RDD-F: <https://github.com/CS6240/hw-3-jill666666/blob/master/src/main/scala/fc/RDD-F.scala>
- RDD-A: <https://github.com/CS6240/hw-3-jill666666/blob/master/src/main/scala/fc/RDD-A.scala>
- DSET: <https://github.com/CS6240/hw-3-jill666666/blob/master/src/main/scala/fc/DSET.scala>

3. Aggregation and Shuffle

* 'toDebugString' method logs the RDD lineage graph. Using this we can find out how the shuffle, partition, and parallelism work. Indentation indicates the shuffle boundary.

** Value members referenced from: <https://spark.apache.org/docs/latest/api/scala/index.html>

a. RDD-G

```
2021-10-22 22:14:08,247 INFO root: (40) MapPartitionsRDD[6] at filter at RDD-G.scala:31 []
| MapPartitionsRDD[5] at map at RDD-G.scala:30 []
| ShuffledRDD[4] at groupByKey at RDD-G.scala:28 []
+--(40) MapPartitionsRDD[3] at map at RDD-G.scala:28 []
| MapPartitionsRDD[2] at map at RDD-G.scala:27 []
| /Users/sunho/Dropbox/Boston/CS6240/twitter-dataset/data/ MapPartitionsRDD[1] at textFile at RDD-G.scala:25 []
| /Users/sunho/Dropbox/Boston/CS6240/twitter-dataset/data/ HadoopRDD[0] at textFile at RDD-G.scala:25 []
```

- First of all, 'groupByKey()' reshuffles the data. Afterward, we use 'RDD.sum()' equivalent in order to sum up the data. Shuffling happens before the aggregation. The operation is very expensive due to the shuffling.
- Ref: [Apache Spark ScalaDoc groupByKey\(\)](#)

b. RDD-R

```
2021-10-22 22:22:26,061 INFO root: (40) MapPartitionsRDD[5] at filter at RDD-R.scala:30 []
| ShuffledRDD[4] at reduceByKey at RDD-R.scala:29 []
+--(40) MapPartitionsRDD[3] at map at RDD-R.scala:28 []
| MapPartitionsRDD[2] at map at RDD-R.scala:27 []
| /Users/sunho/Dropbox/Boston/CS6240/twitter-dataset/data/ MapPartitionsRDD[1] at textFile at RDD-R.scala:25 []
| /Users/sunho/Dropbox/Boston/CS6240/twitter-dataset/data/ HadoopRDD[0] at textFile at RDD-R.scala:25 []
```

- Here, unlike 'groupByKey()', the 'reduceByKey()' work as a combiner. Before shuffling the data, it can combine values for each key on each partition. Basically, it works similarly to a 'combiner' in MapReduce.

- Ref: [Apache Spark ScalaDoc reduceByKey\(\)](#)

c. RDD-F

```
2021-10-22 22:24:39,189 INFO root: (40) MapPartitionsRDD[5] at filter at RDD-F.scala:30 []
| ShuffledRDD[4] at foldByKey at RDD-F.scala:29 []
+- (40) MapPartitionsRDD[3] at map at RDD-F.scala:28 []
| MapPartitionsRDD[2] at map at RDD-F.scala:27 []
| /Users/sunho/Dropbox/Boston/CS6240/twitter-dataset/data/ MapPartitionsRDD[1] at textFile at RDD-F.scala:25 []
| /Users/sunho/Dropbox/Boston/CS6240/twitter-dataset/data/ HadoopRDD[0] at textFile at RDD-F.scala:25 []
```

- ‘foldByKey()’ merges the values for each key using the combine functions. But this does not cause shuffle.
- Ref: [Apache Spark ScalaDoc foldByKey\(\)](#)

d. RDD-A

```
2021-10-22 22:32:33,049 INFO root: (40) MapPartitionsRDD[5] at filter at RDD-A.scala:36 []
| ShuffledRDD[4] at aggregateByKey at RDD-A.scala:34 []
+- (40) MapPartitionsRDD[3] at map at RDD-A.scala:33 []
| MapPartitionsRDD[2] at map at RDD-A.scala:32 []
| /Users/sunho/Dropbox/Boston/CS6240/twitter-dataset/data/ MapPartitionsRDD[1] at textFile at RDD-A.scala:25 []
| /Users/sunho/Dropbox/Boston/CS6240/twitter-dataset/data/ HadoopRDD[0] at textFile at RDD-A.scala:25 []
```

```
def aggregateByKey[U](zeroValue: U, numPartitions: Int)(seqOp: (U, V) => U, combOp: (U, U) => U)(implicit arg0: ClassTag[U]): RDD[(K, U)]
```

Aggregate the values of each key, using given combine functions and a neutral "zero value". This function can return a different result type, U, than the type of the values in this RDD, V. Thus, we need one operation for merging a V into a U and one operation for merging two U's, as in `scala.TraversableOnce`. The former operation is used for merging values within a partition, and the latter is used for merging values between partitions. To avoid memory allocation, both of these functions are allowed to modify and return their first argument instead of creating a new U.

- ‘aggregateByKey()’ use the given combine functions to aggregate both partition and the value pairs, but it will not cause shuffle.
- Ref: [Apache Spark ScalaDoc aggregateByKey\(\)](#)

e. DSET

```
== Physical Plan ==
*(2) Project [_1#3, sum(_2)#11L]
+- *(2) Filter (isNotNull(sum(cast(_2#4 as bigint)))#16L) AND ((sum(cast(_2#4 as bigint)))#16L % 100) = 0))
+- *(2) HashAggregate(keys=[_1#3], functions=[sum(cast(_2#4 as bigint))])
+- Exchange hashpartitioning(_1#3, 200), ENSURE_REQUIREMENTS, [id=#72]
+- *(1) HashAggregate(keys=[_1#3], functions=[partial_sum(cast(_2#4 as bigint))])
+- *(1) SerializeFromObject [staticinvoke(class org.apache.spark.unsafe.types.UTF8String, StringType, fromString, knownNotNull(assertNotNull(input[0, scala.Tuple2, true]))._1, true, false) AS _1#3, knownNotNull(assertNotNull(input[0, scala.Tuple2, true]))._2 AS _2#4]
+- Scan[obj#2]
```

- Above is the physical plan we get by running ‘counts.explain()’ where variable count can be attained by following the steps on the DSET pseudo-code. The ‘groupBy(col)’ doesn’t combine, but instead reshuffle the data. Afterward, ‘agg(col)’ and ‘sum(col)’ aggregate the data.
- Ref: [Apache Spark ScalaDoc groupBy\(\)](#)

4. Pseudo-code for Triangle-counting Programs

a. RS-R

TRIANGLE-COUNT-RS-RDD

```
val textFile <- read input csv file

// create (X -> Y) RDD
val XtoYRDD =
  textFile
    .map(edge => edge.split(","))
    // split edge to get nodes
    .map(nodes => (nodes(0).toInt, nodes(1).toInt))
    // map nodes to (followerID, followedID)
    .filter{case (from, to) => from < MAX_VALUE && to < MAX_VALUE}
    // only get node IDs less than user-defined max value

// create (Y -> Z) RDD
val YtoZRDD =
  XtoYRDD
    .map{case (from, to) => (to, from)} // swap each node's position

// create (Z -> X) RDD
// this represents candidates for valid closing edges to form triangle
val ZtoXRDD = XtoYRDD
    .join(YtoZRDD) // by joining we acquire (X, Z) for each key Y
    .map{case (_, (from, to)) => (from, to)} // (X, Z) -> (Z, X)

// valid closing edges which form triangle with given nodes
// check equality between the edges and possible closing edges
val closingEdges =
  XtoYRDD // contains all edges (can also use YtoZRDD or any edges RDD)
    .join(ZtoXRDD) // ensure followerID equality
    .filter{case (_, (to1, to2)) => to1 == to2} // ensure followedID equality
    .map{case (from, (to1, _)) => (from, to1)} // get closing edge

val triangleCount = closingEdges.count() / 3 // remove duplicate counts

output <- triangleCount
```

b. RS-D

TRIANGLE-COUNT-RS-DATFRAME

```
val textFile <- read input csv file

// create edge RDD
val edgeRDD =
  textFile
    .map(edge => edge.split(",")) // split edge to get nodes
    .map(nodes => (nodes(0).toInt, nodes(1).toInt))
    // map nodes to (followerID, followedID)
    .filter{case (from, to) => from < MAX_VALUE && to < MAX_VALUE}
    // only get node IDs less than user-defined max value

// convert edge RDD to (X->Y) DataFrame
val XtoYDF = edgeRDD.toDF("X", "Y")

// construct (Y -> Z) DataFrame
val YtoZDF = XtoYDF.toDF("Y", "X").withColumnRenamed("X", "Z")

// construct (Z -> X) DataFrame
// this represents candidates for valid closing edges to form triangle
val ZtoXDF =
  XtoYDF
    .join(YtoZDF, XtoYDF("Y") === YtoZDF("Y") && XtoYDF("X") !== YtoZDF("Z"))
    .select("Z", "X")

// valid closing edges which form triangle with given nodes
val closingEdgeDF =
  XtoYDF
    .withColumnRenamed("Y", "Z")
    .as("XtoZ")
    .join(ZtoXDF.as("ZtoX"), $"XtoZ.X" === $"ZtoX.Z" && $"XtoZ.Z" === $"ZtoX.X")
    // ensure followerID & followedID equality

val closingEdgeCount = closingEdgeDF.count()

val triangleCount = closingEdgeCount / 3 // remove duplicate counts

output <- triangleCount
```

c. Rep-R

```

TRIANGLE-COUNT-REP-RDD

val textFile <- read input csv file

// create edge RDD
val edges =
  textFile
    .map(edge => edge.split(",")) // split edge to get nodes
    .map(nodes => (nodes(0).toInt, nodes(1).toInt))
    // map nodes to (followerID, followedID)
    .filter(case (from, to) => from < MAX_VALUE && to < MAX_VALUE)
    // only get node IDs less than custom set max value

// create edge map which has set of followedIDs for each followerID
val edgeMap =
  edges
    .map(case (from, to) => (from, Set(to)))
    .reduceByKey(_ ++ _) // reduce by followerID to get set of followedIDs
    .collectAsMap()

// create broadcast variable map
val broadcastMap = sc.broadcast( edgeMap )

// calculate triangle count
val triangleCount =
  edges.mapPartitions(iter => {
    iter.flatMap {
      case (nodeX, nodeY) => broadcastMap.value.get(nodeY).map {
        // for each Y followed by X, get set of IDs Y is following
        setY => setY.foreach(nodeZ => broadcastMap.value.get(nodeZ).foreach {
          // for each Z followed by Y, get set of IDs Z is following
          setZ => if (setZ.contains(nodeX)) {
            // if Z follows X, triangle has been formed
            accumulator.add(1) // increment the accumulator
          }
        })
      }
    }
  })
  .collect()

output <- accumulator.value / 3 // remove duplicate counts

```

d. Rep-D

```

TRIANGLE-COUNT-REP-DATFRAME

val textFile <- read input csv file

// create edge RDD
val edgeRDD = textFile
  .map(edge => edge.split(",")) // split edge to get nodes
  .map(nodes => (nodes(0).toInt, nodes(1).toInt))
  // map nodes to (followerID, followedID)
  .filter(case (from, to) => from < MAX_VALUE && to < MAX_VALUE)
  // only get node IDs less than custom set max value

// convert edge RDD to (X->Y) DataFrame
val XtoYDF = edgeRDD.toDF("X", "Y")

// construct (Y -> Z) DataFrame
val YtoZDF = XtoYDF.toDF("Y", "X").withColumnRenamed("X", "Z")

// use broadcast join to construct (Z -> X) DataFrame
val ZtoXDF =
  XtoYDF
    .join(broadcast(YtoZDF),
      XtoYDF("Y") <=> YtoZDF("Y") && XtoYDF("X") != YtoZDF("Z"))
    // join on node Y where X is not equal to Z (filter out unnecessary Path2)
    .select("Z", "X")

// use broadcast join to get all valid closing edges
val closingEdgeDF =
  XtoYDF
    .withColumnRenamed("Y", "Z")
    .as("XtoZ")
    .join(broadcast(ZtoXDF.as("ZtoX")),
      $"XtoZ.X" <=> $"ZtoX.Z" && $"XtoZ.Z" === $"ZtoX.X")
    // ensure followerID & followedID equality

val closingEdgeCount = closingEdgeDF.count()

val triangleCount = closingEdgeCount / 3 // remove duplicate counts

output <- triangleCount

```

5. Source code for Triangle-counting Programs

- RS-R: <https://github.com/CS6240/hw-3-jill666666/blob/master/src/main/scala/tc/RS-R.scala>
- RS-D: <https://github.com/CS6240/hw-3-jill666666/blob/master/src/main/scala/tc/RS-D.scala>
- Rep-R: <https://github.com/CS6240/hw-3-jill666666/blob/master/src/main/scala/tc/Rep-R.scala>
- Rep-D: <https://github.com/CS6240/hw-3-jill666666/blob/master/src/main/scala/tc/Rep-D.scala>

6. Run Triangle-counting Programs on EMR (1 Master & 4 Workers)

* The MAX_VALUE values from HW2 were 600, 2,400, and 6,000, but after receiving a feedback that the running time is too short, the values for HW3 have been updated to 600, 2,400, and 42,000. We also use MAX_VALUE of 12,000 to compare between the Spark and MapReduce implementations, which will be discussed in Question 8.

** Due to out-of-memory error, for Rep-D, the greatest MAX_VALUE has been set to 6,000, instead of 42,000 (Also tested with values 28,000 and 12,000, but same errors have been occurred).

a. RS-R

- MAX: 600 / Running Time: 1 min 20 seconds / Triangle Count: 572 / Machine Type: 1 Master (m4.xlarge) & 4 Workers (m4.xlarge)

Summary		Steps		Bootstrap actions	
Master: public DNS: ec2-3-237-81-101.compute-1.amazonaws.com Termination protection: Off Tags: --		Name Triangle Count Setup Hadoop Debugging		Status Completed Completed	
		Start time (UTC-4) 2021-10-22 03:11 (UTC-4) 2021-10-22 03:10 (UTC-4)		Elapsed time 1 minute 20 seconds	
Hardware Master: Terminated 1 m4.xlarge Core: Terminated 4 m4.xlarge Tags: --				No bootstrap actions available	

- MAX: 2,400 / Running Time: 1 min 4 seconds / Triangle Count: 8,798 / Machine Type: 1 Master (m4.xlarge) & 4 Workers (m4.xlarge)

Summary

Master public DNS: ec2-54-160-75-39.compute-1.amazonaws.com

Termination protection: Off

Tags: --

Hardware

Masters Terminated 1 m4.xlarge
 Core Terminated 4 m4.xlarge
 Task: --

Name	Status	Start time (UTC+)	Elapsed time
Triangle Count	Completed	2021-10-22 11:14 (UTC+)	1 minute
Setup Hadoop Debugging	Completed	2021-10-22 11:14 (UTC+)	4 seconds

No bootstrap actions available

- MAX: 42,000 / Running Time: 29 mins 8 seconds / Triangle Count: 4,912,313 / Machine Type: 1 Master (m4.xlarge) & 4 Workers (m4.xlarge)

Master (minidmge) **Worker (minidmge)**

j-2X7XO588ZGYO		Terminated All steps completed	2021-10-22 08:39 (UTC-4)	39 minutes	40
Summary					
Master public DNS: ec2-54-90-15-160.compute-1.amazonaws.com					
Termination protection: Off					
Tags: --					
Hardware					
Master: Terminated 1 m4.xlarge					
Core: Terminated 4 m4.xlarge					
Task: --					
Steps					
Name	Status	Start time (UTC-4)	Elapsed time	Bootstrap actions	
Triangle Count	Completed	2021-10-22 08:47 (UTC-4)	29 minutes		
Setup Hadoop Debugging	Completed	2021-10-22 08:46 (UTC-4)	8 seconds	No bootstrap actions available	

b. RS-D

- MAX: 600 / Running Time: 1 min 4 seconds / Triangle Count: 572 / Machine Type: 1 Master (m4.xlarge) & 4 Workers (m4.xlarge)

RS-D Max 600 Spark Cluster j-2CQJ3A5X9776I Terminated All steps completed 2021-10-22 03:14 (UTC-4) 12 minutes 40

Summary		Steps				Bootstrap actions	
Master public DNS: ec2-34-203-215-205.compute-1.amazonaws.com							
Termination protection: Off							
Tags: --							
		Name	Status	Start time (UTC-4)	Elapsed time	Name	
		Triangle Count	Completed	2021-10-22 03:22 (UTC-4)	1 minute		
		Setup Hadoop Debugging	Completed	2021-10-22 03:22 (UTC-4)	4 seconds	No bootstrap actions available	
Hardware							
Master: Terminated 1 m4.xlarge							
Core: Terminated 4 m4.xlarge							
ZooKeeper:							

- MAX: 2,400 / Running Time: 1 min 16 seconds / Triangle Count: 8,798 / Machine Type: 1 Master (m4.xlarge) & 4 Workers (m4.xlarge)

<div><div></div><div>RS-D Max 2400 Spark Cluster</div></div>		j-2LH3990XHTTP	Terminated All steps completed	2021-10-22 11:17 (UTC-4)	11 minutes	40
<div>Summary</div> <div><div>Master public DNS: ec2-100-25-190-119.compute-1.amazonaws.com</div><div>Termination protection: Off</div><div>Tags: --</div></div> <div>Hardware</div> <div><div>Master: Terminated 1 m4.xlarge</div><div>Core: Terminated 4 m4.xlarge</div></div>		<div>Steps</div> <div><div><div>Name</div><div>Status</div><div>Start time (UTC-4)</div><div>Elapsed time</div></div><div><div>Triangle Count</div><div>Completed</div><div>2021-10-22 11:25 (UTC-4)</div><div>1 minute</div></div><div><div>Setup Hadoop Debugging</div><div>Completed</div><div>2021-10-22 11:25 (UTC-4)</div><div>16 seconds</div></div></div> <div><div>View all interactive jobs</div></div>			<div>Bootstrap actions</div> <div><div>Name</div></div> <div>No bootstrap actions available</div>	

- MAX: 42,000 / Running Time: 3 min 4 seconds / Triangle Count: 4,912,313 / Machine Type: 1 Master (m4.xlarge) & 4 Workers (m4.xlarge)

<div><div></div><div>RS-D Max 42000 Spark Cluster</div></div>	j-199KGJIDYXLUS	Terminated All steps completed	2021-10-22 09:19 (UTC-4)	13 minutes	40												
<div><div>Summary</div><div><div>Master public DNS: ec2-100-25-37-185.compute-1.amazonaws.com</div><div>Termination protection: Off</div><div>Tags: --</div></div><div><div>Hardware</div><div><div>Master: Terminated 1 m4.xlarge</div><div>Core: Terminated 4 m4.xlarge</div><div>Task: --</div></div></div></div>	<div><div>Steps</div><table><thead><tr><th>Name</th><th>Status</th><th>Start time (UTC-4)</th><th>Elapsed time</th></tr></thead><tbody><tr><td>Triangle Count</td><td>Completed</td><td>2021-10-22 09:26 (UTC-4)</td><td>3 minutes</td></tr><tr><td>Setup Hadoop Debugging</td><td>Completed</td><td>2021-10-22 09:26 (UTC-4)</td><td>4 seconds</td></tr></tbody></table></div>	Name	Status	Start time (UTC-4)	Elapsed time	Triangle Count	Completed	2021-10-22 09:26 (UTC-4)	3 minutes	Setup Hadoop Debugging	Completed	2021-10-22 09:26 (UTC-4)	4 seconds	<div><div>View all interactive jobs</div></div>	<div><div>Bootstrap actions</div><table><thead><tr><th>Name</th></tr></thead><tbody><tr><td>No bootstrap actions available</td></tr></tbody></table></div>	Name	No bootstrap actions available
Name	Status	Start time (UTC-4)	Elapsed time														
Triangle Count	Completed	2021-10-22 09:26 (UTC-4)	3 minutes														
Setup Hadoop Debugging	Completed	2021-10-22 09:26 (UTC-4)	4 seconds														
Name																	
No bootstrap actions available																	

c. Rep-R

- MAX: 600 / Running Time: 1 min 4 seconds / Triangle Count: 572 / Machine Type: 1 Master (m4.xlarge) & 4 Workers (m4.xlarge)

<div><div></div><div>Rep-R Max 600 Spark Cluster</div></div>		j-2TJFCZIXRR3C2	Terminated All steps completed	2021-10-22 02:30 (UTC-4)	10 minutes	40												
<div>Summary</div> <div>Master public DNS: ec2-52-91-178-26.compute-1.amazonaws.com</div> <div>Termination protection: Off</div> <div>Tags: --</div> <div>Hardware</div> <div><div>Master: Terminated 1 m4.xlarge</div><div>Core: Terminated 4 m4.xlarge</div><div>Task: --</div></div>		<div>Steps</div> <table><thead><tr><th>Name</th><th>Status</th><th>Start time (UTC-4)</th><th>Elapsed time</th></tr></thead><tbody><tr><td>Triangle Count</td><td>Completed</td><td>2021-10-22 02:38 (UTC-4)</td><td>56 seconds</td></tr><tr><td>Setup Hadoop Debugging</td><td>Completed</td><td>2021-10-22 02:38 (UTC-4)</td><td>8 seconds</td></tr></tbody></table>		Name	Status	Start time (UTC-4)	Elapsed time	Triangle Count	Completed	2021-10-22 02:38 (UTC-4)	56 seconds	Setup Hadoop Debugging	Completed	2021-10-22 02:38 (UTC-4)	8 seconds	<div><div>View all interactive jobs</div><div>Bootstrap actions</div><table><thead><tr><th>Name</th></tr></thead><tbody><tr><td>No bootstrap actions available</td></tr></tbody></table></div>	Name	No bootstrap actions available
Name	Status	Start time (UTC-4)	Elapsed time															
Triangle Count	Completed	2021-10-22 02:38 (UTC-4)	56 seconds															
Setup Hadoop Debugging	Completed	2021-10-22 02:38 (UTC-4)	8 seconds															
Name																		
No bootstrap actions available																		

- MAX: 2,400 / Running Time: 1 min 10 seconds / Triangle Count: 8,798 / Machine Type: 1 Master (m4.xlarge) & 4 Workers (m4.xlarge)

<div><div></div><div>Rep-R Max 2400 Spark Cluster</div></div>		j-21I3840MBFAPX	Terminated All steps completed	2021-10-22 10:33 (UTC-4)	12 minutes	40												
<div>Summary</div> <div>Master public DNS: ec2-54-53-128-249.compute-1.amazonaws.com</div> <div>Termination protection: Off</div> <div>Tags: --</div> <div>Hardware</div> <div><div>Master: Terminated 1 m4.xlarge</div><div>Core: Terminated 4 m4.xlarge</div><div>Task: --</div></div>		<div>Steps</div> <table><thead><tr><th>Name</th><th>Status</th><th>Start time (UTC-4)</th><th>Elapsed time</th></tr></thead><tbody><tr><td>Triangle Count</td><td>Completed</td><td>2021-10-22 10:41 (UTC-4)</td><td>1 minute</td></tr><tr><td>Setup Hadoop Debugging</td><td>Completed</td><td>2021-10-22 10:41 (UTC-4)</td><td>10 seconds</td></tr></tbody></table>		Name	Status	Start time (UTC-4)	Elapsed time	Triangle Count	Completed	2021-10-22 10:41 (UTC-4)	1 minute	Setup Hadoop Debugging	Completed	2021-10-22 10:41 (UTC-4)	10 seconds	<div>View all interactive jobs</div> <div>Bootstrap actions</div> <table><thead><tr><th>Name</th></tr></thead><tbody><tr><td>No bootstrap actions available</td></tr></tbody></table>	Name	No bootstrap actions available
Name	Status	Start time (UTC-4)	Elapsed time															
Triangle Count	Completed	2021-10-22 10:41 (UTC-4)	1 minute															
Setup Hadoop Debugging	Completed	2021-10-22 10:41 (UTC-4)	10 seconds															
Name																		
No bootstrap actions available																		

- MAX: 42,000 / Running Time: 2 mins 12 seconds / Triangle Count: 4,912,313 / Machine Type: 1 Master (m4.xlarge) & 4 Workers (m4.xlarge)

▼		Master-R Max 42000 Spark Cluster		j-3NSVE728KM5YN	Terminated All steps completed	2021-10-22 08:19 (UTC-4)	13 minutes	40
Summary		Steps				View all interactive jobs		Bootstrap actions
Master public DNS: ec2-34-239-253-172.compute-1.amazonaws.com		Name		Status	Start time (UTC-4) ▼	Elapsed time	Name	
Termination protection: Off		Triangle Count		Completed	2021-10-22 08:27 (UTC-4)	2 minutes		
Tags: --		Setup Hadoop Debugging		Completed	2021-10-22 08:27 (UTC-4)	12 seconds	No bootstrap actions available	
Hardware		Master: Terminated 1 m4.xlarge						
		Core: Terminated 4 m4.xlarge						
		Task: --						

d. Rep-D

- MAX: 600 / Running Time: 1 min 28 seconds / Triangle Count: 572 / Machine Type: 1 Master (m4.xlarge) & 4 Workers (m4.xlarge)

Rep-D Max 600 Spark Cluster

j-2J83M7WGL394H

Terminated
All steps completed

2021-10-22 03:27 (UTC-4)

11 minutes

40

Summary

Master public DNS: ec2-107-21-6-109.compute-1.amazonaws.com

Termination protection: Off

Tags: --

Hardware

Master: Terminated 1 m4.xlarge

Core: Terminated 4 m4.xlarge

Task: --

Steps

Name	Status	Start time (UTC-4)	Elapsed time
Triangle Count	Completed	2021-10-22 03:34 (UTC-4)	1 minute
Setup Hadoop Debugging	Completed	2021-10-22 03:34 (UTC-4)	28 seconds

[View all interactive jobs](#)

No bootstrap actions available

Bootstrap actions

Name

- MAX: 2,400 / Running Time: 1 min 14 seconds / Triangle Count: 8,798 / Machine Type: 1 Master (m4.xlarge) & 4 Workers (m4.xlarge)

<div><div></div><div>Rep-D Max 2400 Spark Cluster</div></div>	j-379WV1GBYRC2Z	Terminated All steps completed	2021-10-22 10:21 (UTC-4)	11 minutes	40											
<div><div>Summary</div><div><div>Master public DNS: ec2-18-207-141-225.compute-1.amazonaws.com</div><div>Termination protection: Off</div><div>Tags: --</div></div><div><div>Hardware</div><div><div>Master: Terminated 1 m4.xlarge</div><div>Core: Terminated 4 m4.xlarge</div><div>Task: --</div></div></div></div>	<div><div>Steps</div><table><thead><tr><th>Name</th><th>Status</th><th>Start time (UTC-4)</th><th>Elapsed time</th></tr></thead><tbody><tr><td>Triangle Count</td><td>Completed</td><td>2021-10-22 10:29 (UTC-4)</td><td>1 minute</td></tr><tr><td>Setup Hadoop Debugging</td><td>Completed</td><td>2021-10-22 10:28 (UTC-4)</td><td>14 seconds</td></tr></tbody></table></div>	Name	Status	Start time (UTC-4)	Elapsed time	Triangle Count	Completed	2021-10-22 10:29 (UTC-4)	1 minute	Setup Hadoop Debugging	Completed	2021-10-22 10:28 (UTC-4)	14 seconds	<div><div>View all interactive jobs</div><div>Bootstrap actions</div><table><thead><tr><th>Name</th></tr></thead><tbody><tr><td>No bootstrap actions available</td></tr></tbody></table></div>	Name	No bootstrap actions available
Name	Status	Start time (UTC-4)	Elapsed time													
Triangle Count	Completed	2021-10-22 10:29 (UTC-4)	1 minute													
Setup Hadoop Debugging	Completed	2021-10-22 10:28 (UTC-4)	14 seconds													
Name																
No bootstrap actions available																

- MAX: 6,000** / Running Time: 1 min 14 seconds / Triangle Count: 131,654** / Machine Type: Master 1 (m4.xlarge) & Workers (m4.xlarge)

<div><div></div><div>Reo-D Max 6000 Spark Cluster</div></div>		j-1A6L0HED1W50C	Terminated All steps completed	2021-10-22 09:56 (UTC-4)	11 minutes	40												
<div>Summary</div> <div>Master public DNS: ec2-18-206-233-136.compute-1.amazonaws.com</div> <div>Termination protection: Off</div> <div>Tags: --</div> <div>Hardware</div> <div><div>Master: Terminated 1 m4.xlarge</div><div>Core: Terminated 4 m4.xlarge</div><div>Task: --</div></div>		<div>Steps</div> <table><thead><tr><th>Name</th><th>Status</th><th>Start time (UTC-4)</th><th>Elapsed time</th></tr></thead><tbody><tr><td>Triangle Count</td><td>Completed</td><td>2021-10-22 10:04 (UTC-4)</td><td>1 minute</td></tr><tr><td>Setup Hadoop Debugging</td><td>Completed</td><td>2021-10-22 10:04 (UTC-4)</td><td>14 seconds</td></tr></tbody></table>		Name	Status	Start time (UTC-4)	Elapsed time	Triangle Count	Completed	2021-10-22 10:04 (UTC-4)	1 minute	Setup Hadoop Debugging	Completed	2021-10-22 10:04 (UTC-4)	14 seconds	<div>View all interactive jobs</div> <div>Bootstrap actions</div> <table><thead><tr><th>Name</th></tr></thead><tbody><tr><td>No bootstrap actions available</td></tr></tbody></table>	Name	No bootstrap actions available
Name	Status	Start time (UTC-4)	Elapsed time															
Triangle Count	Completed	2021-10-22 10:04 (UTC-4)	1 minute															
Setup Hadoop Debugging	Completed	2021-10-22 10:04 (UTC-4)	14 seconds															
Name																		
No bootstrap actions available																		

7. EMR Output File std-err (1 Master & 4 Workers)

- RS-R: <https://github.com/CS6240/hw-3-jill666666/tree/master/aws-outputs/RS-R/1M4W>
- RS-D: <https://github.com/CS6240/hw-3-jill666666/tree/master/aws-outputs/RS-D/1M4W>
- Rep-R: <https://github.com/CS6240/hw-3-jill666666/tree/master/aws-outputs/Rep-R/1M4W>
- Rep-D: <https://github.com/CS6240/hw-3-jill666666/tree/master/aws-outputs/Rep-D/1M4W>

8. Run Triangle-counting Programs on EMR (1 Master & 8 Workers)

* The MAX_VALUE values from HW2 were 600, 2,400, and 6,000, but after receiving a feedback that the running time is too short, the values for HW3 have been updated to 600, 2,400, and 42,000. We also use MAX_VALUE of 12,000 to compare between the Spark and MapReduce implementations.

** Every implementations **except** Spark Rep-D has the same setup, MAX_VALUE and machine type, as well as triangle count result. As mentioned in the Question 6, due to out-of-memory error, Rep-D program has been ran using the MAX_VALUE of 6,000. The details including the comparison with the MapReduce are shown below.

1) Spark (HW3)

a. RS-R

- MAX: 12,000
- Running Time: 5 mins 20 seconds
- Triangle Count: 856,482
- Machine Type: 1 Master (m4.large) & 8 Workers (m4.large)

<div><div></div><div>RS-R 8 Workers Max 12000 Spark Cluster</div></div>		j-LW7CF7QD2ANE	Terminated All steps completed	2021-10-22 15:37 (UTC-4)	16 minutes	36													
<div>Summary</div> <div><div>Master public DNS: ec2-3-230-151-206.compute-1.amazonaws.com</div><div>Termination protection: Off</div><div>Tags: --</div></div> <div>Hardware</div> <div><div>Master: Terminated 1 m4.large</div><div>Core: Terminated 8 m4.large</div><div>Task: --</div></div>		<div>Steps</div> <table><thead><tr><th>Name</th><th>Status</th><th>Start time (UTC-4)</th><th>Elapsed time</th></tr></thead><tbody><tr><td>Triangle Count</td><td>Completed</td><td>2021-10-22 15:45 (UTC-4)</td><td>5 minutes</td></tr><tr><td>Setup Hadoop Debugging</td><td>Completed</td><td>2021-10-22 15:45 (UTC-4)</td><td>20 seconds</td></tr></tbody></table>			Name	Status	Start time (UTC-4)	Elapsed time	Triangle Count	Completed	2021-10-22 15:45 (UTC-4)	5 minutes	Setup Hadoop Debugging	Completed	2021-10-22 15:45 (UTC-4)	20 seconds	<div><div>View all interactive jobs</div><div>Bootstrap actions</div><table><thead><tr><th>Name</th></tr></thead><tbody><tr><td>No bootstrap actions available</td></tr></tbody></table></div>	Name	No bootstrap actions available
Name	Status	Start time (UTC-4)	Elapsed time																
Triangle Count	Completed	2021-10-22 15:45 (UTC-4)	5 minutes																
Setup Hadoop Debugging	Completed	2021-10-22 15:45 (UTC-4)	20 seconds																
Name																			
No bootstrap actions available																			

b. RS-D

- MAX: 12,000
- Running Time: 1 min 23 seconds
- Triangle Count: 856,482

<div><div><div></div><div></div></div><div>Replicated Join Max 12000 MR Cluster</div></div>		j-36VUUQQW6WNGS	Terminated All steps completed	2021-10-22 15:19 (UTC-4)	13 minutes	36
<div>Summary</div> <div>Master public DNS: ec2-3-239-27-118.compute-1.amazonaws.com</div> <div>Termination protection: Off</div> <div>Tags: --</div> <div>Hardware</div> <div>Master: Terminated 1 m4.large</div> <div>Core: Terminated 8 m4.large</div> <div>Task: --</div> <div><div>View cluster details</div><div>View monitoring details</div></div>		<div>Steps</div> <div><div><div>Name</div><div>Status</div><div>Start time (UTC-4)</div><div>Elapsed time</div></div><div><div>Custom JAR</div><div>Completed</div><div>2021-10-22 15:27 (UTC-4)</div><div>3 minutes</div></div><div><div>Setup Hadoop Debugging</div><div>Completed</div><div>2021-10-22 15:27 (UTC-4)</div><div>13 seconds</div></div></div>			<div><div>View all interactive jobs</div><div>Bootstrap actions</div><div><div>Name</div><div>No bootstrap actions available</div></div></div>	

9. EMR Output File std-err (1 Master & 8 Workers)

- RS-R: <https://github.com/CS6240/hw-3-jill666666/tree/master/aws-outputs/RS-R/1M8W>
- RS-D: <https://github.com/CS6240/hw-3-jill666666/tree/master/aws-outputs/RS-D/1M8W>
- Rep-R: <https://github.com/CS6240/hw-3-jill666666/tree/master/aws-outputs/Rep-R/1M8W>
- Rep-D: <https://github.com/CS6240/hw-3-jill666666/tree/master/aws-outputs/Rep-D/1M8W>