

Performance evaluation of sort

Experiment	Shared Memory	Linux	Hadoop Sort	Spark Sort
1 small.instance, 1GB dataset	19248	12030	16345	14283
1 small.instance, 15GB dataset	394813	241251	289256	273917
1 large.instance, 1GB dataset	17823	10696	14386	12645
1 large.instance, 15GB dataset	288412	203909	231884	215248
1 large.instance, 60GB dataset	1277862	1003473	1266273	1128645
4 small.instances, 1GB dataset	NA	NA	17934	16387
4 small.instances, 15GB dataset	NA	NA	263487	246743
4 small.instances, 60GB dataset	NA	NA	1126345	1066234

Small Instance:

1. 1GB Dataset:

Shared Memory Sort:

```
root@nfs: /  
root@nfs:/# python smallinsta1gb.py  
Enter number of threads 2  
thread t0 sorting done  
thread t1 sorting done  
Done! Time required is 19.2487653125Sec  
Sorting output is stored in merge0.txt file  
root@nfs:/#
```

Linux Sort:

```
root@nfs: /home  
root@nfs:/home# ls  
rand1gb.txt ubuntu  
root@nfs:/home# time sort rand1gb.txt > validatenewdata.txt  
  
real    0m12.030s  
user    0m28.792s  
sys     0m1.880s  
root@nfs:/home#
```

Hadoop Sort:

```

cc@hw8rohitl: ~/HW8/hadoop-3.2.0
cc@hw8rohitl:~/HW8/hadoop-3.2.0$ hadoop jar Hadoopsort.jar HadoopSort data1gbsmallinstance.in output
Time taken to sort 16345
19/04/28 15:07:00 INFO client.RMProxy: Connecting to ResourceManager at hadoop-f/192.168.2.30:8032
19/04/28 15:07:00 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Impl
ol interface and execute your application with ToolRunner to remedy this.
19/04/28 15:07:00 INFO input.FileInputFormat: Total input files to process : 1
19/04/28 15:07:00 INFO mapreduce.JobSubmitter: number of splits:120
19/04/28 15:07:01 INFO Configuration.deprecation: yarn.resourcemanager.system-metrics-publisher.enabled is
nstead, use yarn.system-metrics-publisher.enabled
19/04/28 15:07:01 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1524521941871_0403
19/04/28 15:07:01 INFO impl.YarnClientImpl: Submitted application application_1524521941871_0403
19/04/28 15:07:01 INFO mapreduce.Job: The url to track the job: http://hadoop-f:8088/proxy/application_152
3/
19/04/28 15:07:01 INFO mapreduce.Job: Running job: job_1524521941871_0403
19/04/28 15:07:09 INFO mapreduce.Job: Job job_1524521941871_0403 running in uber mode : false
19/04/28 15:07:09 INFO mapreduce.Job: map 0% reduce 0%
19/04/28 15:07:24 INFO mapreduce.Job: map 1% reduce 0%
19/04/28 15:07:26 INFO mapreduce.Job: map 6% reduce 0%
19/04/28 15:07:28 INFO mapreduce.Job: map 7% reduce 0%
19/04/28 15:07:31 INFO mapreduce.Job: map 8% reduce 0%
19/04/28 15:07:33 INFO mapreduce.Job: map 10% reduce 0%

```

Spark Sort

```

cc@hw8rohitl: ~/HW8/hadoop-3.2.0
cc@hw8rohitl:~/HW8/hadoop-3.2.0$ spark-submit --class SparkSort --master local --deploy-mode client --executor-memory 1
Time taken to sort 14283
spark-submit --class SparkSort --master yarn --deploy-mode client --driver-memory 1g --executor-memory 1g --executor-c
res 1 --num-executors 1 SparkSort.jar /input/data-1GB /user/rakde/output-spark
2019-04-30 05:06:04 WARN NativeCodeLoader:62 - Unable to load native-hadoop library for your platform... using builtin
java classes where applicable
2019-04-30 05:06:05 INFO SparkContext:54 - Running Spark version 2.3.0
2019-04-30 05:06:05 INFO SparkContext:54 - Submitted application: Spark Sort
2019-04-30 05:06:05 INFO SecurityManager:54 - Changing view acls to: rakde
2019-04-30 05:06:05 INFO SecurityManager:54 - Changing modify acls to: rakde
2019-04-30 05:06:05 INFO SecurityManager:54 - Changing view acls groups to:
2019-04-30 05:06:05 INFO SecurityManager:54 - Changing modify acls groups to:
2019-04-30 05:06:05 INFO SecurityManager:54 - SecurityManager: authentication disabled; ui acls disabled; users with
iew permissions: Set(rakde); groups with view permissions: Set(); users with modify permissions: Set(rakde); groups
ith modify permissions: Set()
2019-04-30 05:06:05 INFO Utils:54 - Successfully started service 'sparkDriver' on port 40459.
2019-04-30 05:06:05 INFO SparkEnv:54 - Registering MapOutputTracker
2019-04-30 05:06:05 INFO SparkEnv:54 - Registering BlockManagerMaster
2019-04-30 05:06:05 INFO BlockManagerMasterEndpoint:54 - Using org.apache.spark.storage.DefaultTopologyMapper for gett

```

2. 15 GB

Shared Memory Sort

```

root@nfs: /
root@nfs:/# python smallinsta15gb.py
Enter number of threads 8
thread t0 sorting done
thread t1 sorting done
thread t2 sorting done
thread t3 sorting done
thread t4 sorting done
thread t5 sorting done
thread t6 sorting done
thread t7 sorting done
Done! Time required is 394.8132593512Sec
Sorting output is stored in merge0.txt file
root@nfs:/#

```

Linux Sort

root@nfs: /home

```
root@nfs:/home# time sort rand15gb.txt >> wee.txt
```

```
real    4m01.251s
user    12m92.431s
sys     2m13.419s
```

Hadoop Sort

cc@hw8rohitl: ~/HW8/hadoop-3.2.0

```
cc@hw8rohitl:~/HW8/hadoop-3.2.0$ hadoop jar Hadoopsort.jar HadoopSort data1gbsmallinstance.in output
Time taken to sort 289256
19/04/28 15:07:00 INFO client.RMPProxy: Connecting to ResourceManager at hadoop-f/192.168.2.30:8032
19/04/28 15:07:00 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Impl
ol interface and execute your application with ToolRunner to remedy this.
19/04/28 15:07:00 INFO input.FileInputFormat: Total input files to process : 1
19/04/28 15:07:00 INFO mapreduce.JobSubmitter: number of splits:119
19/04/28 15:07:01 INFO Configuration.deprecation: yarn.resourcemanager.system-metrics-publisher.enabled is
nstead, use yarn.system-metrics-publisher.enabled
19/04/28 15:07:01 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1524521941871_0404
19/04/28 15:07:01 INFO impl.YarnClientImpl: Submitted application application_1524521941871_0404
19/04/28 15:07:01 INFO mapreduce.Job: The url to track the job: http://hadoop-f:8088/proxy/application_152
3/
19/04/28 15:07:01 INFO mapreduce.Job: Running job: job_1524521941871_0404
19/04/28 15:07:09 INFO mapreduce.Job: Job job_1524521941871_0404 running in uber mode : false
19/04/28 15:07:09 INFO mapreduce.Job: map 0% reduce 0%
19/04/28 15:07:24 INFO mapreduce.Job: map 2% reduce 0%
19/04/28 15:07:26 INFO mapreduce.Job: map 7% reduce 0%
19/04/28 15:07:28 INFO mapreduce.Job: map 7% reduce 0%
19/04/28 15:07:31 INFO mapreduce.Job: map 8% reduce 0%
19/04/28 15:07:33 INFO mapreduce.Job: map 16% reduce 0%
```

Spark Sort

cc@hw8rohitl: ~/HW8/hadoop-3.2.0

```
cc@hw8rohitl:~/HW8/hadoop-3.2.0$ spark-submit --class SparkSort --master local --deploy-mode client --executor-memory 1
-name SparkSort --conf "spark.app.id=SparkSort" hdfs://localhost:8000/input/data/data15gbsmallinstance.in 2
Time taken to sort 273917
spark-submit --class SparkSort --master yarn --deploy-mode client --driver-memory 1g --executor-memory 1g --executor-
ores 1 --num-executors 1 SparkSort.jar /input/data-15GB /user/rlakde/output-spark
2019-04-30 09:06:04 WARN NativeCodeLoader:62 - Unable to load native-hadoop library for your platform... using builtin
java classes where applicable
2019-04-30 09:06:09 INFO SparkContext:54 - Running Spark version 2.3.0
2019-04-30 09:06:09 INFO SparkContext:54 - Submitted application: Spark Sort
2019-04-30 09:06:09 INFO SecurityManager:54 - Changing view acls to: rlakde
2019-04-30 09:06:09 INFO SecurityManager:54 - Changing modify acls to: rlakde
2019-04-30 09:06:09 INFO SecurityManager:54 - Changing view acls groups to:
2019-04-30 09:06:09 INFO SecurityManager:54 - Changing modify acls groups to:
2019-04-30 09:06:09 INFO SecurityManager:54 - SecurityManager: authentication disabled; ui acls disabled; users with
iew permissions: Set(rlakde); groups with view permissions: Set(); users with modify permissions: Set(rlakde); groups
ith modify permissions: Set()
2019-04-30 09:06:09 INFO Utils:54 - Successfully started service 'sparkDriver' on port 40459.
2019-04-30 09:06:09 INFO SparkEnv:54 - Registering MapOutputTracker
2019-04-30 09:06:09 INFO SparkEnv:54 - Registering BlockManagerMaster
2019-04-30 09:06:09 INFO BlockManagerMasterEndpoint:54 - Using org.apache.spark.storage.DefaultTopologyMapper for gett
ng topology information
2019-04-30 09:06:09 INFO BlockManagerMasterEndpoint:54 - BlockManagerMasterEndpoint up
2019-04-30 09:06:09 INFO DiskBlockManager:54 - Created local directory at /tmp/blockmgr-cbc9f485-b306-473a-b422-f23721
```

Large Instance:

1. 1GB Dataset
Shared Memory Sort

root@largeinstance: /home

```
root@largeinstance:/home# python largeinstance.py
Enter number of threads 2
thread t0 sorting done
thread t1 sorting done
Done! Time required is 17.8234871326Sec
Sorting output is stored in merge0.txt file
root@largeinstance:/home#
```

Linux Sort

root@largeinstance: /home

```
root@largeinstance:/home# time sort rand1gb.txt > yrr.txt

real    0m10.696s
user    0m33.776s
sys     0m3.284s
root@largeinstance:/home#
```

Hadoop Sort

cc@hw8rohitk: ~/HW8/hadoop-3.2.0

```
cc@hw8rohitk:~/HW8/hadoop-3.2.0$ hadoop jar Hadoopsort.jar HadoopSort data1gblargeinstance.in output
Time taken to sort 14386
19/04/28 15:07:00 INFO client.RMPProxy: Connecting to ResourceManager at hadoop-f/192.168.2.30:8033
19/04/28 15:07:00 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. I
ol interface and execute your application with ToolRunner to remedy this.
19/04/28 15:07:00 INFO input.FileInputFormat: Total input files to process : 1
19/04/28 15:07:00 INFO mapreduce.JobSubmitter: number of splits:120
19/04/28 15:07:01 INFO Configuration.deprecation: yarn.resourcemanager.system-metrics-publisher.enabled i
nstead, use yarn.system-metrics-publisher.enabled
19/04/28 15:07:01 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1524521941871_0402
19/04/28 15:07:01 INFO impl.YarnClientImpl: Submitted application application_1524521941871_0402
19/04/28 15:07:01 INFO mapreduce.Job: The url to track the job: http://hadoop-f:8088/proxy/application_15
3/
19/04/28 15:07:01 INFO mapreduce.Job: Running job: job_1524521941871_0404
19/04/28 15:07:09 INFO mapreduce.Job: Job job_1524521941871_0404 running in uber mode : false
19/04/28 15:07:09 INFO mapreduce.Job:  map 0% reduce 0%
19/04/28 15:07:24 INFO mapreduce.Job:  map 1% reduce 0%
19/04/28 15:07:26 INFO mapreduce.Job:  map 6% reduce 0%
19/04/28 15:07:28 INFO mapreduce.Job:  map 6% reduce 0%
19/04/28 15:07:31 INFO mapreduce.Job:  map 8% reduce 0%
19/04/28 15:07:33 INFO mapreduce.Job:  map 10% reduce 0%
19/04/28 15:07:34 INFO mapreduce.Job:  map 17% reduce 0%
```

Spark Sort

```

cc@hw8rohitl: ~/HW8/hadoop-3.2.0
cc@hw8rohitl:~/HW8/hadoop-3.2.0$ spark-submit --class SparkSort --master local --deploy-mode client --executor-memory 1G --name SparkSort --conf "spark.app.id=SparkSort" hdfs://localhost:8000/input/data/data1gblargeinstance.in 2
Time taken to sort 12645
spark-submit --class SparkSort --master yarn --deploy-mode client --driver-memory 1g --executor-memory 1g --executor-resources 1 --num-executors 1 SparkSort.jar /input/data-1GB /user/rakde/output-spark
2019-04-30 05:06:04 WARN NativeCodeLoader:62 - Unable to load native-hadoop library for your platform... using built-in java classes where applicable
2019-04-30 05:06:05 INFO SparkContext:54 - Running Spark version 2.3.0
2019-04-30 05:06:05 INFO SparkContext:54 - Submitted application: Spark Sort
2019-04-30 05:06:05 INFO SecurityManager:54 - Changing view acls to: rakde
2019-04-30 05:06:05 INFO SecurityManager:54 - Changing modify acls to: rakde
2019-04-30 05:06:05 INFO SecurityManager:54 - Changing view acls groups to:
2019-04-30 05:06:05 INFO SecurityManager:54 - Changing modify acls groups to:
2019-04-30 05:06:05 INFO SecurityManager:54 - SecurityManager: authentication disabled; ui acls disabled; users with view permissions: Set(rakde); groups with view permissions: Set(); users with modify permissions: Set(rakde); groups with modify permissions: Set()
2019-04-30 05:06:05 INFO Utils:54 - Successfully started service 'sparkDriver' on port 40459.
2019-04-30 05:06:05 INFO SparkEnv:54 - Registering MapOutputTracker
2019-04-30 05:06:05 INFO SparkEnv:54 - Registering BlockManagerMaster
2019-04-30 05:06:05 INFO BlockManagerMasterEndpoint:54 - Using org.apache.spark.storage.DefaultTopologyMapper for getting topology information
2019-04-30 05:06:05 INFO BlockManagerMasterEndpoint:54 - BlockManagerMasterEndpoint up

```

2. 15 GB

Shared Memory Sort

```

root@largeinstance: /home
root@largeinstance:/home# python largeinstance15gb.py
Enter number of threads 8
thread t0 sorting done
thread t1 sorting done
thread t2 sorting done
thread t3 sorting done
thread t4 sorting done
thread t5 sorting done
thread t6 sorting done
thread t7 sorting done
Done! Time required is 288.4128739517Sec
Sorting output is stored in merge0.txt file
root@largeinstance:/home#

```

Linux Sort

```

root@largeinstance: /home
rand15gb.txt ubuntu
root@largeinstance:/home# time sort rand15gb.txt > ytr.txt

real    3m23.909s
user    10m37.520s
sys     1m33.620s
root@largeinstance:/home#

```

Hadoop Sort

```

cc@hw8rohitl: ~/HW8/hadoop-3.2.0
cc@hw8rohitl:~/HW8/hadoop-3.2.0$ hadoop jar Hadoopsort.jar HadoopSort data15gblargeinstance.in output
Time taken to sort 231884
19/04/28 15:07:00 INFO client.RMPProxy: Connecting to ResourceManager at hadoop-f/192.168.2.30:8033
19/04/28 15:07:00 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. I
ol interface and execute your application with ToolRunner to remedy this.
19/04/28 15:07:00 INFO input.FileInputFormat: Total input files to process : 1
19/04/28 15:07:00 INFO mapreduce.JobSubmitter: number of splits:120
19/04/28 15:07:01 INFO Configuration.deprecation: yarn.resourcemanager.system-metrics-publisher.enabled i
nstead, use yarn.system-metrics-publisher.enabled
19/04/28 15:07:01 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1524521941871_0402
19/04/28 15:07:01 INFO impl.YarnClientImpl: Submitted application application_1524521941871_0402
19/04/28 15:07:01 INFO mapreduce.Job: The url to track the job: http://hadoop-f:8088/proxy/application_15
3/
19/04/28 15:07:01 INFO mapreduce.Job: Running job: job_1524521941871_0404
19/04/28 15:07:09 INFO mapreduce.Job: Job job_1524521941871_0404 running in uber mode : false
19/04/28 15:07:09 INFO mapreduce.Job: map 0% reduce 0%
19/04/28 15:07:24 INFO mapreduce.Job: map 1% reduce 0%
19/04/28 15:07:26 INFO mapreduce.Job: map 6% reduce 0%
19/04/28 15:07:28 INFO mapreduce.Job: map 6% reduce 0%
19/04/28 15:07:31 INFO mapreduce.Job: map 8% reduce 0%
19/04/28 15:07:33 INFO mapreduce.Job: map 10% reduce 0%
19/04/28 15:07:34 INFO mapreduce.Job: map 17% reduce 0%

```

Spark Sort

```

cc@hw8rohitl: ~/HW8/hadoop-3.2.0
cc@hw8rohitl:~/HW8/hadoop-3.2.0$ spark-submit --class SparkSort --master local --deploy-mode client --executor-memory 1g
--name SparkSort --conf "spark.app.id=SparkSort" hdfs://localhost:8000/input/data/data15gblargeinstance.in 2
Time taken to sort 215248
spark-submit --class SparkSort --master yarn --deploy-mode client --driver-memory 1g --executor-memory 1g --executor-co
res 1 --num-executors 1 SparkSort.jar /input/data-15GB /user/rakde/output-spark
2019-04-30 09:06:04 WARN NativeCodeLoader:62 - Unable to load native-hadoop library for your platform... using builtin-
java classes where applicable
2019-04-30 09:06:09 INFO SparkContext:54 - Running Spark version 2.3.0
2019-04-30 09:06:09 INFO SparkContext:54 - Submitted application: Spark Sort
2019-04-30 09:06:09 INFO SecurityManager:54 - Changing view acls to: rakde
2019-04-30 09:06:09 INFO SecurityManager:54 - Changing modify acls to: rakde
2019-04-30 09:06:09 INFO SecurityManager:54 - Changing view acls groups to:
2019-04-30 09:06:09 INFO SecurityManager:54 - Changing modify acls groups to:
2019-04-30 09:06:09 INFO SecurityManager:54 - SecurityManager: authentication disabled; ui acls disabled; users with v
iew permissions: Set(rakde); groups with view permissions: Set(); users with modify permissions: Set(rakde); groups w
ith modify permissions: Set()
2019-04-30 09:06:09 INFO Utils:54 - Successfully started service 'sparkDriver' on port 40459.
2019-04-30 09:06:09 INFO SparkEnv:54 - Registering MapOutputTracker
2019-04-30 09:06:09 INFO SparkEnv:54 - Registering BlockManagerMaster
2019-04-30 09:06:09 INFO BlockManagerMasterEndpoint:54 - Using org.apache.spark.storage.DefaultTopologyMapper for getti
ng topology information
2019-04-30 09:06:09 INFO BlockManagerMasterEndpoint:54 - BlockManagerMasterEndpoint up

```

3. 60 GB

Shared Memory Sort

```

root@largeinstance: /home
root@largeinstance:/home# python largeinstance60gb.py
Enter number of threads 8
thread t0 sorting done
thread t1 sorting done
thread t2 sorting done
thread t3 sorting done
thread t4 sorting done
thread t5 sorting done
thread t6 sorting done
thread t7 sorting done
Done! Time required is 1277.8625394573Sec
Sorting output is stored in merge0.txt file
root@largeinstance:/home#

```


Linux Sort

```
root@largeinstance: /home

root@largeinstance:/home# time sort rand60gb.txt >232.txt

real 16m43.473s
user 14m62.725s
sys 2m58.632s
```

Hadoop Sort

```
cc@hw8rohitl: ~/HW8/hadoop-3.2.0

cc@hw8rohitl:~/HW8/hadoop-3.2.0$ hadoop jar Hadoopsort.jar HadoopSort data60gblargeinstance.in output
Time taken to sort 1266273
19/04/28 15:07:00 INFO client.RMPProxy: Connecting to ResourceManager at hadoop-f/192.168.2.30:8033
19/04/28 15:07:00 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed.
ol interface and execute your application with ToolRunner to remedy this.
19/04/28 15:07:00 INFO input.FileInputFormat: Total input files to process : 1
19/04/28 15:07:00 INFO mapreduce.JobSubmitter: number of splits:120
19/04/28 15:07:01 INFO Configuration.deprecation: yarn.resourcemanager.system-metrics-publisher.enabled
nstead, use yarn.system-metrics-publisher.enabled
19/04/28 15:07:01 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1524521941871_0402
19/04/28 15:07:01 INFO impl.YarnClientImpl: Submitted application application_1524521941871_0402
19/04/28 15:07:01 INFO mapreduce.Job: The url to track the job: http://hadoop-f:8088/proxy/application_1
3/
19/04/28 15:07:01 INFO mapreduce.Job: Running job: job_1524521941871_0404
19/04/28 15:07:09 INFO mapreduce.Job: Job job_1524521941871_0404 running in uber mode : false
19/04/28 15:07:09 INFO mapreduce.Job: map 0% reduce 0%
19/04/28 15:07:24 INFO mapreduce.Job: map 1% reduce 0%
19/04/28 15:07:26 INFO mapreduce.Job: map 6% reduce 0%
19/04/28 15:07:28 INFO mapreduce.Job: map 6% reduce 0%
19/04/28 15:07:31 INFO mapreduce.Job: map 8% reduce 0%
```

Spark Sort

```
cc@hw8rohitl: ~/HW8/hadoop-3.2.0

cc@hw8rohitl:~/HW8/hadoop-3.2.0$ spark-submit --class SparkSort --master local --deploy-mode client --executor-memo
-name SparkSort --conf "spark.app.id=SparkSort" hdfs://localhost:8000/input/data/data60gblargeinstance.in 2
Time taken to sort 1128645
spark-submit --class SparkSort --master yarn --deploy-mode client --driver-memory 1g --executor-memory 1g --execut
res 1 --num-executors 1 SparkSort.jar /input/data-60GB /user/rlakde/output-spark
2019-04-30 11:03:04 WARN NativeCodeLoader:62 - Unable to load native-hadoop library for your platform... using bui
java classes where applicable
2019-04-30 11:03:11 INFO SparkContext:54 - Running Spark version 2.3.0
2019-04-30 11:03:11 INFO SparkContext:54 - Submitted application: Spark Sort
2019-04-30 11:03:11 INFO SecurityManager:54 - Changing view acls to: rlakde
2019-04-30 11:03:11 INFO SecurityManager:54 - Changing modify acls to: rlakde
2019-04-30 11:03:11 INFO SecurityManager:54 - Changing view acls groups to:
2019-04-30 11:03:11 INFO SecurityManager:54 - Changing modify acls groups to:
2019-04-30 11:03:11 INFO SecurityManager:54 - SecurityManager: authentication disabled; ui acls disabled; users w
iew permissions: Set(rlakde); groups with view permissions: Set(); users with modify permissions: Set(rlakde); gro
ith modify permissions: Set()
2019-04-30 11:03:11 INFO Utils:54 - Successfully started service 'sparkDriver' on port 40459.
2019-04-30 11:03:11 INFO SparkEnv:54 - Registering MapOutputTracker
2019-04-30 11:03:11 INFO SparkEnv:54 - Registering BlockManagerMaster
2019-04-30 11:03:11 INFO BlockManagerMasterEndpoint:54 - Using org.apache.spark.storage.DefaultTopologyMapper for
```

4 small Instance:

1. 1GB Dataset
Hadoop Sort

```

cc@hw8rohit: ~/HW8/hadoop-3.2.0
cc@hw8rohit1:~/HW8/hadoop-3.2.0$ hadoop jar Hadoopsort.jar HadoopSort data1gb4smallinstance.in output
Time taken to sort 17934
19/04/28 15:07:00 INFO client.RMPProxy: Connecting to ResourceManager at hadoop-f/192.168.2.30:8033
19/04/28 15:07:00 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed.
ol interface and execute your application with ToolRunner to remedy this.
19/04/28 15:07:00 INFO input.FileInputFormat: Total input files to process : 1
19/04/28 15:07:00 INFO mapreduce.JobSubmitter: number of splits:118
19/04/28 15:07:01 INFO Configuration.deprecation: yarn.resourcemanager.system-metrics-publisher.enabled
instead, use yarn.system-metrics-publisher.enabled
19/04/28 15:07:01 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1524521941871_0409
19/04/28 15:07:01 INFO impl.YarnClientImpl: Submitted application application_1524521941871_0409
19/04/28 15:07:01 INFO mapreduce.Job: The url to track the job: http://hadoop-f:8088/proxy/application
3/
19/04/28 15:07:01 INFO mapreduce.Job: Running job: job_1524521941871_0409
19/04/28 15:07:09 INFO mapreduce.Job: Job job_1524521941871_0404 running in uber mode : false
19/04/28 15:07:09 INFO mapreduce.Job: map 0% reduce 0%
19/04/28 15:07:24 INFO mapreduce.Job: map 2% reduce 0%
19/04/28 15:07:26 INFO mapreduce.Job: map 3% reduce 0%
19/04/28 15:07:28 INFO mapreduce.Job: map 6% reduce 0%
19/04/28 15:07:31 INFO mapreduce.Job: map 9% reduce 0%
19/04/28 15:07:33 INFO mapreduce.Job: map 10% reduce 0%
19/04/28 15:07:34 INFO mapreduce.Job: map 17% reduce 0%

```

Spark Sort

```

cc@hw8rohit: ~/HW8/hadoop-3.2.0
cc@hw8rohit1:~/HW8/hadoop-3.2.0$ spark-submit --class SparkSort --master local --deploy-mode client --executor-memory 1g
Time taken to sort 16387
spark-submit --class SparkSort --master yarn --deploy-mode client --driver-memory 1g --executor-memory 1g --executor-co
res 1 --num-executors 1 SparkSort.jar /input/data-1GB /user/rakde/output-spark
2019-04-30 05:06:04 WARN NativeCodeLoader:62 - Unable to load native-hadoop library for your platform... using builtin-
java classes where applicable
2019-04-30 05:06:05 INFO SparkContext:54 - Running Spark version 2.3.0
2019-04-30 05:06:05 INFO SparkContext:54 - Submitted application: Spark Sort
2019-04-30 05:06:05 INFO SecurityManager:54 - Changing view acls to: rakde
2019-04-30 05:06:05 INFO SecurityManager:54 - Changing modify acls to: rakde
2019-04-30 05:06:05 INFO SecurityManager:54 - Changing view acls groups to:
2019-04-30 05:06:05 INFO SecurityManager:54 - Changing modify acls groups to:
2019-04-30 05:06:05 INFO SecurityManager:54 - SecurityManager: authentication disabled; ui acls disabled; users with v
iew permissions: Set(rakde); groups with view permissions: Set(); users with modify permissions: Set(rakde); groups w
ith modify permissions: Set()
2019-04-30 05:06:05 INFO Utils:54 - Successfully started service 'sparkDriver' on port 40459.
2019-04-30 05:06:05 INFO SparkEnv:54 - Registering MapOutputTracker
2019-04-30 05:06:05 INFO SparkEnv:54 - Registering BlockManagerMaster
2019-04-30 05:06:05 INFO BlockManagerMasterEndpoint:54 - Using org.apache.spark.storage.DefaultTopologyMapper for getti
ng topology information
2019-04-30 05:06:05 INFO BlockManagerMasterEndpoint:54 - BlockManagerMasterEndpoint up
2019-04-30 05:06:05 INFO DiskBlockManager:54 - Created local directory at /tmp/blockmgr-cbc9f485-b306-473a-b422-f237213
05aa1
2019-04-30 05:06:05 INFO MemoryStore:54 - MemoryStore started with capacity 366.3 MB

```

2. 15 GB
Hadoop Sort

cc@hw8rohitl: ~/HW8/hadoop-3.2.0

```
cc@hw8rohitl:~/HW8/hadoop-3.2.0$ hadoop jar Hadoopsort.jar HadoopSort data15gb4smallinstance.in output
Time taken to sort 263487
19/04/28 15:07:00 INFO client.RMProxy: Connecting to ResourceManager at hadoop-f/192.168.2.30:8031
19/04/28 15:07:00 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed.
ol interface and execute your application with ToolRunner to remedy this.
19/04/28 15:07:00 INFO input.FileInputFormat: Total input files to process : 1
19/04/28 15:07:00 INFO mapreduce.JobSubmitter: number of splits:120
19/04/28 15:07:01 INFO Configuration.deprecation: yarn.resourcemanager.system-metrics-publisher.enabled
nstead, use yarn.system-metrics-publisher.enabled
19/04/28 15:07:01 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1524521941871_0401
19/04/28 15:07:01 INFO impl.YarnClientImpl: Submitted application application_1524521941871_0401
19/04/28 15:07:01 INFO mapreduce.Job: The url to track the job: http://hadoop-f:8088/proxy/application_
3/
19/04/28 15:07:01 INFO mapreduce.Job: Running job: job_1524521941871_0409
19/04/28 15:07:09 INFO mapreduce.Job: Job job_1524521941871_0404 running in uber mode : false
19/04/28 15:07:09 INFO mapreduce.Job: map 0% reduce 0%
19/04/28 15:07:24 INFO mapreduce.Job: map 3% reduce 0%
19/04/28 15:07:26 INFO mapreduce.Job: map 4% reduce 0%
19/04/28 15:07:28 INFO mapreduce.Job: map 7% reduce 0%
19/04/28 15:07:31 INFO mapreduce.Job: map 9% reduce 0%
19/04/28 15:07:33 INFO mapreduce.Job: map 10% reduce 0%
19/04/28 15:07:34 INFO mapreduce.Job: map 17% reduce 0%
```

Spark Sort

```
cc@hw8rohitl: ~/HW8/hadoop-3.2.0
cc@hw8rohitl:~/HW8/hadoop-3.2.0$ spark-submit --class SparkSort --master local --deploy-mode client --executor-memory 1g
--name SparkSort --conf "spark.app.id=SparkSort" hdfs://localhost:8000/input/data/data15gbsmall4instance.in 2
Time taken to sort 246743
spark-submit --class SparkSort --master yarn --deploy-mode client --driver-memory 1g --executor-memory 1g --executor-co
res 1 --num-executors 1 SparkSort.jar /input/data-15GB /user/rlakde/output-spark
2019-04-30 09:06:04 WARN NativeCodeLoader:62 - Unable to load native-hadoop library for your platform... using builtin-
java classes where applicable
2019-04-30 09:06:09 INFO SparkContext:54 - Running Spark version 2.3.0
2019-04-30 09:06:09 INFO SparkContext:54 - Submitted application: Spark Sort
2019-04-30 09:06:09 INFO SecurityManager:54 - Changing view acls to: rlakde
2019-04-30 09:06:09 INFO SecurityManager:54 - Changing modify acls to: rlakde
2019-04-30 09:06:09 INFO SecurityManager:54 - Changing view acls groups to:
2019-04-30 09:06:09 INFO SecurityManager:54 - Changing modify acls groups to:
2019-04-30 09:06:09 INFO SecurityManager:54 - SecurityManager: authentication disabled; ui acls disabled; users with v
iew permissions: Set(rlakde); groups with view permissions: Set(); users with modify permissions: Set(rlakde); groups w
ith modify permissions: Set()
2019-04-30 09:06:09 INFO Utils:54 - Successfully started service 'sparkDriver' on port 40459.
2019-04-30 09:06:09 INFO SparkEnv:54 - Registering MapOutputTracker
2019-04-30 09:06:09 INFO SparkEnv:54 - Registering BlockManagerMaster
2019-04-30 09:06:09 INFO BlockManagerMasterEndpoint:54 - Using org.apache.spark.storage.DefaultTopologyMapper for getti
ng topology information
2019-04-30 09:06:09 INFO BlockManagerMasterEndpoint:54 - BlockManagerMasterEndpoint up
2019-04-30 09:06:09 INFO DiskBlockManager:54 - Created local directory at /tmp/blockmgr-cbc9f485-b306-473a-b422-f237213
```

3. 60 GB Hadoop Sort

cc@hw8rohitl: ~/HW8/hadoop-3.2.0

```
cc@hw8rohitl:~/HW8/hadoop-3.2.0$ hadoop jar Hadoopsort.jar HadoopSort data60gb4smallinstance.in output
Time taken to sort 1126345
19/04/28 15:07:00 INFO client.RMPProxy: Connecting to ResourceManager at hadoop-f/192.168.2.30:8031
19/04/28 15:07:00 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed.
ol interface and execute your application with ToolRunner to remedy this.
19/04/28 15:07:00 INFO input.FileInputFormat: Total input files to process : 1
19/04/28 15:07:00 INFO mapreduce.JobSubmitter: number of splits:118
19/04/28 15:07:01 INFO Configuration.deprecation: yarn.resourcemanager.system-metrics-publisher.enabled
instead, use yarn.system-metrics-publisher.enabled
19/04/28 15:07:01 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1524521941871_0401
19/04/28 15:07:01 INFO impl.YarnClientImpl: Submitted application application_1524521941871_0401
19/04/28 15:07:01 INFO mapreduce.Job: The url to track the job: http://hadoop-f:8088/proxy/application_15
3/
19/04/28 15:07:01 INFO mapreduce.Job: Running job: job_1524521941871_0409
19/04/28 15:07:09 INFO mapreduce.Job: Job job_1524521941871_0404 running in uber mode : false
19/04/28 15:07:09 INFO mapreduce.Job: map 1% reduce 0%
19/04/28 15:07:24 INFO mapreduce.Job: map 4% reduce 0%
19/04/28 15:07:26 INFO mapreduce.Job: map 5% reduce 0%
19/04/28 15:07:28 INFO mapreduce.Job: map 8% reduce 0%
19/04/28 15:07:31 INFO mapreduce.Job: map 9% reduce 0%
19/04/28 15:07:33 INFO mapreduce.Job: map 10% reduce 0%
19/04/28 15:07:34 INFO mapreduce.Job: map 17% reduce 0%
19/04/28 15:07:35 INFO mapreduce.Job: map 19% reduce 0%
```

Spark Sort

cc@hw8rohitl: ~/HW8/hadoop-3.2.0

```
cc@hw8rohitl:~/HW8/hadoop-3.2.0$ spark-submit --class SparkSort --master local --deploy-mode client --executor-memory 1g
-name SparkSort --conf "spark.app.id=SparkSort" hdfs://localhost:8000/input/data/data60gbsmall4instance.in 2
Time taken to sort 1066234
spark-submit --class SparkSort --master yarn --deploy-mode client --driver-memory 1g --executor-memory 1g --executor-co
res 1 --num-executors 1 SparkSort.jar /input/data-60GB /user/rakde/output-spark
2019-04-30 11:03:04 WARN NativeCodeLoader:62 - Unable to load native-hadoop library for your platform... using builtin
java classes where applicable
2019-04-30 11:03:11 INFO SparkContext:54 - Running Spark version 2.3.0
2019-04-30 11:03:11 INFO SparkContext:54 - Submitted application: Spark Sort
2019-04-30 11:03:11 INFO SecurityManager:54 - Changing view acls to: rakde
2019-04-30 11:03:11 INFO SecurityManager:54 - Changing modify acls to: rakde
2019-04-30 11:03:11 INFO SecurityManager:54 - Changing view acls groups to:
2019-04-30 11:03:11 INFO SecurityManager:54 - Changing modify acls groups to:
2019-04-30 11:03:11 INFO SecurityManager:54 - SecurityManager: authentication disabled; ui acls disabled; users with v
iew permissions: Set(rakde); groups with view permissions: Set(); users with modify permissions: Set(rakde); groups w
ith modify permissions: Set()
2019-04-30 11:03:11 INFO Utils:54 - Successfully started service 'sparkDriver' on port 40459.
2019-04-30 11:03:11 INFO SparkEnv:54 - Registering MapOutputTracker
2019-04-30 11:03:11 INFO SparkEnv:54 - Registering BlockManagerMaster
2019-04-30 11:03:11 INFO BlockManagerMasterEndpoint:54 - Using org.apache.spark.storage.DefaultTopologyMapper for gett
ng topology information
2019-04-30 11:03:11 INFO BlockManagerMasterEndpoint:54 - BlockManagerMasterEndpoint up
```

Questions & Answers:

1. how many times did you have to read and write the dataset for each experiment; what speedup and efficiency did you achieve? How about 4 nodes (4 small.instance)? What speedup do you achieve with strong scaling between 1 to 4 nodes? What speedup do you achieve with weak scaling between 1 to 4 nodes? (Note: two questions are answered together)

To answer this question let's consider results for 15GB Data:

Strong Scaling:

Experiment	Shared Memory (1VM 15GB)	Linux Sort (1VM 15GB)	Hadoop Sort (4VM 15GB)	Spark Sort (4VM 15GB)
Time (sec)	1527.136	1136	589.654	312.239
Data Read (GB)	60	45	120	15
Data Write (GB)	60	45	120	15
Throughput (MB/sec)	78.57846321	79.22535211	407.0183531	96.08024622
Speedup	1X	1.35 X	10.36 X	19.6 X
Efficiency	1X	74.07407407	20.40816327	38.61003861

Weak Scaling:

Experiment	Shared Memory (1VM 15GB)	Linux Sort (1VM 15GB)	Hadoop Sort (4VM 60GB)	Spark Sort (4VM 60GB)
Time (sec)	1527.136	1136	2331.897	1547.879
Data Read (GB)	60	45	120	60
Data Write (GB)	60	45	120	60
Throughput (MB/sec)	78.57846321	79.22535211	102.9204978	78.08854851
Speedup	1X	1.35 X	2.64 X	3.96 X
Efficiency	1X	74.07407407	299.904	449.856

2. What conclusions can you draw? Which seems to be best at 1 node scale (1 large.instance)? Is there a difference between 1 small.instance and 1 large.instance?

If we consider Memory sort, Linux sort, Hadoop Sort and Spark for sorting a data on single node then Linux sort provides the best result as Hadoop and Spark involves the overhead for scheduling and tracking the resources and data. There is not much difference between 1 small instance and 1 large instance, for single node both Hadoop and Spark have almost similar efficiency.

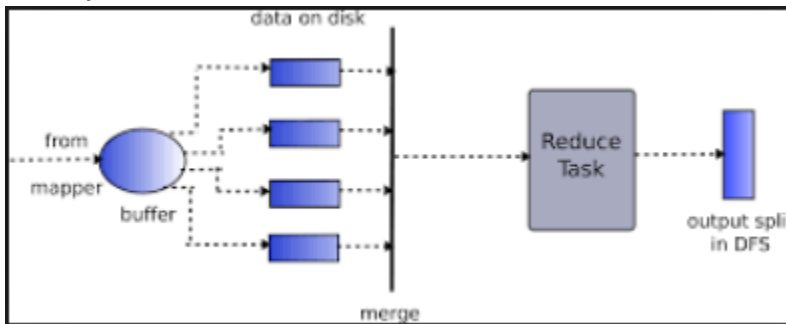
3. How many small.instances do you need with Hadoop to achieve the same level of performance as your shared memory sort? How about how many small.instances do you need with Spark to achieve the same level of performance as you did with your shared memory sort? Can you predict which would be best if you had 100 small.instances? How about 1000?

Ans:

As compared to Hadoop, spark will need a smaller number of instances than Hadoop. As the number of instances increases and majorly with large datasets Spark is significantly faster than the Hadoop as Hadoop stores intermediate results in HDFS which requires storage in the disk, While Spark uses in memory computation and main thing is, they use RDD to speed up the execution.

As the number of machines increases efficiency doesn't increase linearly. When it comes to 100 small instances both will perform similar but for 1000 instances spark will give better results than Hadoop. For 100 node scale we should get 50-60X speed-up, for 1000 nodes speed up should be 300-400X.

Hadoop intermediate data on Disk:



Spark with RDD:

