## Single-node Hadoop Setup

After installation, we successfully access the port 50070 of localhost.





We try to run the Terasort. First, we generate a list of random numbers.

20/09/12 08:09:30 INFO terasort.TeraGen: Generating 100000 using 2
20/09/12 08:09:30 INFO mapreduce.JobSubmitter: number of splits:2
20/09/12 08:09:30 INFO Configuration.deprecation: yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use
yarn.system-metrics-publisher.enabled
20/09/12 08:09:31 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1599898159880_0001
20/09/12 08:09:32 INFO impl.YarnClientImpl: Submitted application application_1599898159880_0001
20/09/12 08:09:32 INFO mapreduce.Job: The url to track the job: http://instance-1:8088/proxy/application_1599898159880_0001/
20/09/12 08:09:32 INFO mapreduce.Job: Running job: job_1599898159880_0001
20/09/12 08:09:40 INFO mapreduce.Job: Job job_1599898159880_0001 running in uber mode : false
20/09/12 08:09:40 INFO mapreduce.Job:  map 0% reduce 0%
20/09/12 08:09:49 INFO mapreduce.Job:  map 50% reduce 0%
20/09/12 08:09:50 INFO mapreduce.Job:  map 100% reduce 0%
20/09/12 08:09:50 INFO mapreduce.Job: Job job_1599898159880_0001 completed successfully
20/09/12 08:09:50 INFO mapreduce.Job: Counters: 31
        File System Counters
                FILE: Number of bytes read=0
                FILE: Number of bytes written=396300
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=164
                HDFS: Number of bytes written=10000000
                HDFS: Number of read operations=8
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=4
        Job Counters
                Launched map tasks=2
                Other local map tasks=2
                Total time spent by all maps in occupied slots (ms)=13215
                Total time spent by all reduces in occupied slots (ms)=0
                Total time spent by all map tasks (ms)=13215
                Total vcore-milliseconds taken by all map tasks=13215
                Total megabyte-milliseconds taken by all map tasks=13532160
        Map-Reduce Framework
                Map input records=100000
                Map output records=100000
                Input split bytes=164
                Spilled Records=0
                Failed Shuffles=0
                Merged Map outputs=0
                GC time elapsed (ms)=369
                CPU time spent (ms)=2650
                Physical memory (bytes) snapshot=389361664
                Virtual memory (bytes) snapshot=3854913536
                Total committed heap usage (bytes)=2317735296
        org.apache.hadoop.examples.terasort.TeraGen$Counters
                CHECKSUM=2145749851290000
        File Input Format Counters
                Bytes Read=0
        File Output Format Counters
                Bytes Written=10000000
[jessechan@instance-1 hadoop]$

Then, Sorting the data.

Instead use the hdfs command for it.

ls: 'terasort/check': No such file or directory
[jessechan@instance-1 hadoop]$ bin/hadoop jar share/hadoop/mapreduce/hadoop-mapreduce-examples-2.9.2.jar terasort terasort/input terasort/output
20/09/12 08:32:02 INFO terasort.TeraSort: starting
20/09/12 08:32:03 INFO input.FileInputFormat: Total input files to process : 2
Spent 126ms computing base-splits.
Spent 3ms computing TeraScheduler splits.
Computing input splits took 130ms
Sampling 2 splits of 2
Making 1 from 100000 sampled records
Computing parititions took 684ms
Spent 817ms computing partitions.
20/09/12 08:32:03 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
20/09/12 08:32:04 WARN hdfs.DataStreamer: Caught exception
java.lang.InterruptedException
        at java.lang.Object.wait(Native Method)
        at java.lang.Thread.join(Thread.java:1252)
        at java.lang.Thread.join(Thread.java:1326)
        at org.apache.hadoop.hdfs.DataStreamer.closeResponder(DataStreamer.java:980)
        at org.apache.hadoop.hdfs.DataStreamer.endBlock(DataStreamer.java:630)
        at org.apache.hadoop.hdfs.DataStreamer.run(DataStreamer.java:807)
20/09/12 08:32:04 INFO mapreduce.JobSubmitter: number of splits:2
20/09/12 08:32:04 INFO Configuration.deprecation: yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enable
d
20/09/12 08:32:04 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1599898159880_0006
20/09/12 08:32:05 INFO impl.YarnClientImpl: Submitted application application_1599898159880_0006
20/09/12 08:32:05 INFO mapreduce.Job: The url to track the job: http://instance-1:8088/proxy/application_1599898159880_0006/
20/09/12 08:32:05 INFO mapreduce.Job: Running job: job_1599898159880_0006
20/09/12 08:32:12 INFO mapreduce.Job: Job job_1599898159880_0006 running in uber mode : false
20/09/12 08:32:12 INFO mapreduce.Job:  map 0% reduce 0%
20/09/12 08:32:22 INFO mapreduce.Job:  map 100% reduce 0%
20/09/12 08:32:29 INFO mapreduce.Job:  map 100% reduce 100%
20/09/12 08:32:30 INFO mapreduce.Job: Job job_1599898159880_0006 completed successfully
20/09/12 08:32:30 INFO mapreduce.Job: Counters: 50
        File System Counters
                FILE: Number of bytes read=10400006
                FILE: Number of bytes written=21398682
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=10000258
                HDFS: Number of bytes written=10000000
                HDFS: Number of read operations=9
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=2
        Job Counters
                Killed map tasks=1
                Launched map tasks=2
                Launched reduce tasks=1
                Data-local map tasks=2
                Total time spent by all maps in occupied slots (ms)=15092

Finish and now checking,

Screenshot 1 (PuTTY 34.92.40.164):

```
                Total time spent by all reduces in occupied slots (ms)=4791
                Total time spent by all map tasks (ms)=15092
                Total time spent by all reduce tasks (ms)=4791
                Total vcore-milliseconds taken by all map tasks=15092
                Total vcore-milliseconds taken by all reduce tasks=4791
                Total megabyte-milliseconds taken by all map tasks=15454208
                Total megabyte-milliseconds taken by all reduce tasks=4905984
        Map-Reduce Framework
                Map input records=100000
                Map output records=100000
                Map output bytes=10200000
                Map output materialized bytes=10400012
                Input split bytes=258
                Combine input records=0
                Combine output records=0
                Reduce input groups=100000
                Reduce shuffle bytes=10400012
                Reduce input records=100000
                Reduce output records=100000
                Spilled Records=200000
                Shuffled Maps =2
                Failed Shuffles=0
                Merged Map outputs=2
                GC time elapsed (ms)=456
                CPU time spent (ms)=6120
                Physical memory (bytes) snapshot=795328512
                Virtual memory (bytes) snapshot=5778935808
                Total committed heap usage (bytes)=523763712
        Shuffle Errors
                BAD_ID=0
                CONNECTION=0
                IO_ERROR=0
                WRONG_LENGTH=0
                WRONG_MAP=0
                WRONG_REDUCE=0
        File Input Format Counters
                Bytes Read=10000000
        File Output Format Counters
                Bytes Written=10000000
20/09/12 08:32:30 INFO terasort.TeraSort: done
[jessechan@instance-1 hadoop]$ bin/hadoop jar share/hadoop/mapreduce/hadoop-mapreduce-examples-2.9.2.jar teravalidate terasort/output terasort/check
20/09/12 08:32:50 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
20/09/12 08:32:51 INFO input.FileInputFormat: Total input files to process : 1
Spent 17ms computing base-splits.
Spent 3ms computing TeraScheduler splits.
20/09/12 08:32:51 INFO mapreduce.JobSubmitter: number of splits:1
20/09/12 08:32:51 INFO Configuration.deprecation: yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enable
d
20/09/12 08:32:51 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1599898159880_0007
20/09/12 08:32:51 INFO impl.YarnClientImpl: Submitted application application_1599898159880_0007
20/09/12 08:32:52 INFO mapreduce.Job: The url to track the job: http://instance-1:8088/proxy/application_1599898159880_0007/
20/09/12 08:32:52 INFO mapreduce.Job: Running job: job_1599898159880_0007
```

Screenshot 2 (PuTTY 34.92.40.164):

```
                Bytes Written=10000000
20/09/12 08:32:30 INFO terasort.TeraSort: done
[jessechan@instance-1 hadoop]$ bin/hadoop jar share/hadoop/mapreduce/hadoop-mapreduce-examples-2.9.2.jar teravalidate terasort/output terasort/check
20/09/12 08:32:50 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
20/09/12 08:32:51 INFO input.FileInputFormat: Total input files to process : 1
Spent 17ms computing base-splits.
Spent 3ms computing TeraScheduler splits.
20/09/12 08:32:51 INFO mapreduce.JobSubmitter: number of splits:1
20/09/12 08:32:51 INFO Configuration.deprecation: yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enable
d
20/09/12 08:32:51 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1599898159880_0007
20/09/12 08:32:51 INFO impl.YarnClientImpl: Submitted application application_1599898159880_0007
20/09/12 08:32:52 INFO mapreduce.Job: The url to track the job: http://instance-1:8088/proxy/application_1599898159880_0007/
20/09/12 08:32:52 INFO mapreduce.Job: Running job: job_1599898159880_0007
20/09/12 08:32:59 INFO mapreduce.Job: Job job_1599898159880_0007 running in uber mode : false
20/09/12 08:32:59 INFO mapreduce.Job:  map 0% reduce 0%
20/09/12 08:33:05 INFO mapreduce.Job:  map 100% reduce 0%
20/09/12 08:33:10 INFO mapreduce.Job:  map 100% reduce 100%
20/09/12 08:33:10 INFO mapreduce.Job: Job job_1599898159880_0007 completed successfully
20/09/12 08:33:10 INFO mapreduce.Job: Counters: 49
        File System Counters
                FILE: Number of bytes read=92
                FILE: Number of bytes written=397313
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=10000130
                HDFS: Number of bytes written=22
                HDFS: Number of read operations=6
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=2
        Job Counters
                Launched map tasks=1
                Launched reduce tasks=1
                Data-local map tasks=1
                Total time spent by all maps in occupied slots (ms)=3519
                Total time spent by all reduces in occupied slots (ms)=3150
                Total time spent by all map tasks (ms)=3519
                Total time spent by all reduce tasks (ms)=3150
                Total vcore-milliseconds taken by all map tasks=3519
                Total vcore-milliseconds taken by all reduce tasks=3150
                Total megabyte-milliseconds taken by all map tasks=3603456
                Total megabyte-milliseconds taken by all reduce tasks=3225600
        Map-Reduce Framework
                Map input records=100000
                Map output records=3
                Map output bytes=80
                Map output materialized bytes=92
                Input split bytes=130
                Combine input records=0
                Combine output records=0
                Reduce input groups=3
```

Use command to validate the sorting result,



Screenshot 3:

```
                WRONG_LENGTH=0
                WRONG_MAP=0
                WRONG_REDUCE=0
        File Input Format Counters
                Bytes Read=10000000
        File Output Format Counters
                Bytes Written=22
[jessechan@instance-1 hadoop]$ ./bin/hadoop dfs -ls terasort/check
DEPRECATED: Use of this script to execute hdfs command is deprecated.
Instead use the hdfs command for it.

Found 2 items
-rw-r--r--   3 jessechan supergroup          0 2020-09-12 08:33 terasort/check/_SUCCESS
-rw-r--r--   3 jessechan supergroup         22 2020-09-12 08:33 terasort/check/part-r-00000
[jessechan@instance-1 hadoop]$
```
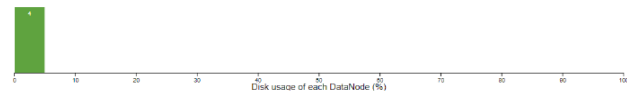
## Multi-node Hadoop Cluster Setup

Using One name-node as master and 3 data-node as slave

## Datanode Information



In the 2GB TeraSort example:



Figure 1 2GB Gen.

Logged in as: dr.who

**hadoop**

# Application application_1600064201081_0002

**Application Overview**

| | |
|---|---|
| User: | hduser |
| Name: | TeraSort |
| Application Type: | MAPREDUCE |
| Application Tags: | |
| Application Priority: | 0 (Higher Integer value indicates higher priority) |
| YarnApplicationState: | FINISHED |
| Queue: | default |
| FinalStatus Reported by AM: | SUCCEEDED |
| Started: | Mon Sep 14 06:19:58 +0000 2020 |
| Elapsed: | 1mins, 42sec |
| Tracking URL: | History |
| Log Aggregation Status: | DISABLED |
| Application Timeout (Remaining Time): | Unlimited |
| Diagnostics: | |
| Unmanaged Application: | false |
| Application Node Label expression: | <Not set> |
| AM container Node Label expression: | <DEFAULT_PARTITION> |

**Application Metrics**

| | |
|---|---|
| Total Resource Preempted: | <memory:0, vCores:0> |
| Total Number of Non-AM Containers Preempted: | 0 |
| Total Number of AM Containers Preempted: | 0 |
| Resource Preempted from Current Attempt: | <memory:0, vCores:0> |
| Number of Non-AM Containers Preempted from Current Attempt: | 0 |
| Aggregate Resource Allocation: | 904307 MB-seconds, 766 vcore-seconds |
| Aggregate Preempted Resource Allocation: | 0 MB-seconds, 0 vcore-seconds |

Show 20 entries    Search:

| Attempt ID | Started | Node | Logs | Nodes blacklisted by the app | Nodes blacklisted by the system |
|---|---|---|---|---|---|
| appattempt_1600064201081_0002_000001 | Mon Sep 14 14:19:58 +0800 2020 | http://datanode1-1.asia-east2-a.c.driven-plexus-289109.internal:8042 | Logs | 0 | 0 |

Showing 1 to 1 of 1 entries    First Previous 1 Next Last

Figure 2 2GB Sort

Logged in as: dr.who

**hadoop**

# Application application_1600064201081_0003

**Application Overview**

| | |
|---|---|
| User: | hduser |
| Name: | TeraValidate |
| Application Type: | MAPREDUCE |
| Application Tags: | |
| Application Priority: | 0 (Higher Integer value indicates higher priority) |
| YarnApplicationState: | FINISHED |
| Queue: | default |
| FinalStatus Reported by AM: | SUCCEEDED |
| Started: | Mon Sep 14 06:22:37 +0000 2020 |
| Elapsed: | 30sec |
| Tracking URL: | History |
| Log Aggregation Status: | DISABLED |
| Application Timeout (Remaining Time): | Unlimited |
| Diagnostics: | |
| Unmanaged Application: | false |
| Application Node Label expression: | <Not set> |
| AM container Node Label expression: | <DEFAULT_PARTITION> |

**Application Metrics**

| | |
|---|---|
| Total Resource Preempted: | <memory:0, vCores:0> |
| Total Number of Non-AM Containers Preempted: | 0 |
| Total Number of AM Containers Preempted: | 0 |
| Resource Preempted from Current Attempt: | <memory:0, vCores:0> |
| Number of Non-AM Containers Preempted from Current Attempt: | 0 |
| Aggregate Resource Allocation: | 96798 MB-seconds, 57 vcore-seconds |
| Aggregate Preempted Resource Allocation: | 0 MB-seconds, 0 vcore-seconds |

Show 20 entries    Search:

| Attempt ID | Started | Node | Logs | Nodes blacklisted by the app | Nodes blacklisted by the system |
|---|---|---|---|---|---|
| appattempt_1600064201081_0003_000001 | Mon Sep 14 14:22:37 +0800 2020 | http://datanode3-1.asia-east2-a.c.driven-plexus-289109.internal:8042 | Logs | 0 | 0 |

Showing 1 to 1 of 1 entries    First Previous 1 Next Last

Figure 3 2GB Check

In the 20GB TeraSort example:



Figure 4 20GB Gen.



Figure 5 20GB Sort.

Figure 6 20GB Check

As conclusion, The 20GB TeraSort program take much longer time than 2GB TeraSort program due to the data size. The gen., sort. and check process are similar to above, below figures attached as references.



Running the Python Code on Hadoop

The result of the Python Wordcount Script.

Figure 7 Python Wordcount Program

## Compiling the Java WordCount program for MapReduce

The result is identical to the above python program.



Figure 8 Java Wordcount Program

As a conclusion, the java version of Wordcount program run faster than python program. I guess it due to Java is less dynamic than Python and It makes it more efficient on VM.