



(P.9)

File Management in Hadoop and Benchmarking.

Ans →

1. File Management Task:-

- Adding files and directories.
- Retrieving files
- Deleting files.

2. Benchmarking & stress Testing on Apache Hadoop cluster.

We will be performing these task on Hadoop running in Pseudo-distributed mode on a Windows environment, using WSL. Make sure Hadoop is already installed & configured.

Part(i): Implementing File Management Task in Hadoop

1. Adding Files and Directories to HDFS:-

To add files and directories to the Hadoop Distributed file System (HDFS). We use the hdfs fs command.

• Create a Directory in HDFS: To create a directory in HDFS, use the following command

• `hadoop fs-mkdir [user]/hadoop/new-directory`

• Copy a local file to HDFS:- if you want to add file from your local file system to HDFS. Use the -copyFromLocal command

• `hadoop fs-copyFromLocal /path/to/local/file.txt /user/hadoop/new-directory/`



2. Retrieving Files from HDFS:-

You can retrieve files from HDFS bucket to your local file system using the -get or -copyTo command.

• Retrieve a file from HDFS:-

- hdfs fs -get /user/hadoop/new-directory /path/to/local/destination/

• List the contents of a directory in HDFS:-

To list the contents of a directory in HDFS, you can use the -ls command

• hdfs fs -ls /user/hadoop/new-directory

3. Deleting Files & Directories from HDFS

To delete file and directories from HDFS, use the -rm and -rf flags.

• Delete a file from HDFS:

-> hdfs fs -rm /user/hadoop/new-directory/file.txt

• Delete a Directory & all its Contents:

-> hdfs fs -rm -r /user/hadoop/new-directory



Part (ii): Benchmarking and Stress Testing an Apache Hadoop Cluster

1. Benchmarking Hadoop Using the terasort Benchmark
→ Apache Hadoop comes with a built-in benchmark called Terasort, which is a sorting benchmark that tests the scalability and performance of your Hadoop cluster.

Step 1: Generate Data for the Benchmark

- You need to generate some test data before running the benchmark.
- first, create a directory to store the test data:
 - hadoop fs -mkdir /user/hadoop/terasort-input
 - hadoop jar \$HADOOP_HOME/share/hadoop/mapreduce/hadoop-mapreduce-examples-* jar terasort 10000000 /user/hadoop/terasort output

Step 2: Run the Benchmark

- The terasort command runs the sorting process and place the output in /user/hadoop/terasort output.
- hadoop jar \$HADOOP_HOME/share/hadoop/mapreduce/hadoop-mapreduce-examples-* jar terasort /user/hadoop/terasort-input /user/hadoop/terasort-output



Step 3: Monitor the Benchmarking Job

1. Monitor via Web UI:

- YARN Resource Manager: `http://localhost:8088` to monitor job progress & performance.
- HDFS NameNode web UI: `http://localhost:50070` to monitor HDFS usage and health.

2. Job Completion:

- Once the job completes, you will see the sorted output in the `/user/hadoop/terysort_output` directory.
- You can verify the output running:
 - `hadoop fs -ls /user/hadoop/terysort_output`

2. Stress Testing Hadoop cluster:

Stress testing is important to evaluate the stability and performance of the Hadoop cluster under heavy load. You can use multiple methods to stress Hadoop, such as running multiple MapReduce jobs, adding large datasets, etc.



Step 1: Run Multiple Concurrent jobs.

- Generate multiple datasets:

- hadoop jar \$HADOOP_HOME/share/hadoop/mapreduce/hadoop-mapreduce-examples*.jar teragen 10000000/user/hadoop/stress-input-1
- hadoop jar \$HADOOP_HOME/share/hadoop/mapreduce/hadoop-mapreduce-examples*.jar teragen 10000000/user/hadoop/stress-input-2
- hadoop jar \$HADOOP_HOME/share/hadoop/mapreduce/hadoop-mapreduce-examples*.jar teragen 10000000/user/hadoop/stress-input-3

- Run multiple sort operation on each dataset concurrently.

- hadoop jar \$Hadoop_HOME/share/hadoop/mapreduce/hadoop-mapreduce-examples*.jar terasort /user/hadoop/stress-input-1 /user/hadoop/stress-output-1 &
- hadoop jar \$Hadoop_HOME/share/hadoop/mapreduce/hadoop-mapreduce-examples*.jar terasort /user/hadoop/stress-input-2 /user/hadoop/stress-output-2 &
- hadoop jar \$Hadoop_HOME/share/hadoop/mapreduce/hadoop-mapreduce-examples*.jar terasort /user/hadoop/stress-input-3 /user/hadoop/stress-output-3 &



→ Step 2: Monitor system performance

- Use the Yarn Resource Manager UI: `http://localhost:8088` to monitor the job status and resource consumption (CPU)
- Check the HDFS NameNode UI: `http://localhost:50070` for information on HDFS health and block usage.

→ Track the Hadoop cluster's Health & Performance

Step 3:

1. HDFS Web UI: Check the usage, file system health, and block description.
2. Resource Manager: View the status of running job, CPU, memory usage, and overall job completion statistics.

→ Step 4: Evaluate the Results:-

- After completing the benchmark & stress test job, check the Hadoop logs to assess the overall performance & identify bottlenecks. Hadoop logs can be found in the following directory:
- `/usr/local/hadoop/logs`
- The logs can help you identify areas where the system may be underperforming or experiencing issues, such as slow disk I/O, network bottlenecks or insufficient memory.



→ Conclusion :-

by following the above steps, you have successfully performed the following tasks:

1. Added files and directories to HDFS.
2. Retrieved and deleted files from HDFS.
3. Run benchmarking tasks using the Terasort example to assess Hadoop's performance.
4. Conducted a stress test by running multiple jobs concurrently and monitoring system performance.