

Lab 12: Map Reduce using Hadoop

Instructor: Vivek Sarkar

Resource Summary

Course wiki: <https://wiki.rice.edu/confluence/display/PARPROG/COMP322>

Staff Email: comp322-staff@mailman.rice.edu

Clear Login: ssh *your-netid*@ssh.clear.rice.edu and then login with your password

Sugar Login: ssh *your-netid*@sugar.rice.edu and then login with your password

Linux Tutorial visit <http://www.rcsg.rice.edu/tutorials/>

IMPORTANT: Please refer to the tutorial on Linux and SUGAR from Lab 5, as needed. Also, if you edit files on a PC or laptop, be sure to transfer them to SUGAR before you compile and execute them (otherwise you may compile and execute a stale/old version on SUGAR).

As in past labs, create a text file named lab_12_written.txt in the lab_12 directory, and enter your timings and observations there.

1 Map and Reduce Operations on Sets of Key-Value Pairs

Map Reduce is a simple data processing paradigm introduced in Lecture 35. To process a data set, the user just needs to provide implementations of map and reduce functions. A single stage of a map-reduce computation typically consists of a map phase, a shuffle phase and a reduce phase. The input and output of a Map Reduce computation is represented as (key, value) pairs. (The input and output pairs can be of different types.) Hadoop Map Reduce is an open source implementation of the MapReduce paradigm originally proposed by Google. It is currently the go to platform for many Big Data applications.

2 Hadoop Map Reduce Environment Setup

1. Copy the tar file to your home directory using the following command

```
"cp /home/yz17/comp322-lab12.tar ."
```

2. Run the following command in the lab_12/ directory to set up the environment for executing mpiJava programs, "tar -xvf comp322-lab12.tar".

You should now see a directory named "hadoop-1.0.3". You will work in this directory for the rest of the lab.

3 Write a Hadoop Map Reduce Program

1. Study the WordCount.java file, and fill in the code for the map() and reduce() functions as indicated in the comments labeled TODO.
2. Compile the program by using the command "./compileWordCount.sh". The script will compile the Java file into class files and package them into a jar file. You should see a jar file named WordCount.jar after successfully running the script.

4 Run Hadoop Map Reduce in Standalone mode

1. Obtain a Sugar compute node using the `qsub` command as in past labs.
2. Go to the `hadoop-1.0.3` directory.
3. Change your distribution to Standalone mode using the following command
“`./changeToStandAlone.sh`”
4. Run the `WordCount` program in Standalone mode using the following command:
“`bin/hadoop jar WordCount.jar WordCount inputFiles output`”
In Standalone mode, the Hadoop FileSystem is not started. It is a mode often used for debugging.
5. Check your output from your word count program with “`cat output/*`”

5 Run Hadoop Map Reduce in Pseudo Distributed mode

In this section, we will run `wordcount` in a “Pseudo Distributed” mode. This will create a full fledged Hadoop MapReduce system with multiple processes on a single Sugar node. We will demonstrate utilization of multi-cores in this section.

1. First, convert your Hadoop MapReduce system from Standalone mode to Pseudo Distributed Mode
“`./changeToPD.sh`”
2. Run your `wordcount` program using the following script “`./runWordCountPD.sh`”
3. Record the running time output from the script
Hadoop utilizes multi-cores in a machine by having the user specify a parameter in the configuration file. The `mapred.tasktracker.map.tasks.maximum` property determines how many map tasks can run in parallel in the system, thereby using multiple cores.
4. Open “`conf/mapred-site.xml`” using any text editor, and change the value of the property “`mapred.tasktracker.map.tasks`” to 2
5. Run your `wordcount` program using the following script and record the running time of the program
“`./runWordCountPD.sh`”
This running time is for executing two map tasks in parallel.
6. Repeat the process and set the value of the parameter to be 4, 8 and record the running times in `lab_12_written.txt`.
7. Take a few minutes to look at the `runWordCountPD.sh` file. It shows how you work with the Hadoop File System. Hadoop Map Reduce is built on top of HDFS.

6 Turning in your lab work

1. Check in the `WordCount.java` and `lab_12_written.txt` files. You should have recorded running time results in the textfile.
2. Use the turn-in script to submit the contents of the `lab_12.zip` file as a new `lab_12` directory in your turnin repository as explained in Lab 1. You can always examine the most recent contents of your svn repository by visiting <https://svn.rice.edu/r/comp322/turnin/S13/your-netid>.