# Report: Experiment on Cloud Computing

Victrid

*dept. of Computer Science and Engineering*
*Shanghai Jiao Tong University*
Shanghai, China
github.com/Victrid

## I. Introduction

As computer technology continues to cross into various fields, the resulting huge amount of data makes it challenging to analyze and process the data. Traditional single computer programs have difficulty processing these data, while mainframe clusters are expensive and require complex operations and controls, making them limited to large analytics companies.

Today, cloud computing technology is reaching into every corner of the world. As the superiority of controlled cost, elastic management, and rapid migration is gradually being demonstrated, even ordinary people are able to access these resources and take advantage of them. The development of distributed technology has also made it possible to process huge amounts of data on inexpensive computer arrays. Today's cloud computing centers have become an important part of the Internet's infrastructure and will occupy an even more important position in the future.

This report on Cloud Computing Experiment contains these parts:

- Create VM cluster and configure Hadoop and Spark framework on Huawei Cloud.
- Run The Examples from Hadoop and Spark.
- Solve actual problems with GraphX API with Spark.

## II. Configuration of Virtual Machine Cluster

### A. Mirror Configuration

Fast provisioning and delivery, automation, and ease of scaling are key to today's elastic computing. However, setting up instances one by one, configuring their cumbersome dependencies and environment variables, and installing them by downloading source code and binary packages from the web and configuring them by copy and paste is typical of UNIX mainframe administrators in the 1980s.

Compared to the handbook [1] suggests that using an Ubuntu image and configure them manually, we've built a custom mirror, which already has Hadoop, Spark, and SBT packed up and modified. We referenced the image creation tool provided by Arch Linux on [2] with the Cloud-init tool and rewrote the compiling script of the Hadoop package on [3] to match the `JAVA_HOME` lookup process. A system installation script was written to suit the needs of this experiment. This series of configuration files are placed in the submitted files' `image-buildscript` folder.

Many cloud providers has provided custom image service, so do Huawei Cloud. By utilizing its IMS services [4] as figure 1, we can now create virtual machines with our prebuilt images.



Fig. 1: Huawei Cloud IMS

### B. Virtual Machine Configuration

To perform VM configuration as designed in our mirror, adjusting roles while the VM is created, we use Huawei Cloud SDK and its OpenAPI to perform creation. The core part is to utilize OpenStack `user_data`. When the `user_data` is provided as bash scripts, it will be run at VM creation by cloud-init. The script we created performs Hadoop configuration according to Hadoop Documentation [5], worker appointment, and SSH key configuration. This series of configuration scripts and API call Python scripts are placed in the submitted files' `HuaweiCloud-openAPI` folder.

### C. Hadoop Initiation

By running the commands below: (Our Hadoop instance is installed under `/usr/lib`)

```
cd /usr/lib/hadoop
bin/hdfs namenode -format
sbin/start-all.sh
```

The Hadoop instance will be set up, and can be viewed via `http://master:9870` as in figure 2.

Fig. 2: DataNode Information on NameNode

## III. HADOOP EXAMPLE: WORDCOUNT

We've prepared a dummy text file, containing *lorem ipsum*, an industry standard dummy text in the printing and typesetting since 1500s. The text file contains 150 paragraphs, 13547 words, and is used to test.

The text file is transmitted to the master node via `scp` at `/root/lorem.txt`.

by running the commands below (at `/root` folder):

```
hadoop fs -mkdir /input
hadoop fs -put lorem.txt /input
hadoop jar /usr/lib/hadoop/share/hadoop/\
mapreduce/\
hadoop-mapreduce-examples-3.3.1.jar \
wordcount /input /output
hadoop fs -cat /output/part-r-00000
```

The running results are shown in figure 3. Both input file and results are put in the `WordCount` folder.



Fig. 3: Word Count Results

By running the commands below, we can clear the input files and results. This will be helpful for further experiments.

```
hadoop fs -rm -f -r /output
hadoop fs -rm -f -r /input
```

## IV. GRAPHICX EXAMPLE: CONNECTED COMPONENT

We've copied the needed file and organized as [6]. The package needs more dependencies than WordCount program, and the scala used by Spark 3.2.0 should be 2.12.15. We modified the simple build tools script as below:

```
#    [Trailed for PDF typesetting]
scalaVersion := "2.12.15"

libraryDependencies+="org.apache.spark" \
    %% "spark-core" % "3.2.0"
libraryDependencies+="org.apache.spark" \
    %% "spark-sql" % "3.2.0"
libraryDependencies+="org.apache.spark" \
    %% "spark-graphx" % "3.2.0"
```

After compiling as in figure 4b, we upload these 2 graph txt files to data/graphx folder on HDFS, and call the class by `org.apache.spark.examples.graphx.ConnectedComponentsExample` as in figure 4c and received the result as in figure 4d:

```
(justinbieber,1)
(matei_zaharia,3)
(ladygaga,1)
(BarackObama,1)
(jeresig,3)
(odersky,3)
```

The necessary files and source code are put in the `GraphX` folder.

## V. PAGERANK ALGORITHM

Here we use the Wikipedia Vote Network dataset [7] as our processing source. The requested PageRank algorithm is similar to the *Connected Component* one, and after checking the spark source code, we're sure that the `GraphLoader` can also be used to process the wiki votes.

PageRank algorithm is named after both the term "web page" and Google co-founder Larry Page. The key to the ranking is a probabilistic balance between nodes. At the beginning of the computational process the pagerank for each node is randomized, and for each iteration, the rank is computed as:

$$PR(p_i) = \frac{1-d}{N} + d \sum_{p_j \in M(p_i)} \frac{PR(p_j)}{L(p_j)} \qquad (1)$$

An approximation of real page rank value can be calculated with several iterations.

After researching on GraphX Programming Guide [8], the graph build with GraphLoader has implemented PageRank algorithm. In our implementation, it is called by

```
val ranksGraph = graph.pageRank(0.0001)
```

The original source file does not contain name information, but we want to form a more intuitive result as *Connected Component*, which adds an additional user–node linkage. We used `sed` scripts to pre-process the original vote data from [9], and generate the user–node list as in `users.txt`.

```
sed '/^[\s]*#/d;
/^\s*$/d;
/^[ET].*/d;
/^N\t-1\tUNKNOWN.*/d;
s/^[UN]\t\(.*\)\t\(.*\)$/\1,\2/g;
s/^V\t.*\t\(.*\)\t.*\t\(.*\)$/\1,\2/g' \
original.txt | sort -g | uniq > users.txt
```

The user–node list generated would be like:

```
3,ludraman
4,gzornenplatz
5,orthogonal
6,andrevan
7,texture
8,lst27
9,mirv
...
```

After compiling and submitting like we've done in the *Connected Component* part, the results are shown below and as figure 4e:

```
PR: 32.78, ID:4037, Name:elonka
PR: 26.18, ID:  15, Name:danny
PR: 25.52, ID:6634, Name:tenpoundhammer
PR: 23.36, ID:2625, Name:_clown_will_eat_me
PR: 18.56, ID:2398, Name:werdna
PR: 17.96, ID:2470, Name:alex_bakharev
PR: 17.76, ID:2237, Name:khoikhoi
PR: 16.14, ID:4191, Name:ryulong
PR: 15.44, ID:7553, Name:dihydrogen_monoxide
PR: 15.30, ID:5254, Name:gracenotes
PR: 14.51, ID:2328, Name:phaedriel
PR: 14.48, ID:1186, Name:william_m._connolley
PR: 13.84, ID:1297, Name:robchurch
PR: 13.78, ID:4335, Name:mer-c
PR: 13.75, ID:7620, Name:cobi
PR: 13.65, ID:5412, Name:protectionbot
PR: 13.57, ID:7632, Name:redirectcleanupbot
PR: 13.33, ID:4875, Name:earle_martin
PR: 12.87, ID:6946, Name:useight
PR: 12.69, ID:3352, Name:crzrussian
```

Although we do not know how Wikipedia administrators are selected, after checking their usernames, those we checked were all engaged in Wikipedia administration during 2008. This shows that distributed computing is not just an airy idea in papers or academics, but can also work well in solving practical problems.

All necessary files and source code are put in the `PageRank` folder.

## VI. CONCLUSION

*Omitted.*

## REFERENCES

[1] C. Li *et al.*, "Cloud computing course experiment handbook," 2021.
[2] K. Klausen *et al.*, "Arch linux cloud image build script," 2022. [Online]. Available: https://gitlab.archlinux.org/archlinux/arch-boxes
[3] C. Severance *et al.*, "Arch linux packaging script for hadoop," 2021. [Online]. Available: https://aur.archlinux.org/packages/hadoop
[4] huaweicloud.com, "Ims documentation," 2021. [Online]. Available: https://support.huaweicloud.com/ims/index.html
[5] Apache Software Foundation, "Hadoop documentation: Hadoop cluster setup," 2021. [Online]. Available: https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-common/ClusterSetup.html
[6] Z. Lin, "Wordcount tutorial," 2017. [Online]. Available: http://dblab.xmu.edu.cn/blog/1311-2/
[7] J. Leskovec, D. Huttenlocher, and J. Kleinberg, "Signed networks in social media," in *Proceedings of the SIGCHI conference on human factors in computing systems*, 2010, pp. 1361–1370. [Online]. Available: http://snap.stanford.edu/data/wiki-Vote.html
[8] Apache Software Foundation, "Graphx programming guide," 2021. [Online]. Available: https://spark.apache.org/docs/latest/graphx-programming-guide.html
[9] J. Leskovec, D. Huttenlocher, and J. Kleinberg, "Signed networks in social media," in *Proceedings of the SIGCHI conference on human factors in computing systems*, 2010, pp. 1361–1370. [Online]. Available: http://snap.stanford.edu/data/wiki-Elec.html

(a) Building Images



(b) SBT compilation



(c) Spark Submission



(d) Spark GraphX results



(e) Page Rank results

Fig. 4: Supplementary Figures