## BÁO CÁO LAB2 LƯU TRỮ XỬ LÝ DỮ LIỆU LỚN

Nhóm: Squad Game

Thành viên: Lại Ngọc Thăng Long  20183581

Nguyễn Đình Dũng     20183506

Nguyễn Thành Long   20183586

Nguyễn Khương Duy 20183513

BÀI LÀM:

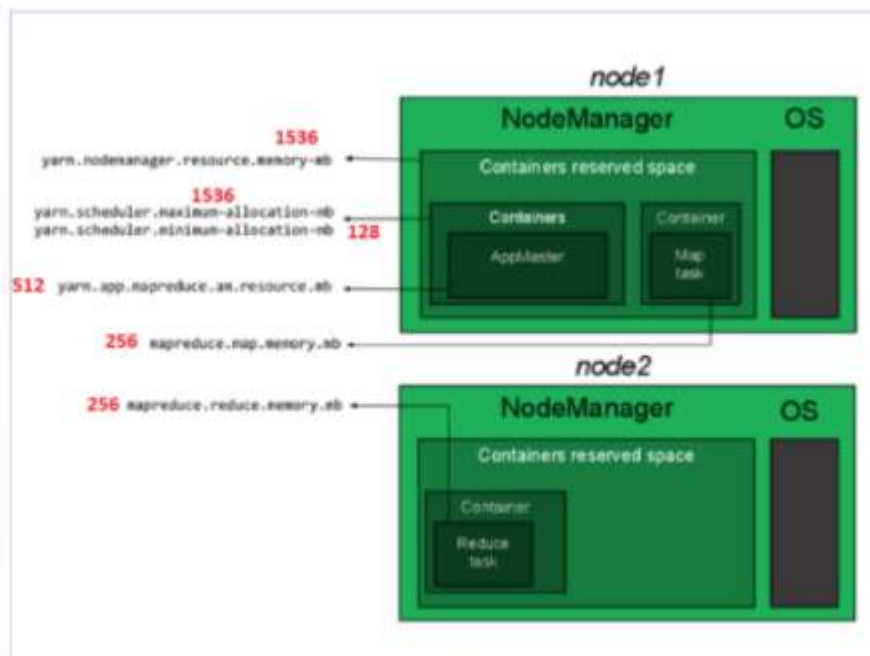# I.    Trên 3 máy master, slave1, slave2



Figure 8: Sample config for 2GB Nodes (will change in next lab)

## 1st Step: Configure yarn

```
export HADOOP_HOME="/usr/local/hadoop"
export HADOOP_COMMON_HOME=$HADOOP_HOME
export HADOOP_CONF_DIR=$HADOOP_HOME/etc/hadoop
export HADOOP_HDFS_HOME=$HADOOP_HOME
export HADOOP_MAPRED_HOME=$HADOOP_HOME
export HADOOP_YARN_HOME=$HADOOP_HOME
```

```
hadoopuser@hadoop-master:~$ export HADOOP_HOME="/usr/local/hadoop"
hadoopuser@hadoop-master:~$ export HADOOP_COMMON_HOME=$HADOOP_HOME
hadoopuser@hadoop-master:~$ export HADOOP_CONF_DIR=$HADOOP_HOME/etc/hadoop
hadoopuser@hadoop-master:~$ export HADOOP_HDFS_HOME=$HADOOP_HOME
hadoopuser@hadoop-master:~$ export HADOOP_MAPRED_HOME=$HADOOP_HOME
hadoopuser@hadoop-master:~$ export HADOOP_YARN_HOME=$HADOOP_HOME
hadoopuser@hadoop-master:~$
```

# 2nd Step: Configure mapred-site.xml

```
sudo nano /usr/local/hadoop/etc/hadoop/mapred-site.xml
```

```
<configuration>

<property>

        <name>mapreduce.framework.name</name>

        <value>yarn</value>

   </property>

   <property>

        <name>yarn.app.mapreduce.am.env</name>


   <value>HADOOP_MAPRED_HOME=$HADOOP_HOME</value>

   </property>

   <property>

        <name>mapreduce.map.env</name>


   <value>HADOOP_MAPRED_HOME=$HADOOP_HOME</value>

   </property>
```

```xml
    <property>
            <name>mapreduce.reduce.env</name>
    <value>HADOOP_MAPRED_HOME=$HADOOP_HOME</value>
    </property>

<property>
        <name>yarn.app.mapreduce.am.resource.mb</name>
        <value>512</value>
</property>

<property>
        <name>mapreduce.map.memory.mb</name>
        <value>256</value>
</property>

<property>
        <name>mapreduce.reduce.memory.mb</name>
        <value>256</value>
</property>
</configuration>
```

```
  GNU nano 4.8          /usr/local/hadoop/etc/hadoop/mapred-site.xml
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
<!--
  Licensed under the Apache License, Version 2.0 (the "License");
  you may not use this file except in compliance with the License.
  You may obtain a copy of the License at

    http://www.apache.org/licenses/LICENSE-2.0

  Unless required by applicable law or agreed to in writing, software
  distributed under the License is distributed on an "AS IS" BASIS,
  WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
  See the License for the specific language governing permissions and
  limitations under the License. See accompanying LICENSE file.
-->

<!-- Put site-specific property overrides in this file. -->

<configuration>
<property>
            <name>mapreduce.framework.name</name>
            <value>yarn</value>
    </property>
    <property>
            <name>yarn.app.mapreduce.am.env</name>
            <value>HADOOP_MAPRED_HOME=$HADOOP_HOME</value>
    </property>
    <property>
            <name>mapreduce.map.env</name>
            <value>HADOOP_MAPRED_HOME=$HADOOP_HOME</value>
    </property>
    <property>
            <name>mapreduce.reduce.env</name>
            <value>HADOOP_MAPRED_HOME=$HADOOP_HOME</value>
    </property>

<property>
        <name>yarn.app.mapreduce.am.resource.mb</name>
        <value>512</value>
</property>

<property>
        <name>mapreduce.map.memory.mb</name>
        <value>256</value>
</property>

<property>
        <name>mapreduce.reduce.memory.mb</name>
        <value>256</value>
</property>
</configuration>
```

```
^G Get Help   ^O Write Out  ^W Where Is   ^K Cut Text   ^J Justify    ^C Cur Pos
^X Exit       ^R Read File   ^\ Replace    ^U Paste Text ^T To Spell   ^  Go To Line
```

# 3rd Step: Config yarn-site.xml

```
sudo nano /usr/local/hadoop/etc/hadoop/yarn-site.xml


<configuration>


<!-- Site specific YARN configuration properties -->
<property>

        <name>yarn.acl.enable</name>

        <value>0</value>

  </property>
<property>
<name>yarn.resourcemanager.hostname</name>
<value>hadoop-master</value>
</property>
 <property>

        <name>yarn.nodemanager.aux-services</name>

        <value>mapreduce_shuffle</value>

  </property>


<property>

      <name>yarn.nodemanager.resource.memory-mb</name>

      <value>1536</value>
</property>
```

```xml
<property>
        <name>yarn.scheduler.maximum-allocation-mb</name>
        <value>1536</value>
</property>

<property>
        <name>yarn.scheduler.minimum-allocation-mb</name>
        <value>128</value>
</property>

<property>
        <name>yarn.nodemanager.vmem-check-enabled</name>
        <value>false</value>
</property>
</configuration>
```

```
GNU nano 4.8          /usr/local/hadoop/etc/hadoop/yarn-site.xml
<?xml version="1.0"?>
<!--
  Licensed under the Apache License, Version 2.0 (the "License");
  you may not use this file except in compliance with the License.
  You may obtain a copy of the License at

    http://www.apache.org/licenses/LICENSE-2.0

  Unless required by applicable law or agreed to in writing, software
  distributed under the License is distributed on an "AS IS" BASIS,
  WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
  See the License for the specific language governing permissions and
  limitations under the License. See accompanying LICENSE file.
-->
<configuration>

<!-- Site specific YARN configuration properties -->
<property>
            <name>yarn.acl.enable</name>
            <value>0</value>
    </property>
<property>
<name>yarn.resourcemanager.hostname</name>
<value>hadoop-master</value>
</property>
 <property>
            <name>yarn.nodemanager.aux-services</name>
            <value>mapreduce_shuffle</value>
    </property>

<property>
        <name>yarn.nodemanager.resource.memory-mb</name>
        <value>1536</value>
</property>

<property>
        <name>yarn.scheduler.maximum-allocation-mb</name>
        <value>1536</value>
</property>

<property>
        <name>yarn.scheduler.minimum-allocation-mb</name>
        <value>128</value>
</property>

<property>
        <name>yarn.nodemanager.vmem-check-enabled</name>
        <value>false</value>
</property>
</configuration>
                        [ Read 50 lines ]
^G Get Help  ^O Write Out ^W Where Is  ^K Cut Text  ^J Justify   ^C Cur Pos
^X Exit      ^R Read File ^\ Replace   ^U Paste Text^T To Spell  ^  Go To Line
```

## II.    Trên máy master

## 1st Step: Start yarn

```
start-yarn.sh
```



Lên web, truy cập [hadoop-master:8088/cluster](hadoop-master:8088/cluster)



# 2nd Step: Tạo MapReduce Job

Chuẩn bị data



-     Updata lên HDFS  folder /WordCountLab02/Input

Browse Directory

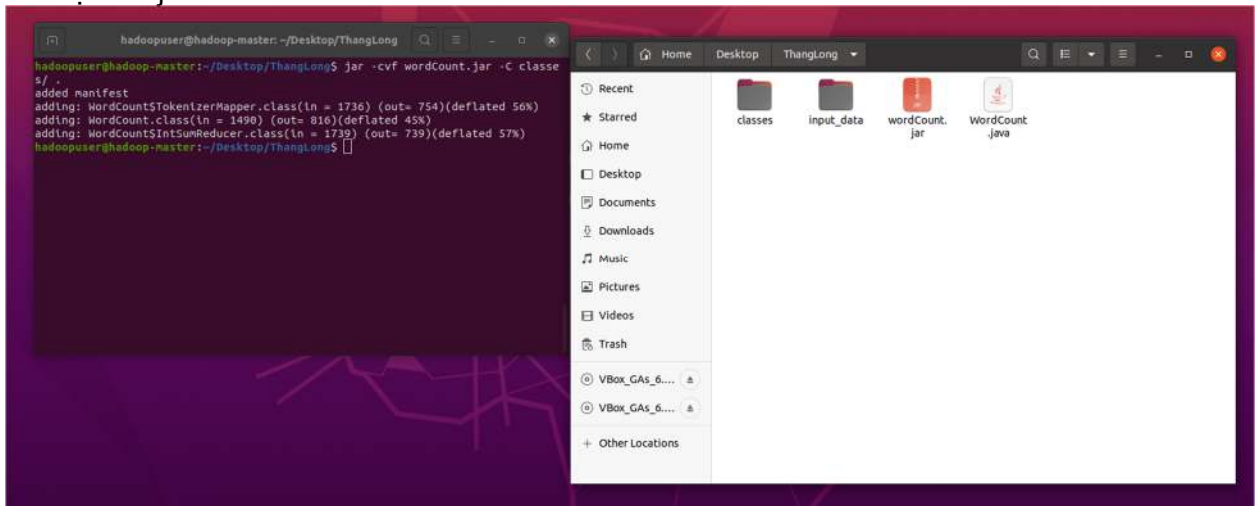| | Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name | |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | -rw-r--r-- | hadoopuser | supergroup | 201.52 MB | Oct 21 22:19 | 2 | 128 MB | purchases_2012.txt | 🗑 |
| ☐ | -rw-r--r-- | hadoopuser | supergroup | 201.52 MB | Oct 21 22:19 | 2 | 128 MB | purchases_2013.txt | 🗑 |
| ☐ | -rw-r--r-- | hadoopuser | supergroup | 201.52 MB | Oct 21 22:19 | 2 | 128 MB | purchases_2014.txt | 🗑 |
| ☐ | -rw-r--r-- | hadoopuser | supergroup | 201.52 MB | Oct 21 22:19 | 2 | 128 MB | purchases_2015.txt | 🗑 |
| ☐ | -rw-r--r-- | hadoopuser | supergroup | 201.52 MB | Oct 21 22:19 | 2 | 128 MB | purchases_2016.txt | 🗑 |
| ☐ | -rw-r--r-- | hadoopuser | supergroup | 201.52 MB | Oct 21 22:19 | 2 | 128 MB | purchases_2017.txt | 🗑 |

Showing 1 to 6 of 6 entries

Hadoop, 2019.

- Thực thi file wordcount.java
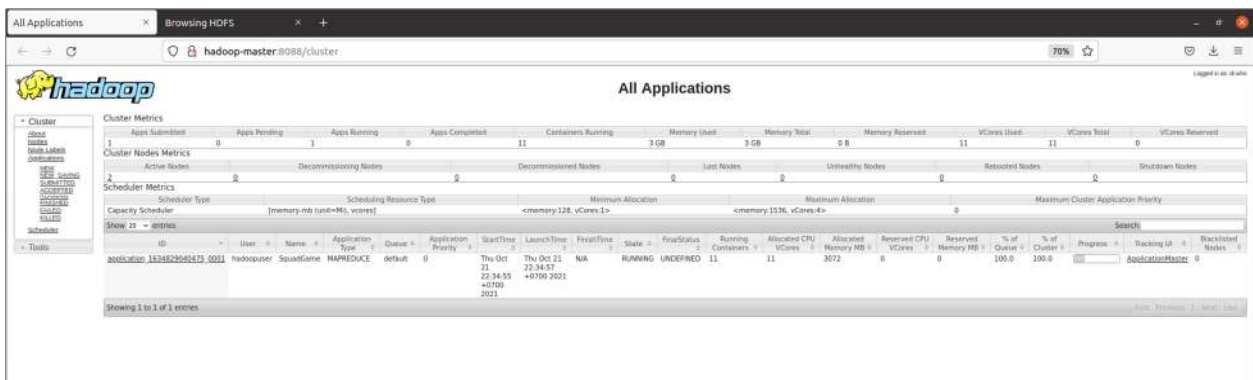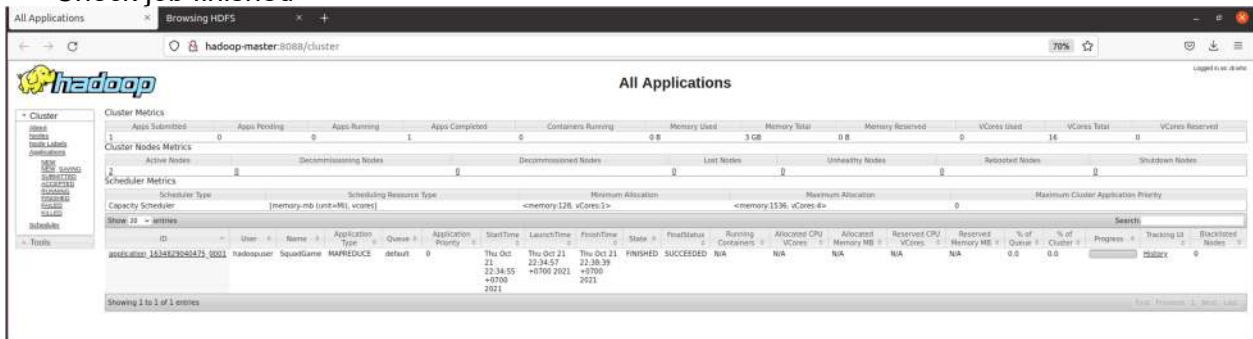- Xuất hiện 3 file class trong folder /classes

- Tạo file jar



- Sử dụng file jar thực thi chạy job mapreduce

- Lên web, check job đang chạy



- Check job finished



- Xem kết quả

Browsing HDFS

hadoop-master:9870/explorer.html#/WordCountLab02/Output

Hadoop    Overview    Datanodes    Datanode Volume Failures    Snapshot    Startup Progress    Utilities ▾

## Browse Directory

/WordCountLab02/Output                                          Go!

Show 25 ∨ entries                                        Search:

| | ↓↑ | Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name | ↓↑ |
|---|---|---|---|---|---|---|---|---|---|---|
| ☐ | | -rw-r--r-- | hadoopuser | supergroup | 0 B | Oct 21 22:38 | 2 | 128 MB | _SUCCESS | 🗑 |
| ☐ | | -rw-r--r-- | hadoopuser | supergroup | 565.36 KB | Oct 21 22:38 | 2 | 128 MB | part-r-00000 | 🗑 |

Showing 1 to 2 of 2 entries                        Previous  1  Next

Hadoop, 2019.

Chạy lệnh:
Hdfs dfs -cat /WordCountLab02/Output/*



```
Orlando 241170
Orleans 239484
Paso    239292
Paul    240960
Pet     1375332
Petersburg      240558
Philadelphia    244488
Phoenix 241998
Pittsburgh      242148
Plano   241020
Portland        240390
Raleigh 241566
Reno    241524
Richmond        239898
Riverside       239778
Rochester       242730
Rouge   242322
Sacramento      243366
Saint   240960
San     1200120
Santa   241836
Scottsdale      241038
Seattle 239196
Spokane 241332
Sporting        1379592
Springs 242334
St.     480450
Stockton        239976
Supplies        1375332
Tampa   240816
Toledo  240834
Toys    1379784
Tucson  239220
Tulsa   241482
Vegas   481068
Video   1381422
Virginia        241014
Visa    4963326
Vista   240480
Washington      243018
Wayne   242634
Wichita 242532
Winston-Salem   241248
Women's 1380300
Worth   242016
York    242184
and     1378002
hadoopuser@hadoop-master:~/Desktop/ThangLong$
```