

## 华东师范大学数据学院上机实践报告

课程名称：分布式模型与编程 年级：2018 上机实践成绩：  
指导教师：徐辰 姓名：孙秋实  
上机实践名称：Hadoop2.x 部署 学号：10185501402 上机实践日期：2021/3/14  
上机实践编号：Lab4 组号：Group5 上机实践时间：

---

### Part 1

#### 实验目的

- (1) 学习 Hadoop2.x 的部署，理解单机集中式、单机伪分布式部署与分布式部署的区别
  - (2) 学会通过查找系统日志中的错误来解决系统部署中遇到的问题
  - (3) 学会使用基本的 HDFS Shell 操作命令
  - (4) 通过系统部署理解 Hadoop2.x 的体系架构，以及 Hadoop1.x 和 Hadoop2.x 之间的差异
- 

### Part 2

#### 实验任务

- (1) 完成 Hadoop2.x 的单机集中式、单机伪分布式部署与分布式部署（分布式部署为小组合作）
  - (2) 在三种部署方式下均能成功运行示例程序
- 

### Part 3

#### 使用环境

- (1) 操作系统：Ubuntu 18.04
  - (2) JDK 版本：1.8
  - (3) Hadoop 版本：2.10.1
- 

### Part 4

#### 实验过程

### Section 1

#### 单机集中式部署

在实验开始前，先登录 dase-local 本地用户，解压 hadoop-2.10.1 并检查 hadoop 的版本，确认为 hadoop-2.10.1

**Task 1****实验：HDFS**

```
dase-local@10-24-21-12:~/softwares/hadoop-2.10.1$ ./bin/hadoop jar ./share/hadoop/mapreduce/hadoop-mapreduce-examples-2.10.1.jar grep ~/input/grep ~/output/grep 'dfs[a-z.]+'
21/03/18 15:53:55 INFO Configuration.deprecation: session.id is deprecated. Instead, use dfs.metrics.session-id
21/03/18 15:53:55 INFO jvm.JvmMetrics: Initializing JVM Metrics with processName=JobTracker, sessionId=
21/03/18 15:53:55 INFO input.FileInputFormat: Total input files to process : 8
21/03/18 15:53:55 INFO mapreduce.JobSubmitter: number of splits:8
21/03/18 15:53:55 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local1615833063_0001
```

图 1: Login and Check Hadoop Version

运行 MapReduce 的应用程序，提交 jar 包任务并且运行 grep 示例

```
dase-local@10-24-21-12:~/softwares$ tar -xvf hadoop-2.10.1.tar.gz
dase-local@10-24-21-12:~/softwares$ ls
flink-1.12.1-bin-scala_2.11.tgz  hadoop-2.10.1.tar.gz
hadoop-1.2.1                      jdk-8u171-linux-x64.tar.gz
hadoop-1.2.1.tar.gz                scala-2.11.12.tgz
hadoop-2.10.1                      spark-2.4.7-bin-without-hadoop.tgz
dase-local@10-24-21-12:~/softwares$ cd hadoop-2.10.1
dase-local@10-24-21-12:~/softwares/hadoop-2.10.1$ ls
bin  include  libexec  NOTICE.txt  sbin
etc  lib     LICENSE.txt  README.txt  share
dase-local@10-24-21-12:~/softwares/hadoop-2.10.1$ ./bin/hadoop version
Hadoop 2.10.1
Subversion https://github.com/apache/hadoop -r 1827467c9a56f133025f28557bfc2c562d78e816
Compiled by centos on 2020-09-14T13:17Z
Compiled with protoc 2.5.0
From source with checksum 3114edef868f1f3824e7d0f68be03650
This command was run using /home/dase-local/softwares/hadoop-2.10.1/share/hadoop/common/hadoop-common-2.10.1.jar
```

图 2: 提交 jar 包，运行 grep 示例

运行结束后检查输出结果

```
dase-local@10-24-21-12:~/softwares/hadoop-2.10.1$ cat ~/output/grep/*
1      dfsadmin
```

图 3: 输出结果

如实验手册所示运行一个 wordcount 示例

运行过程如下所示，另起一个 terminal，使用 jps 查看进程，Runjar 是一个单独 java 进程，也是单机集中式部署的特点

```
dase-local@10-24-21-12: ~/softwares/hadoop-2.10.1
File Edit View Search Terminal Help
21/03/18 16:02:02 INFO mapred.MapTask: soft limit at 83886080
21/03/18 16:02:02 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
21/03/18 16:02:02 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
21/03/18 16:02:02 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
21/03/18 16:02:02 INFO mapreduce.Job: map 100% reduce 0%
21/03/18 16:02:03 INFO mapred.MapTask: Spilling map output I
21/03/18 16:02:03 INFO mapred.MapTask: bufstart = 0; bufend = 31339557; bufvoid = 600
21/03/18 16:02:03 INFO mapred.MapTask: kvstart = 26214396(104857584); kvend = 1302311088; length = 13136625/6553600
21/03/18 16:02:03 INFO mapred.MapTask: (EQUATOR) 41825314 kvi 10456324(41825296)
21/03/18 16:02:07 INFO mapred.MapTask: Finished spill 0
21/03/18 16:02:07 INFO mapred.MapTask: (RESET) equator 41825314 kv 10456324(41825296) i 7834896(31339584)
dase-local@10-24-21-12: ~
File Edit View Search Terminal Help
dase-local@10-24-21-12:~$ jps
3202 Jps
3156 RunJar
```

图 4: 运行 WordCount

Wordcount 任务运行结束，在单机配置下耗时约9min（稍后和分布式环境下运行效率进行比较）

```
dase-local@10-24-21-12: ~/softwares/hadoop-2.10.1
File Edit View Search Terminal Help
Combine input records=389987792
Combine output records=2091672
Reduce input groups=247183
Reduce shuffle bytes=28239625
Reduce input records=2091672
Reduce output records=247183
Spilled Records=6275016
Shuffled Maps =65
Failed Shuffles=0
Merged Map outputs=65
GC time elapsed (ms)=4241
Total committed heap usage (bytes)=30669799424
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=2173827625
File Output Format Counters
Bytes Written=3020133
```

图 5: WordCount 运行结束

检查运行后的目录情况

```
dase-local@10-24-21-12:~/softwares/hadoop-2.10.1$ ls
bin etc include lib libexec LICENSE.txt NOTICE.txt README.txt sbin share
dase-local@10-24-21-12:~/softwares/hadoop-2.10.1$
```

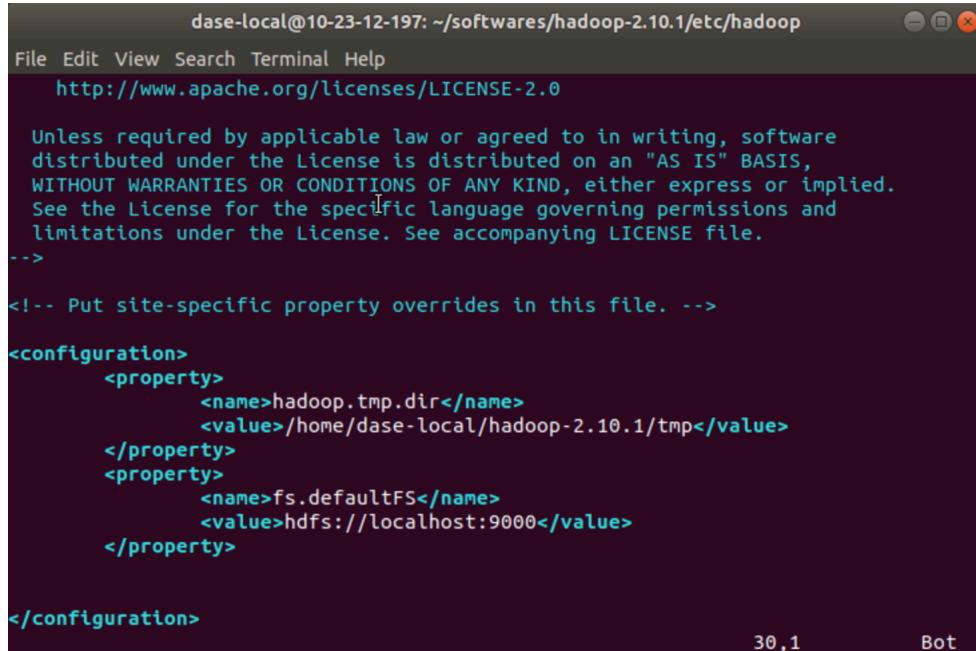
图 6: WordCount 运行结束后的目录

## Section 2

### 单机伪分布式部署

在实验开始先，首先修改配置文件，一些配置文件的功能在上一次报告中有讲到，这次就不再赘述了。

首先是 core-site.html



```
dase-local@10-23-12-197: ~/softwares/hadoop-2.10.1/etc/hadoop
File Edit View Search Terminal Help
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. -->

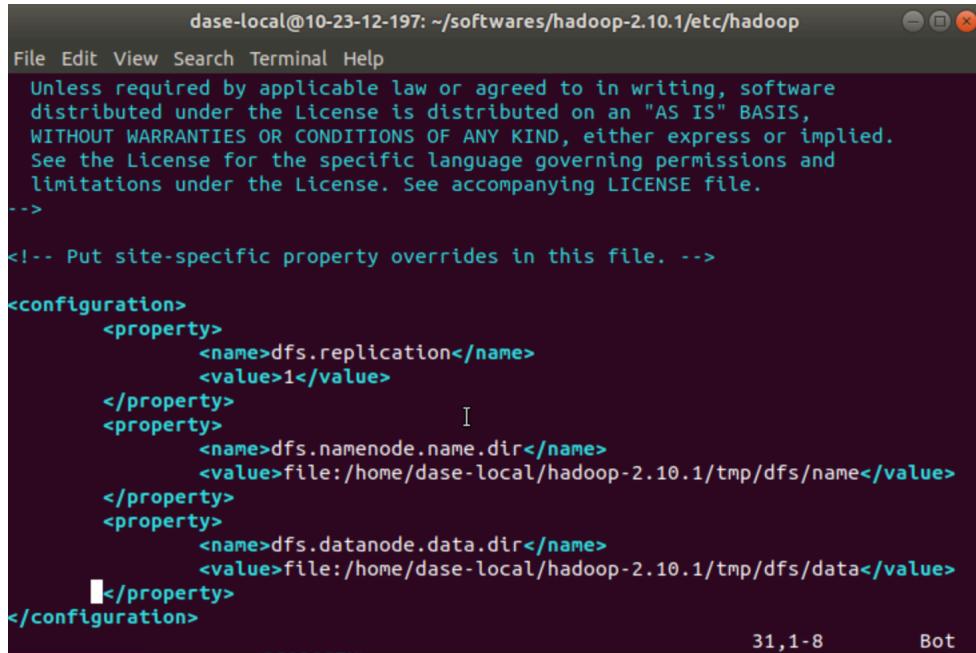
<configuration>
    <property>
        <name>hadoop.tmp.dir</name>
        <value>/home/dase-local/hadoop-2.10.1/tmp</value>
    </property>
    <property>
        <name>fs.defaultFS</name>
        <value>hdfs://localhost:9000</value>
    </property>

</configuration>
```

30,1 Bot

图 7: core-site.html

其次是修改 hdfs-site.html



```
dase-local@10-23-12-197: ~/softwares/hadoop-2.10.1/etc/hadoop
File Edit View Search Terminal Help
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.
-->

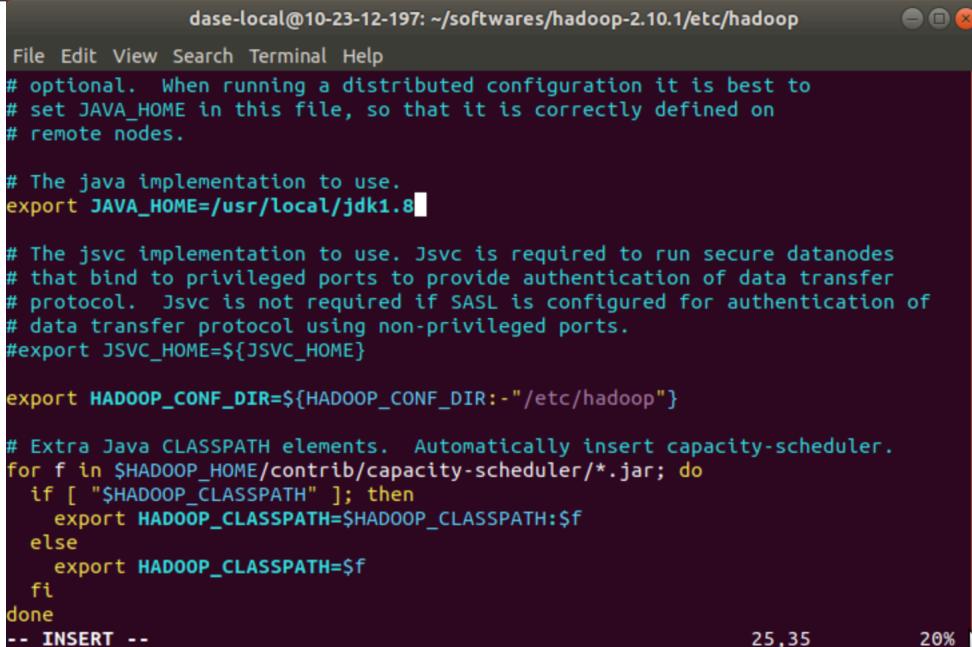
<!-- Put site-specific property overrides in this file. -->

<configuration>
    <property>
        <name>dfs.replication</name>
        <value>1</value>
    </property>
    <property>
        <name>dfs.namenode.name.dir</name>
        <value>file:/home/dase-local/hadoop-2.10.1/tmp/dfs/name</value>
    </property>
    <property>
        <name>dfs.datanode.data.dir</name>
        <value>file:/home/dase-local/hadoop-2.10.1/tmp/dfs/data</value>
    </property>
</configuration>
```

31,1-8 Bot

图 8: hdfs-site.html

最后是修改 hadoop-env.sh



```

dase-local@10-23-12-197: ~/softwares/hadoop-2.10.1/etc/hadoop
File Edit View Search Terminal Help
# optional. When running a distributed configuration it is best to
# set JAVA_HOME in this file, so that it is correctly defined on
# remote nodes.

# The java implementation to use.
export JAVA_HOME=/usr/local/jdk1.8

# The jsvc implementation to use. Jsvc is required to run secure datanodes
# that bind to privileged ports to provide authentication of data transfer
# protocol. Jsvc is not required if SASL is configured for authentication of
# data transfer protocol using non-privileged ports.
#export JSVC_HOME=${JSVC_HOME}

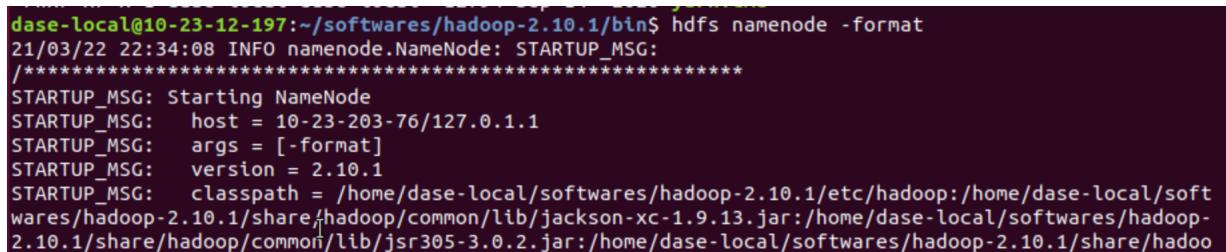
export HADOOP_CONF_DIR=${HADOOP_CONF_DIR:-"/etc/hadoop"}

# Extra Java CLASSPATH elements. Automatically insert capacity-scheduler.
for f in $HADOOP_HOME/contrib/capacity-scheduler/*.jar; do
    if [ "$HADOOP_CLASSPATH" ]; then
        export HADOOP_CLASSPATH=$HADOOP_CLASSPATH:$f
    else
        export HADOOP_CLASSPATH=$f
    fi
done
-- INSERT --

```

图 9: hadoop-env.sh

接下来启动 hdfs 服务，首先要格式化 NameNode，这里需要注意的是如果这台机器之前运行过 NameNode 话，要修改 ssh-key 才能正确启动 namenode 和 datanode，否则只能看到 SecondaryNameNode 被启动了（这个地方，关掉的时候也记得先 kill 掉免得出 bug）

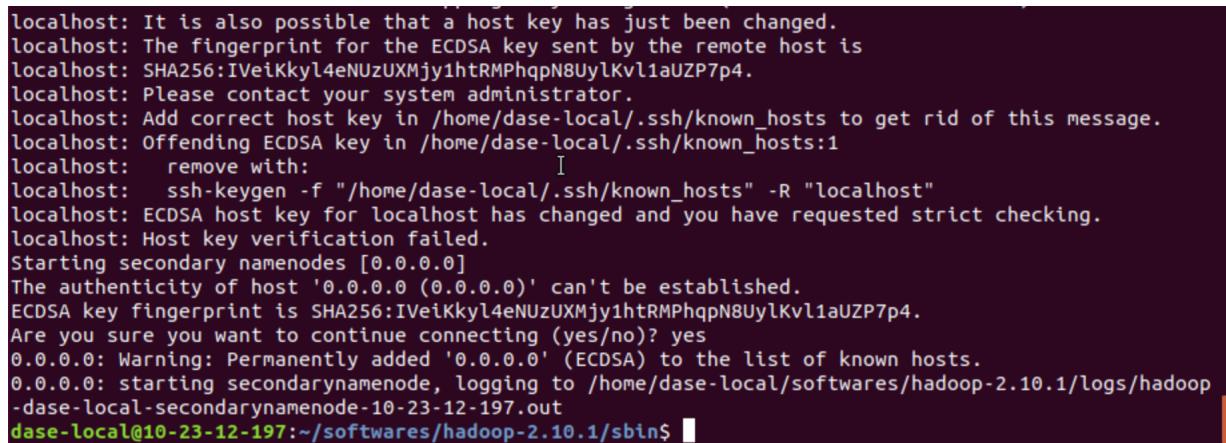


```

dase-local@10-23-12-197:~/softwares/hadoop-2.10.1/bin$ hdfs namenode -format
21/03/22 22:34:08 INFO namenode.NameNode: STARTUP_MSG:
/*****STARTUP_MSG: Starting NameNode
STARTUP_MSG:   host = 10-23-203-76/127.0.1.1
STARTUP_MSG:   args = [-format]
STARTUP_MSG:   version = 2.10.1
STARTUP_MSG:   classpath = /home/dase-local/softwares/hadoop-2.10.1/etc/hadoop:/home/dase-local/softwares/hadoop-2.10.1/share/hadoop/common/lib/jackson-xc-1.9.13.jar:/home/dase-local/softwares/hadoop-2.10.1/share/hadoop/common/lib/jsr305-3.0.2.jar:/home/dase-local/softwares/hadoop-2.10.1/share/hadoop

```

图 10: 格式化 NameNode



```

localhost: It is also possible that a host key has just been changed.
localhost: The fingerprint for the ECDSA key sent by the remote host is
localhost: SHA256:IVeIKkyl4eNUzUXMjy1htRMPPhqpN8UylKvl1aUZP7p4.
localhost: Please contact your system administrator.
localhost: Add correct host key in /home/dase-local/.ssh/known_hosts to get rid of this message.
localhost: Offending ECDSA key in /home/dase-local/.ssh/known_hosts:1
localhost: remove with:
localhost:   ssh-keygen -f "/home/dase-local/.ssh/known_hosts" -R "localhost"
localhost: ECDSA host key for localhost has changed and you have requested strict checking.
localhost: Host key verification failed.
Starting secondary namenodes [0.0.0.0]
The authenticity of host '0.0.0.0 (0.0.0.0)' can't be established.
ECDSA key fingerprint is SHA256:IVeIKkyl4eNUzUXMjy1htRMPPhqpN8UylKvl1aUZP7p4.
Are you sure you want to continue connecting (yes/no)? yes
0.0.0.0: Warning: Permanently added '0.0.0.0' (ECDSA) to the list of known hosts.
0.0.0.0: starting secondarynamenode, logging to /home/dase-local/softwares/hadoop-2.10.1/logs/hadoop-dase-local-secondarynamenode-10-23-12-197.out
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1/sbin$ 

```

图 11: 执行 ssh-keygen

随后重复上述操作，成功启动 hdfs 服务，可以使用 jps 检查 HDFS 是否启动成功

## 华东师范大学数据科学与工程学院学生实验报告

```
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1/sbin$ start-dfs.sh
Starting namenodes on [localhost]
The authenticity of host 'localhost (::1)' can't be established.
ECDSA key fingerprint is SHA256:IVeiKkyl4eNUzUXMjy1htRMPHqpN8UylKvl1aUZP7p4.
Are you sure you want to continue connecting (yes/no)? yes
localhost: Warning: Permanently added 'localhost' (ECDSA) to the list of known hosts.
localhost: starting namenode, logging to /home/dase-local/softwares/hadoop-2.10.1/logs/hadoop-dase-local-namenode-10-23-12-197.out
localhost: starting datanode, logging to /home/dase-local/softwares/hadoop-2.10.1/logs/hadoop-dase-local-datanode-10-23-12-197.out
Starting secondary namenodes [0.0.0.0]
0.0.0.0: starting secondarynamenode, logging to /home/dase-local/softwares/hadoop-2.10.1/logs/hadoop-dase-local-secondarynamenode-10-23-12-197.out
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1/sbin$ jps
5025 DataNode
5399 Jps
4823 NameNode
5273 SecondaryNameNode
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1/sbin$
```

图 12: 启动 HDFS

打开浏览器查看 Web UI，在 localhost 中查看 HDFS 的信息，可以看到有一个 Live Node 正在运行中（就是本机作为这个 Live Node）

The screenshot shows a Linux desktop environment with a terminal window titled "主机显示终端" and a Firefox browser window titled "Namenode information". The Firefox window displays the HDFS Web UI at the URL "localhost:50070/dfshealth.html#tab-overview". The page content includes:

Security is off.  
Safemode is off.  
1 files and directories, 0 blocks = 1 total filesystem object(s).  
Heap Memory used 69.38 MB of 179.5 MB Heap Memory. Max Heap Memory is 889 MB.  
Non Heap Memory used 42.14 MB of 43.44 MB Committed Non Heap Memory. Max Non Heap Memory is <unbounded>.

Configured Capacity:	49.09 GB
DFS Used:	24 KB (0%)
Non DFS Used:	11.25 GB
DFS Remaining:	35.62 GB (72.56%)
Block Pool Used:	24 KB (0%)
DataNodes usages% (Min/Median/Max/stdDev):	0.00% / 0.00% / 0.00% / 0.00%
Live Nodes	1 (Decommissioned: 0, In Maintenance: 0)
Dead Nodes	0 (Decommissioned: 0, In Maintenance: 0)
Decommissioning Nodes	0
Entering Maintenance Nodes	0
Total Datanode Volume Failures	0 (0 B)

图 13: 通过 Web UI 查看 Live Nodes

接下来要熟悉 HDFS Shell 的一些命令，第一次使用时为当前用户创建一个用户根目录 /user/dase-local，然后进行一些操作，如新建/删除目录

## 华东师范大学数据科学与工程学院学生实验报告

```
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ ./bin/dfs dfs -mkdir -p /user/dase-local
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ ./bin/dfs dfs -ls .
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ ./bin/dfs dfs -ls /user/dase-local
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ ./bin/dfs dfs -mkdir input
```

图 14: HDFS Shell 命令练习

```
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ ./bin/dfs dfs -mkdir /input
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ ./bin/dfs dfs -rm -r /input
Deleted /input
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ █
```

图 15: HDFS Shell 命令练习

查看下目录，然后练习在 HDFS 目录中新建/删除文件夹，并且上传文件

```
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ ls
bin etc include lib libexec LICENSE.txt logs NOTICE.txt README.txt sbin share
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ ./bin/dfs dfs -mkdir input
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ ./bin/dfs dfs -put README.txt input/
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ █
```

图 16: HDFS Shell 命令练习

接下来是文件操作，HDFS 中有以下文件操作：上传文件至 HDFS，下载文件到本地，移动文件等  
以下是从本地的文件系统中下载 README.txt，将其上传到 HDFS 之中

```
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ ./bin/dfs dfs -mkdir input
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ ./bin/dfs dfs -put README.txt input/
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ ./bin/dfs dfs -cat input/README.txt
For the latest information about Hadoop, please visit our website at:

    http://hadoop.apache.org/core/
and our wiki, at:
    http://wiki.apache.org/hadoop/
This distribution includes cryptographic software. The country in
which you currently reside may have restrictions on the import,
possession, use, and/or re-export to another country, of
encryption software. BEFORE using any encryption software, please
check your country's laws, regulations and policies concerning the
import, possession, or use, and re-export of encryption software, to
see if this is permitted. See <http://www.wassenaar.org/> for more
information.
```

图 17: HDFS Shell 命令练习

```
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ ./bin/dfs dfs -put ~/input/pd.train input/
dase-local@10-23-12-197:~$ jps
5025 DataNode
6644 Jps
4823 NameNode
6584 FsShell
5273 SecondaryNameNode
dase-local@10-23-12-197:~$ █
```

图 18: 上传文件并使用 jps 查看进程状态

这里还可以通过 get 命令把 HDFS 的文件下载到本地文件系统中，并查看

```
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ bin/hdfs dfs -get input/README.txt ~/Downloads
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ cat ~/Downloads/README.txt
For the latest information about Hadoop, please visit our website at:

  http://hadoop.apache.org/core/
and our wiki, at:
  http://wiki.apache.org/hadoop/
This distribution includes cryptographic software. The country in
which you currently reside may have restrictions on the import,
possession, use, and/or re-export to another country, of
```

图 19: 把 HDFS 文件下载到本地，并查看

最后通过 cp 命令把文件从一个 HDFS 文件夹中的一个目录拷贝到另一个目录之中

```
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ bin/hdfs dfs -cp input/README.txt .
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1$ bin/hdfs dfs -cat README.txt .
For the latest information about Hadoop, please visit our website at:

  http://hadoop.apache.org/core/
and our wiki, at:
  http://wiki.apache.org/hadoop/
This distribution includes cryptographic software. The country in
which you currently reside may have restrictions on the import,
```

图 20: cp 命令拷贝文件夹

结束上述练习后，关闭 HDFS，并使用 jps 查看是否已经停止服务，不再能看到 NameNode，DataNode，SecondaryNameNode 进程表示服务已经停止

```
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1/sbin$ stop-dfs.sh
Stopping namenodes on [localhost]
localhost: stopping namenode
localhost: stopping datanode
Stopping secondary namenodes [0.0.0.0]
0.0.0.0: stopping secondarynamenode
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1/sbin$ jps
7523 Jps
dase-local@10-23-12-197:~/softwares/hadoop-2.10.1/sbin$ █
```

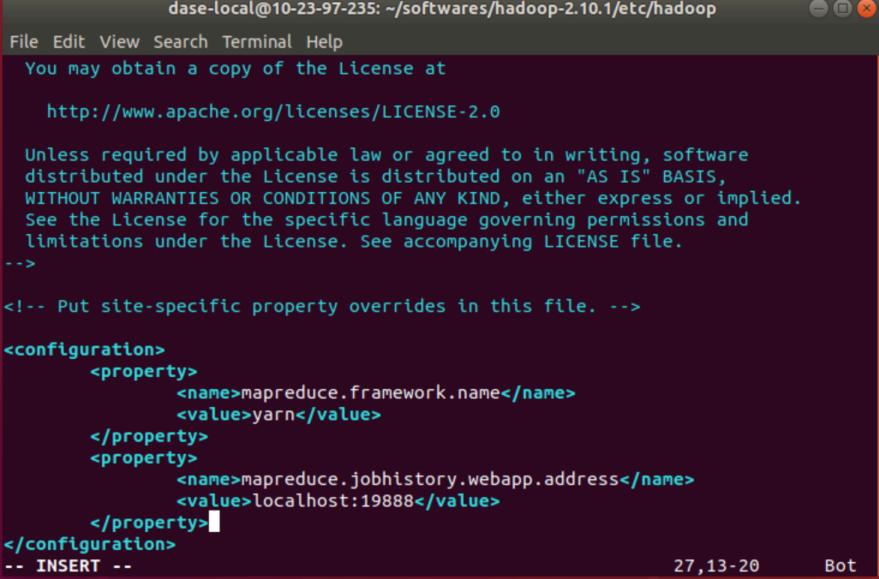
图 21: 关闭 HDFS

可以看到 HDFS 已经被正常关闭，准备进行下一阶段实验

### 实验：MapReduce

在实验开始之前，先修改配置文件

首先是 mapred-site.xml



```
dase-local@10-23-97-235: ~/softwares/hadoop-2.10.1/etc/hadoop
File Edit View Search Terminal Help
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.

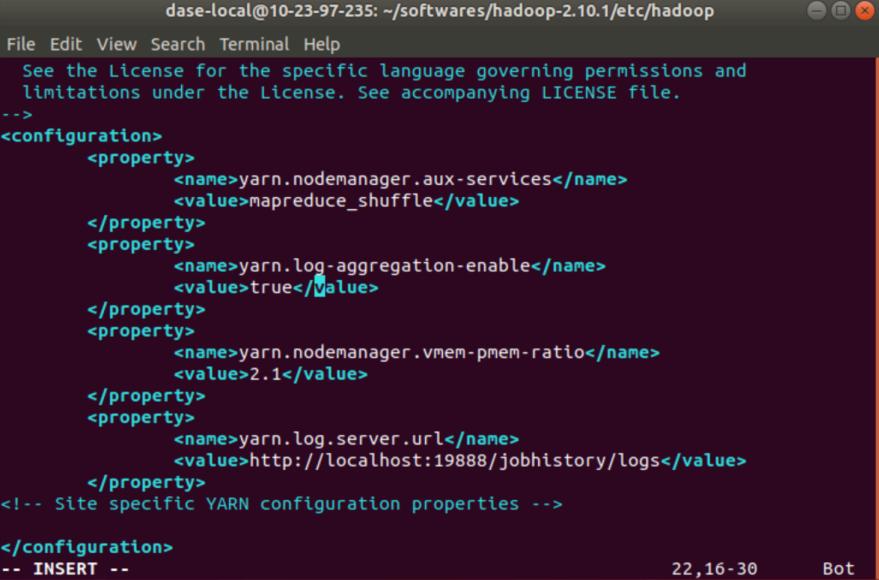
-->

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

<configuration>
    <property>
        <name>mapreduce.framework.name</name>
        <value>yarn</value>
    </property>
    <property>
        <name>mapreduce.jobhistory.webapp.address</name>
        <value>localhost:19888</value>
    </property>
</configuration>
-- INSERT --
27,13-20      Bot
```

图 22: 修改 mapred-site.xml

其次是修改 yarn-site.xml



```
dase-local@10-23-97-235: ~/softwares/hadoop-2.10.1/etc/hadoop
File Edit View Search Terminal Help
See the License for the specific language governing permissions and
limitations under the License. See accompanying LICENSE file.

-->
<configuration>
    <property>
        <name>yarn.nodemanager.aux-services</name>
        <value>mapreduce_shuffle</value>
    </property>
    <property>
        <name>yarn.log-aggregation-enable</name>
        <value>true</value>
    </property>
    <property>
        <name>yarn.nodemanager.vmem-pmem-ratio</name>
        <value>2.1</value>
    </property>
    <property>
        <name>yarn.log.server.url</name>
        <value>http://localhost:19888/jobhistory/logs</value>
    </property>
<!-- Site specific YARN configuration properties -->
</configuration>
-- INSERT --
22,16-30      Bot
```

图 23: 修改 yarn-site.xml

启动 Yarn 服务和 HDFS 服务，同样记得要修改 ssh-key！

```
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1/sbin$ ssh-keygen -f "/home/dase-local/.ssh/known_hosts" -R "0.0.0.0"
# Host 0.0.0.0 found: line 1
/home/dase-local/.ssh/known_hosts updated.
Original contents retained as /home/dase-local/.ssh/known_hosts.old
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1/sbin$ start-dfs.sh
Starting namenodes on [localhost]
localhost: namenode running as process 6023. Stop it first.
localhost: datanode running as process 6232. Stop it first.
Starting secondary namenodes [0.0.0.0]
The authenticity of host '0.0.0.0 (0.0.0.0)' can't be established.
ECDSA key fingerprint is SHA256:FeLMEROR/UQYMpiLKeyCL2LGACfleY0l037yzWd1a8g.
Are you sure you want to continue connecting (yes/no)? yes
0.0.0.0: Warning: Permanently added '0.0.0.0' (ECDSA) to the list of known hosts
.
0.0.0.0: starting secondarynamenode, logging to /home/dase-local/softwares/hadoop-2.10.1/logs/hadoop-dase-local-secondarynamenode-10-23-97-235.out
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1/sbin$
```

图 24: 再次修改 ssh-key

重新启动后可以查看 Yarn 的服务信息

```
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1/sbin$ jps
5346 ResourceManager
6023 NameNode
6232 DataNode
5533 NodeManager
7053 Jps
6927 SecondaryNameNode
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1/sbin$
```

图 25: Yarn 的服务信息

访问 Yarn 的 Web 界面，查看集群信息可以发现有 active 的节点（端口号 8088）

The screenshot shows a Firefox browser window titled 'Firefox Web Browser' with the URL 'localhost:8088/cluster'. The page displays the Hadoop YARN Web Interface. On the left, there's a sidebar with navigation links for Cluster (About, Nodes, Node Labels, Applications), Application Status (NEW, NEW\_SAVING, SUBMITTED, ACCEPTED, RUNNING, FINISHED, FAILED, KILLED), and Scheduler. The main content area has three sections: 'Cluster Metrics' (Shows 0 Apps Submitted, 0 Apps Pending, 0 Apps Running, 0 Apps Completed), 'Cluster Nodes Metrics' (Shows 1 Active Node, 0 Decommissioning Nodes, 0 Decommissioned), and 'Scheduler Metrics' (Shows Capacity Scheduler details). Below these is a table titled 'Show 20 entries' with columns for ID, User, Name, Application Type, Queue, Application Priority, Start Time, Launch Time, Finish Time, State, and Firm.

图 26: Yarn 的 Web 界面

接下来我们提交一个 MapReduce 程序

```

File Edit View Search Terminal Help
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1$ cd >>
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1$ ./bin/hdfs dfs -mkdir -p input/grep
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1$ ./bin/hdfs dfs -put etc/hadoop/* input/grep
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1$ ./bin/hdfs dfs -rm -r output/grep
rm: 'output/grep': No such file or directory
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1$ ./bin/yarn jar share/hadoop/mapreduce/hadoop-mapreduce-examples-2.10.1.jar grep input/grep output/grep 'dfs[a-z.]+'
21/03/23 12:45:00 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
21/03/23 12:45:01 INFO input.FileInputFormat: Total input files to process : 30
21/03/23 12:45:01 INFO mapreduce.JobSubmitter: number of splits:30
21/03/23 12:45:02 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1616474075838_0001
21/03/23 12:45:02 INFO conf.Configuration: resource-types.xml not found
21/03/23 12:45:02 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
21/03/23 12:45:02 INFO resource.ResourceUtils: Adding resource type - name = memory-mb, units = Mi, type = COUNTABLE
21/03/23 12:45:02 INFO resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE
21/03/23 12:45:02 INFO impl.YarnClientImpl: Submitted application application_1616474075838_0001
21/03/23 12:45:02 INFO mapreduce.Job: The url to track the job: http://10-23-203-76:8088/proxy/application_1616474075838_0001/
21/03/23 12:45:02 INFO mapreduce.Job: Running job: job_1616474075838_0001
21/03/23 12:45:09 INFO mapreduce.Job: Job job_1616474075838_0001 running in uber mode : false
21/03/23 12:45:09 INFO mapreduce.Job: map 0% reduce 0%

```

图 27: 提交 MapReduce 程序

执行完 jar 任务后可以查看结果

```

WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0
File Input Format Counters
Bytes Read=564
File Output Format Counters
Bytes Written=268
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1$ ./bin/hdfs dfs -cat output/grep/p*
6      dfs.audit.logger
4      dfs.class
3      dfs.logger
3      dfs.server.namenode.
2      dfs.audit.log.maxbackupindex
2      dfs.period
2      dfs.audit.log.maxfilesize
1      dfs.replication
1      dfs.log
1      dfs.file
1      dfs.datanode.data.dir
1      dfs.servers
1      dfsadmin
1      dfsmetrics.log
1      dfs.namenode.name.dir

```

图 28: 查看结果

查看执行过程，运行一个 Wordcount 示例，并且查看目前正在运行的进程

```
dase-local@10-23-97-235: ~/softwares/hadoop-2.10.1
File Edit View Search Terminal Help
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1$ cd >>
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1$ ./bin/hdfs dfs -mkdir -p input/grep
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1$ ./bin/hdfs dfs -put etc/hadoop/* input/grep
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1$ ./bin/hdfs dfs -rm -r output/grep
rm: 'output/grep': No such file or directory
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1$ ./bin/yarn jar share/hadoop/mapreduce/hadoop-mapreduce-examples-2.10.1.jar grep input/grep output/grep 'dfs[a-z.]+'
21/03/23 12:45:00 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
21/03/23 12:45:01 INFO input.FileInputFormat: Total input files to process : 30
21/03/23 12:45:01 INFO mapreduce.JobSubmitter: number of splits:30
21/03/23 12:45:02 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1616474075838_0001
21/03/23 12:45:02 INFO conf.Configuration: resource-types.xml not found
21/03/23 12:45:02 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
21/03/23 12:45:02 INFO resource.ResourceUtils: Adding resource type - name = memory-mb, units = Mi, type = COUNTABLE
21/03/23 12:45:02 INFO resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE
21/03/23 12:45:02 INFO impl.YarnClientImpl: Submitted application application_1616474075838_0001
21/03/23 12:45:02 INFO mapreduce.Job: The url to track the job: http://10-23-203-76:8088/proxy/application_1616474075838_0001/
21/03/23 12:45:02 INFO mapreduce.Job: Running job: job_1616474075838_0001
21/03/23 12:45:09 INFO mapreduce.Job: Job job_1616474075838_0001 running in uber mode : false
21/03/23 12:45:09 INFO mapreduce.Job: map 0% reduce 0%
```

图 29: Wordcount 示例

可以看到三类之前没有出现过的进程: RunJar, MRAppMaster, YarnChild, 可见在伪分布式模式下, 系统的确在使用多进程并行计算, 在图形界面中也可以看到执行的 MapReduce 任务

The screenshot shows the Hadoop Web UI interface. At the top, there's a navigation bar with tabs for 'All Applications' and a search bar set to 'localhost:8088/cluster'. Below the header is the Hadoop logo. On the left, a sidebar menu under 'Cluster' includes 'About', 'Nodes', 'Node Labels', 'Applications' (selected), 'NEW', 'NEW\_SAVING', 'SUBMITTED', 'ACCEPTED', 'RUNNING', 'FINISHED', 'FAILED', 'KILLED', and 'Scheduler'. Under 'Scheduler', there are links for 'Tools' and 'Capacity Scheduler'. The main content area is titled 'Cluster Metrics' and displays four tables: 'Apps Submitted' (3), 'Apps Pending' (0), 'Apps Running' (1), and 'Apps Completed' (2). Below this is the 'Cluster Nodes Metrics' section with 'Active Nodes' (1) and 'Decommissioning Nodes' (0). The 'Scheduler Metrics' section shows the 'Capacity Scheduler' configuration with 'memory-mb' as the default unit and 'COUNTABLE' as the type. The 'Scheduler' table lists three running applications:

ID	User	Name	Application Type	Queue	Application Priority	StartTime
application_1616474075838_0003	dase-local	word count	MAPREDUCE	default	0	Tue Mar 23 12:49:21 +0800 2021
application_1616474075838_0002	dase-local	grep sort	MAPREDUCE	default	0	Tue Mar 23 12:45:57 +0800 2021
application_1616474075838_0001	dase-local	grep search	MAPREDUCE	default	0	Tue Mar 23 12:45:02 +0800 2021

图 30: 通过 Web UI 查看程序运行信息

启动任务日志历史服务器, 可以检查 MapReduce 应用程序日志

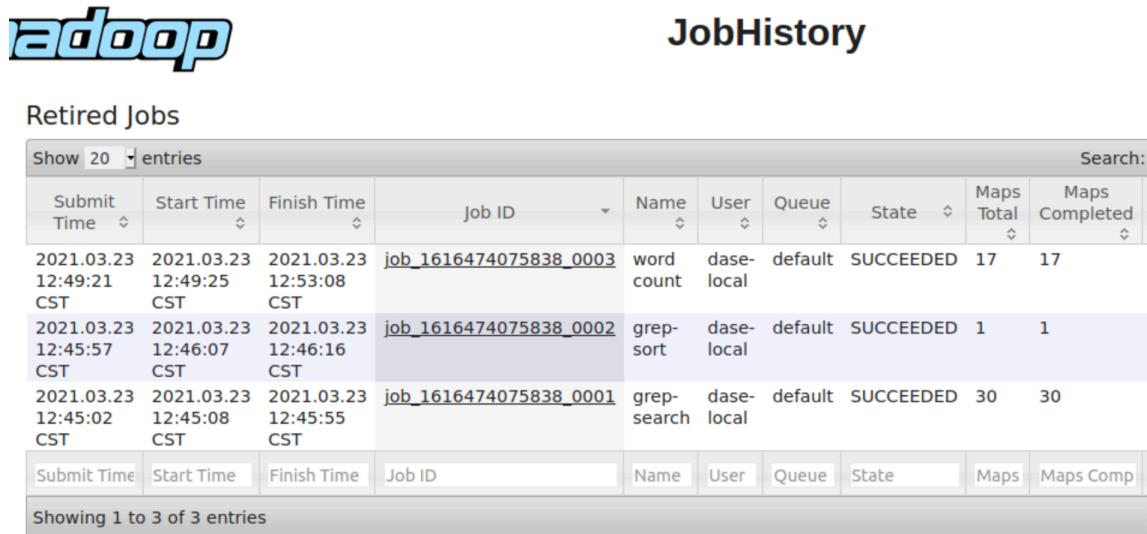
Remark: 其实应该在任务开始前启动, 然而一开始的命令出现了一点小问题, 所以在这里补一下

## 华东师范大学数据科学与工程学院学生实验报告

```
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1$ sbin/mr-jobhistory-daemon.sh start historyserver
starting historyserver, logging to /home/dase-local/softwares/hadoop-2.10.1/logs/mapred-dase-local-hi
storyserver-10-23-97-235.out
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1$
```

图 31: 启动 HistoryServer

稍微让我感到疑惑的是，虽然是后启动的历史记录服务器，但依然可以看到之前的行为，稍后再探索在可视化界面中可以看到程序的历史运行记录



The screenshot shows the Hadoop JobHistory interface. At the top, there's a logo for 'adoop' and the title 'JobHistory'. Below that, a section titled 'Retired Jobs' displays a table of completed jobs. The table has columns for Submit Time, Start Time, Finish Time, Job ID, Name, User, Queue, State, Maps Total, and Maps Completed. There are three entries listed:

Submit Time	Start Time	Finish Time	Job ID	Name	User	Queue	State	Maps Total	Maps Completed
2021.03.23 12:49:21 CST	2021.03.23 12:49:25 CST	2021.03.23 12:53:08 CST	job_1616474075838_0003	word count	dase-local	default	SUCCEEDED	17	17
2021.03.23 12:45:57 CST	2021.03.23 12:46:07 CST	2021.03.23 12:46:16 CST	job_1616474075838_0002	grep-sort	dase-local	default	SUCCEEDED	1	1
2021.03.23 12:45:02 CST	2021.03.23 12:45:08 CST	2021.03.23 12:45:55 CST	job_1616474075838_0001	grep-search	dase-local	default	SUCCEEDED	30	30

At the bottom of the table, there are buttons for 'Submit Time', 'Start Time', 'Finish Time', 'Job ID', 'Name', 'User', 'Queue', 'State', 'Maps', and 'Maps Comp'. A message at the bottom says 'Showing 1 to 3 of 3 entries'.

图 32: 启动 HistoryServer

最后，停止 Yarn 服务，并且使用 jps 检查是否成功停止服务

```
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1$ sbin/stop-dfs.sh
Stopping namenodes on [localhost]
localhost: stopping namenode
localhost: stopping datanode
Stopping secondary namenodes [0.0.0.0]
0.0.0.0: stopping secondarynamenode
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1$ sbin/mr-jobhistory-daemon.sh stop historyserver
stopping historyserver
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1$ sbin/stop-yarn.sh
stopping yarn daemons
no resourcemanager to stop
localhost: no nodemanager to stop
no proxyserver to stop
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1$ jps
12469 Jps
dase-local@10-23-97-235:~/softwares/hadoop-2.10.1$
```

图 33: 停止 Yarn 服务

正常停止后，发现 ResourceManager、NodeManager、JobHistoryServer 等进程都不再出现，说明已经正常关闭了

### Section 3

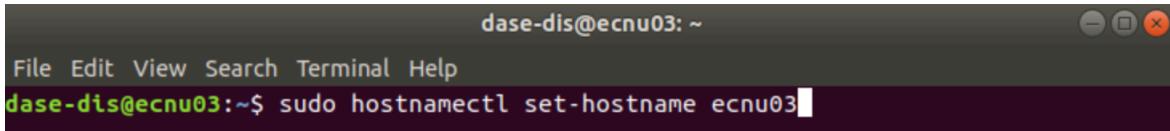
#### 分布式部署

这里需要开始团队协作完成了，准备四台机器分别担任主节点，客户端和两个从节点，并且使用 dase-dis 用户进行实验，我使用的是从节点 ecnu03

在实验开始前，使用命令 HOSTNAMECTL 配置主机名

```
hostnamectl set-hostname name [option...]
```

其中 option 是 `-pretty`, `-static` 或者 `-transient` 中的一个或多个选项。Static 和 transient 主机名会简化为 pretty 主机名格式。使用 “-” 替换空格，并删除特殊字符。如果 `-static` 或 `-transient` 选项与 `-pretty` 选项一同使用，则会将 static 和 transient 主机名简化为 pretty 主机名格式。使用 “-” 替换空格，并删除特殊字符。如果未使用 `-pretty` 选项，则不会发生简化



```
dase-dis@ecnu03: ~
File Edit View Search Terminal Help
dase-dis@ecnu03:~$ sudo hostnamectl set-hostname ecnu03
```

图 34: `sudo hostnamectl set-hostname ecnu0x`

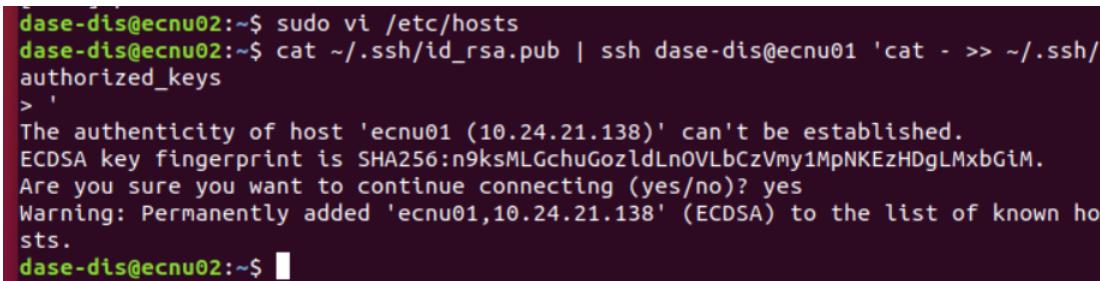
### 调整四台机器名称

随后进行一些分布式部署的准备工作

首先实现多台机器之间实现**免密登录**

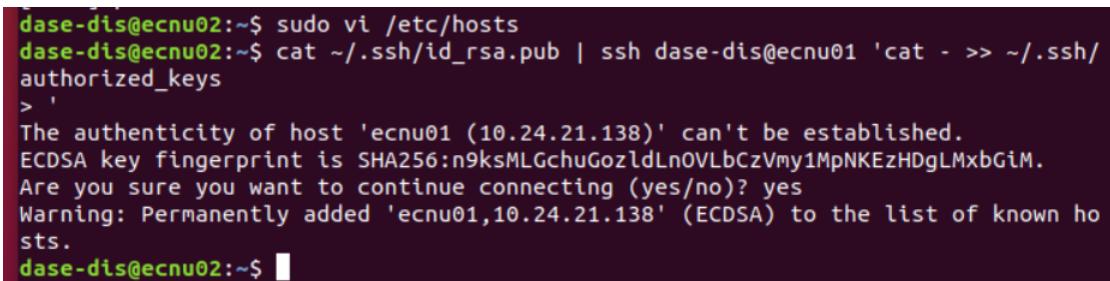
- 创建相同的用户 `dase-dis`, 命名主机
- 在每台机器都修改 `etc/hosts` 文件，添加主机名与 IP 地址的映射
- 将其它三台主机的公钥发到主节点上的授权文件

修改 `etc/hosts` 文件，添加主机名与 IP 地址的映射



```
dase-dis@ecnu02:~$ sudo vi /etc/hosts
dase-dis@ecnu02:~$ cat ~/.ssh/id_rsa.pub | ssh dase-dis@ecnu01 'cat - > ~/.ssh/
authorized_keys
> '
The authenticity of host 'ecnu01 (10.24.21.138)' can't be established.
ECDSA key fingerprint is SHA256:n9ksMLGchuGozldLnOVLbCzVmy1MpNKEzHDgLMxbGiM.
Are you sure you want to continue connecting (yes/no)? yes
Warning: Permanently added 'ecnu01,10.24.21.138' (ECDSA) to the list of known ho
sts.
dase-dis@ecnu02:~$
```

图 35: `ecnu01` (这里是主节点同学忘记改名了): 修改 `etc/hosts` 文件，添加主机名与 IP 地址的映射



```
dase-dis@ecnu02:~$ sudo vi /etc/hosts
dase-dis@ecnu02:~$ cat ~/.ssh/id_rsa.pub | ssh dase-dis@ecnu01 'cat - > ~/.ssh/
authorized_keys
> '
The authenticity of host 'ecnu01 (10.24.21.138)' can't be established.
ECDSA key fingerprint is SHA256:n9ksMLGchuGozldLnOVLbCzVmy1MpNKEzHDgLMxbGiM.
Are you sure you want to continue connecting (yes/no)? yes
Warning: Permanently added 'ecnu01,10.24.21.138' (ECDSA) to the list of known ho
sts.
dase-dis@ecnu02:~$
```

图 36: 从节点 `ecnu02`: 修改 `etc/hosts` 文件，添加主机名与 IP 地址的映射

## 华东师范大学数据科学与工程学院学生实验报告

```
dase-dis@ecnu03:~$ sudo vim /etc/hosts
[sudo] password for dase-dis:
dase-dis@ecnu03:~$ cat ~/.ssh/id_rsa.pub | ssh dase-dis@ecnu01 'cat - >> ~/.ssh/authorized_keys'
The authenticity of host 'ecnu01 (10.24.21.138)' can't be established.
ECDSA key fingerprint is SHA256:n9ksMLGchuGozldLn0VLbCzVmy1MpNKEzHDgLMxbGiM.
Are you sure you want to continue connecting (yes/no)? yes
Warning: Permanently added 'ecnu01,10.24.21.138' (ECDSA) to the list of known hosts.
```

图 37: 从节点 ecnu03: 修改 etc/hosts 文件, 添加主机名与 IP 地址的映射

随后, 主节点将授权文件发给其它三台主机, 覆盖原有文件

```
dase-dis@10-24-21-138:~$ cat ~/.ssh/id_rsa.pub | ssh dase-dis@ecnu02 'cat - >> ~/.ssh/authorized_keys'
The authenticity of host 'ecnu02 (10.23.92.145)' can't be established.
ECDSA key fingerprint is SHA256:jF1MGAXSBwP/PAyQKR6PMetS1+rAjpSqC/0c56C12Kk.
Are you sure you want to continue connecting (yes/no)? yes
Warning: Permanently added 'ecnu02,10.23.92.145' (ECDSA) to the list of known hosts.
dase-dis@10-24-21-138:~$ cat ~/.ssh/id_rsa.pub | ssh dase-dis@ecnu03 'cat - >> ~/.ssh/authorized_keys'
The authenticity of host 'ecnu03 (10.23.192.116)' can't be established.
ECDSA key fingerprint is SHA256:HTAmLd86tG9eZmpt/09kPtQqsiygPXJBXA+va0GPctk.
Are you sure you want to continue connecting (yes/no)? yes
Warning: Permanently added 'ecnu03,10.23.192.116' (ECDSA) to the list of known hosts.
dase-dis@10-24-21-138:~$ cat ~/.ssh/id_rsa.pub | ssh dase-dis@ecnu04 'cat - >> ~/.ssh/authorized_keys'
The authenticity of host 'ecnu04 (10.24.21.14)' can't be established.
```

图 38: 主节点将授权文件发给其它三台主机

主节点将授权文件发给其它三台主机进行验证验证

```
dase-dis@10-24-21-138:~$ ssh dase-dis@ecnu03
Welcome to Ubuntu 18.04.5 LTS (GNU/Linux 4.19.0-9.ucloud x86_64)

 * Documentation:  https://help.ubuntu.com
 * Management:    https://landscape.canonical.com
 * Support:        https://ubuntu.com/advantage

 System information as of Thu Mar 25 15:36:12 CST 2021

 System load:  0.02           Processes:      308
 Usage of /:   27.1% of 49.09GB  Users logged in:  2
 Memory usage: 51%           IP address for eth0: 10.23.192.116
 Swap usage:   0%
```

图 39: Authentication

将之前下载的 Hadoop-2.10.1.tar.gz 拷贝到用户目录下并解压

```
dase-dis@ecnu01:~$ tar -xzf hadoop-2.10.1.tar.gz
dase-dis@ecnu01:~$ ls
Desktop  Downloads      hadoop-2.10.1      Music      Public      Videos
Documents examples.desktop  hadoop-2.10.1.tar.gz  Pictures  Templates
dase-dis@ecnu01:~$ cd hadoop-2.10.1
dase-dis@ecnu01:~/hadoop-2.10.1$ ls
bin  etc  include  lib  libexec  LICENSE.txt  NOTICE.txt  README.txt  sbin  share
dase-dis@ecnu01:~/hadoop-2.10.1$ ./bin/hadoop version
Hadoop 2.10.1
Subversion https://github.com/apache/hadoop -r 1827467c9a56f133025f28557bfc2c562d78e816
Compiled by centos on 2020-09-14T13:17Z
Compiled with protoc 2.5.0
From source with checksum 3114edef868f1f3824e7d0f68be03650
This command was run using /home/dase-dis/hadoop-2.10.1/share/hadoop/common/hadoop-common-2.10.1.jar
dase-dis@ecnu01:~/hadoop-2.10.1$
```

图 40: 主节点解压 Hadoop-2.10.1.tar.gz

```
dase-dis@ecnu03:~/hadoop-2.10.1/bin$ ls -l
total 868
-rwxr-xr-x 1 dase-dis dase-dis 371688 Mar 25 15:56 container-executor
-rwxr-xr-x 1 dase-dis dase-dis 6656 Mar 25 15:56 hadoop
-rwxr-xr-x 1 dase-dis dase-dis 8786 Mar 25 15:56 hadoop.cmd
-rwxr-xr-x 1 dase-dis dase-dis 13032 Mar 25 15:56 hdfs
-rwxr-xr-x 1 dase-dis dase-dis 8371 Mar 25 15:56 hdfs.cmd
-rwxr-xr-x 1 dase-dis dase-dis 6237 Mar 25 15:56 mapred
-rwxr-xr-x 1 dase-dis dase-dis 6310 Mar 25 15:56 mapred.cmd
-rwxr-xr-x 1 dase-dis dase-dis 1776 Mar 25 15:56 rcc
-rwxr-xr-x 1 dase-dis dase-dis 410072 Mar 25 15:56 test-container-executor
-rwxr-xr-x 1 dase-dis dase-dis 15747 Mar 25 15:56 yarn
-rwxr-xr-x 1 dase-dis dase-dis 12794 Mar 25 15:56 yarn.cmd
dase-dis@ecnu03:~/hadoop-2.10.1/bin$
```

图 41: 从节点查看文件目录

接下来，在主节点修改 HDFS 配置

- 修改 slaves 文件，注意把 localhost 删掉
- 修改 core-site.xml
- 修改 hdfs-site.xml 文件
- 修改 hadoop-env.sh 文件
- 将整个 hadoop-2.10.1 文件夹拷贝至其余机器

#### 修改配置文件

注意：修改之前[务必停止正在运行的 HDFS 服务](#)

在主节点执行以下操作

· 修改 slaves 文件（文件路径： ~/hadoop-2.10.1/etc/hadoopaves）

ecnu02

ecnu03

· 修改 core-site.xml(文件路径： ~/hadoop-2.10.1/etc/hadoop/core-site.xml)

```
<configuration>
<property>
<name>hadoop.tmp.dir</name>
<value>/home/dase-dis/hadoop-2.10.1/tmp<alue>
</property>
<property>
<name>fs.defaultFS</name>
<value>hdfs://ecnu01:9000<alue>
</property>
</configuration>
```

· 修改hdfs-site.xml(文件路径： ~/hadoop-2.10.1/etc/hadoopfs-site.xml)

```
<configuration>
<property>
<name>dfs.replication</name>
<value>2<value>
</property>
```

```

<property>
<name>dfs.namenode.name.dir</name>
<value>file:/home/dase-dis/hadoop-2.10.1/tmp/dfs/name</value>
</property>
<property>
<name>dfs.datanode.name.dir</name>
<value>file:/home/dase-dis/hadoop-2.10.1/tmp/dfs/data</value>
</property>
</configuration>

```

tips: 这个地方一定要注意 typo, 我们因为 `<name>fs.defaultFS</name>` 的 FS 大小写错误在这里卡了很久, 浪费了很多时间, 除此之外, 配置文件性质的顺序也要注意不要错。

修改 `hadoop-env.sh` 文件

找到 `#export JAVA_HOME = JAVA_HOME` 这行

将此行修改为 `export JAVA_HOME = /USR/LOCAL/JDK1.8` 时务必去掉注释。

### 主节点启动 HDFS 服务

由主节点启动 HDFS 服务, 从节点可以通过 `jps` 查看自己的状态

```

dase-dis@10-24-21-138:~/hadoop-2.10.1$ ~/hadoop-2.10.1/sbin/start-dfs.sh
Starting namenodes on [ecnu01]
The authenticity of host 'ecnu01 (10.24.21.138)' can't be established.
ECDSA key fingerprint is SHA256:n9ksMLGchUGozldLnoVLBCzVmy1MpNKEzHDgLMxbGiM.
Are you sure you want to continue connecting (yes/no)? yes
ecnu01: Warning: Permanently added 'ecnu01,10.24.21.138' (ECDSA) to the list of known hosts.
ecnu01: starting namenode, logging to /home/dase-dis/hadoop-2.10.1/logs/hadoop-d
ase-dis-namenode-ecnu01.out
ecnu03: starting datanode, logging to /home/dase-dis/hadoop-2.10.1/logs/hadoop-d
ase-dis-datanode-ecnu03.out
ecnu02: starting datanode, logging to /home/dase-dis/hadoop-2.10.1/logs/hadoop-d
ase-dis-datanode-ecnu02.out
Starting secondary namenodes [0.0.0.0]
0.0.0.0: starting secondarynamenode, logging to /home/dase-dis/hadoop-2.10.1/log
s/hadoop-dase-dis-secondarynamenode-ecnu01.out

```

图 42: 主节点启动 HDFS

```

dase-dis@10-24-21-138:~/hadoop-2.10.1$ jps
6976 Jps
6850 SecondaryNameNode
6570 NameNode

```

图 43: 主节点使用 `jps` 检查状态

实验的进程日志记录在 `/hadoop-2.10.1/logs` 路径下, 后缀为.log 文件中  
-NameNode 进程日志:

默认位置: `/hadoop-2.10.1/logs/hadoop-* -namenode-* .log`

-DataNode 进程日志:

默认位置: `/hadoop-2.10.1/logs/hadoop-* -datanode-* .log`

-SecondaryNameNode 进程日志:

默认位置: `/hadoop-2.10.1/logs/hadoop-* -secondarynamenode-* .log`

与之前的单机模式一样, 同样可以通过 Web UI 来查看 HDFS 的信息

## 华东师范大学数据科学与工程学院学生实验报告

The screenshot shows the Hadoop Web UI Overview page. At the top, there's a navigation bar with tabs for 'Overview', 'Datanodes', 'Datanode Volume Failures', 'Snapshot', and 'Startup Progress'. Below the navigation bar is a table with the following data:

<b>Started:</b>	Thu Mar 25 19:26:22 +0800 2021
<b>Version:</b>	2.10.1, r1827467c9a56f133025f28557bfc2c562d78e816
<b>Compiled:</b>	Mon Sep 14 21:17:00 +0800 2020 by centos from branch-2.10.1
<b>Cluster ID:</b>	CID-3c777a61-6afb-4c81-85f4-968cecb8a91c
<b>Block Pool ID:</b>	BP-1547215439-10.24.21.138-1616660711504

图 44: 主节点使用 Web UI 查看信息

The screenshot shows the Hadoop Web UI Summary page. It displays various system statistics and resource usage details. A table provides a summary of configured and used capacity:

<b>Configured Capacity:</b>	9318 GB
<b>DFS Used:</b>	64 KB (0%)
<b>Non DFS Used:</b>	27.62 GB
<b>DFS Remaining:</b>	66.13 GB (67.36%)
<b>Block Pool Used:</b>	64 KB (0%)
<b>DataNodes usages% (Min/Median/Max/stdDev):</b>	0.00% / 0.00% / 0.00% / 0.00%
<b>Live Nodes</b>	2 (Decommissioned: 0, In Maintenance: 0)

图 45: 主节点使用 Web UI 查看摘要

```
dase-dis@ecnu02:~/hadoop-2.10.1$ vi etc/hadoop/hdfs-site.xml
dase-dis@ecnu02:~/hadoop-2.10.1$ jps
14903 DataNode
14989 Jps
dase-dis@ecnu02:~/hadoop-2.10.1$
```

图 46: 从节点 ecnu02 查看状态

```
dase-dis@ecnu03:~$ jps
12340 DataNode
12426 Jps
dase-dis@ecnu03:~$
```

图 47: 从节点 ecnu03 查看状态

接下来在分布式环境下执行常用的 HDFS Shell 操作 HDFS 中的文件操作包括：上传文件，下载文件，移动文件等。示例如下：

```
dase-dis@ecnu04:~/hadoop-2.10.1$ ./bin/hdfs dfs -ls /user/dase-dis
dase-dis@ecnu04:~/hadoop-2.10.1$ ./bin/hdfs dfs -mkdir input
dase-dis@ecnu04:~/hadoop-2.10.1$ ./bin/hdfs dfs -rm -r input
Deleted input
dase-dis@ecnu04:~/hadoop-2.10.1$ ./bin/hdfs dfs -mkdir /input
dase-dis@ecnu04:~/hadoop-2.10.1$ ./bin/hdfs dfs -rm -r /input
Deleted /input
```

图 48: 分布式环境 HDFS 目录操作

```
dase-dis@ecnu04:~/hadoop-2.10.1$ cd ~/hadoop-2.10.1
dase-dis@ecnu04:~/hadoop-2.10.1$ ./bin/hdfs dfs -mkdir input
dase-dis@ecnu04:~/hadoop-2.10.1$ ./bin/hdfs dfs -put README.txt input/
dase-dis@ecnu04:~/hadoop-2.10.1$ ./bin/hdfs dfs -cat input/README.txt
For the latest information about Hadoop, please visit our website at:

  http://hadoop.apache.org/core/
and our wiki, at:
  http://wiki.apache.org/hadoop/
```

This distribution includes cryptographic software. The country in which you currently reside may have restrictions on the import, possession, use, and/or re-export to another country, of encryption software. BEFORE using any encryption software, please check your country's laws, regulations and policies concerning the import, possession, or use, and re-export of encryption software, to

图 49: 分布式环境 HDFS 文件操作

从本地文件系统向 HDFS 中上传文件将之前下载保存至 /input 的文件 pd.train 上传至 HDFS 把 HDFS 中的文件下载到本地文件系统中把文件从 HDFS 中的一个目录拷贝至 HDFS 中的另一个目录

```
dase-dis@ecnu04:~/hadoop-2.10.1$ scp dase-local@localhost:~/input/pd.train ~/input
dase-local@localhost's password:
pd.train                                         100% 2073MB  77.9MB/s   00:26
dase-dis@ecnu04:~/hadoop-2.10.1$ ./bin/hdfs dfs -put ~/input/pd.train input/
```

图 50: 布式环境 HDFS 上传 pd.train

## 华东师范大学数据科学与工程学院学生实验报告

```
dase-dis@ecnu04:~/hadoop-2.10.1$ ~/hadoop-2.10.1/bin/hdfs dfs -cp input/README.txt .
dase-dis@ecnu04:~/hadoop-2.10.1$ ~/hadoop-2.10.1/bin/hdfs dfs -cat README.txt
For the latest information about Hadoop, please visit our website at:

    http://hadoop.apache.org/core/

and our wiki, at:

    http://wiki.apache.org/hadoop/

This distribution includes cryptographic software. The country in
which you currently reside may have restrictions on the import,
possession, use, and/or re-export to another country, of
encryption software. BEFORE using any encryption software, please
check your country's laws, regulations and policies concerning the
import, possession, or use, and re-export of encryption software, to
see if this is permitted. See <http://www.wassenaar.org/> for more
information.

The U.S. Government Department of Commerce, Bureau of Industry and
Security (BIS), has classified this software as Export Commodity
Control Number (ECCN) 5D002.C.1, which includes information security
software using or performing cryptographic functions with asymmetric
algorithms. The form and manner of this Apache Software Foundation
distribution makes it eligible for export under the License Exception
ENC Technology Software Unrestricted (TSU) exception (see the BIS
Export Administration Regulations, Section 740.13) for both object
code and source code.

The following provides more details on the included cryptographic
software:
  Hadoop Core uses the SSL libraries from the Jetty project written
by mortbay.org.
dase-dis@ecnu04:~/hadoop-2.10.1$
```

图 51: 布式环境 HDFS 文件拷贝

```
dase-dis@ecnu04:~/hadoop-2.10.1$ ~/hadoop-2.10.1/bin/hdfs dfs -get input/README.txt ~/Downloads
dase-dis@ecnu04:~/hadoop-2.10.1$ cat ~/Downloads/README.txt
For the latest information about Hadoop, please visit our website at:

    http://hadoop.apache.org/core/

and our wiki, at:

    http://wiki.apache.org/hadoop/

This distribution includes cryptographic software. The country in
which you currently reside may have restrictions on the import,
possession, use, and/or re-export to another country, of
encryption software. BEFORE using any encryption software, please
check your country's laws, regulations and policies concerning the
import, possession, or use, and re-export of encryption software, to
see if this is permitted. See <http://www.wassenaar.org/> for more
information.

The U.S. Government Department of Commerce, Bureau of Industry and
Security (BIS), has classified this software as Export Commodity
Control Number (ECCN) 5D002.C.1, which includes information security
software using or performing cryptographic functions with asymmetric
algorithms. The form and manner of this Apache Software Foundation
distribution makes it eligible for export under the License Exception
ENC Technology Software Unrestricted (TSU) exception (see the BIS
Export Administration Regulations, Section 740.13) for both object
code and source code.

The following provides more details on the included cryptographic
software:
  Hadoop Core uses the SSL libraries from the Jetty project written
by mortbay.org.
dase-dis@ecnu04:~/hadoop-2.10.1$
```

图 52: 布式环境 HDFS 下载文件到本地

最后停止 HDFS 服务停止命令查看进程，验证是否成功停止服务使用jps 在主节点和从节点查看进程，不再出现 NameNode,DataNode,SecondaryNameNode，则表示服务停止

主节点关闭后，从节点也应当被完全关闭

```
dase-dis@ecnu02:~$ jps  
19119 Jps
```

图 53: ecnu02 已经被关闭

```
dase-dis@ecnu03:~$ jps  
15810 Jps  
dase-dis@ecnu03:~$
```

图 54: ecnu03 已经被关闭

### 分布式 MapReduce

同样地，在开始实验前需要修改相关配置文件

在主节点 ecnu01 进行修改 mapred-site.xml,yarn-site.xml，然后将配置文件拷贝到从节点和客户端。（这里的步骤的 HDFS 的非常类似，就不再赘述了）

```
dase-dis@ecnu01:~$ ~/hadoop-2.10.1/sbin/start-yarn.sh  
starting yarn daemons  
starting resourcemanager, logging to /home/dase-dis/hadoop-2.10.1/logs/yarn-dase-  
-dis-resourcemanager-ecnu01.out  
ecnu03: starting nodemanager, logging to /home/dase-dis/hadoop-2.10.1/logs/yarn-  
dase-dis-nodemanager-ecnu03.out  
ecnu02: starting nodemanager, logging to /home/dase-dis/hadoop-2.10.1/logs/yarn-  
dase-dis-nodemanager-ecnu02.out  
dase-dis@ecnu01:~$ ~/hadoop-2.10.1/sbin/mr-jobhistory-daemon.sh start historyser  
ver  
starting historyserver, logging to /home/dase-dis/hadoop-2.10.1/logs/mapred-dase-  
-dis-historyserver-ecnu01.out  
dase-dis@ecnu01:~$ jps  
9511 Jps  
9179 ResourceManager  
9470 JobHistoryServer
```

图 55: ecnu01 启动 Yarn 服务

```
dase-dis@ecnu02:~/ssh$ jps  
4532 DataNode  
5158 NodeManager  
5306 Jps
```

图 56: ecnu02 使用 jps 查看进程状态

```
dase-dis@ecnu03:~$ jps  
13016 DataNode  
12793 NodeManager  
13131 Jps  
dase-dis@ecnu03:~$
```

图 57: ecnu03 使用 jps 查看进程状态

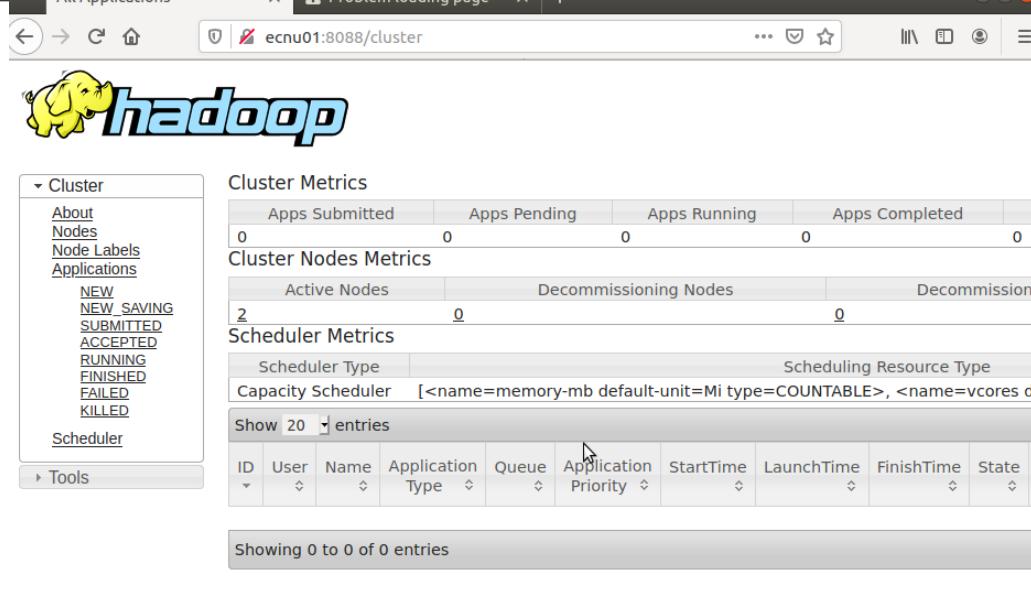


图 58: 使用 Web UI 查看服务状态

```
dase-dis@ecnu04:~/hadoop-2.10.1$ cd ~/hadoop-2.10.1
dase-dis@ecnu04:~/hadoop-2.10.1$ ./bin/hdfs dfs -mkdir -p input/grep
dase-dis@ecnu04:~/hadoop-2.10.1$ ./bin/hdfs dfs -put etc/hadoop/* input/grep
dase-dis@ecnu04:~/hadoop-2.10.1$ ./bin/hdfs dfs -rm -r output/grep
rm: `output/grep': No such file or directory
```

图 59: 客户端 ecnu04 提交 jar 命令

```
dase-dis@ecnu04:~/hadoop-2.10.1$ ./bin/yarn jar share/hadoop/mapreduce/hadoop-mapreduce-examples-2.10.1.jar grep input/grep output/grep 'dfs[a-z.]+'
21/03/25 21:09:05 INFO client.RMProxy: Connecting to ResourceManager at ecnu01/10.24.21.138:8032
21/03/25 21:09:06 INFO input.FileInputFormat: Total input files to process : 30
21/03/25 21:09:06 INFO mapreduce.JobSubmitter: number of splits:30
21/03/25 21:09:06 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1616671487701_0001
21/03/25 21:09:06 INFO conf.Configuration: resource-types.xml not found
21/03/25 21:09:06 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
21/03/25 21:09:06 INFO resource.ResourceUtils: Adding resource type - name = memory-mb, units = Mi, type = COUNTABLE
21/03/25 21:09:06 INFO resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE
21/03/25 21:09:07 INFO impl.YarnClientImpl: Submitted application application_1616671487701_0001
21/03/25 21:09:07 INFO mapreduce.Job: The url to track the job: http://ecnu01:8088/proxy/application_1616671487701_0001/
21/03/25 21:09:07 INFO mapreduce.Job: Running job: job_1616671487701_0001
21/03/25 21:09:13 INFO mapreduce.Job: Job job_1616671487701_0001 running in uber mode : false
21/03/25 21:09:13 INFO mapreduce.Job: map 0% reduce 0%
21/03/25 21:09:21 INFO mapreduce.Job: map 13% reduce 0%
21/03/25 21:09:22 INFO mapreduce.Job: map 33% reduce 0%
21/03/25 21:09:28 INFO mapreduce.Job: map 43% reduce 0%
21/03/25 21:09:29 INFO mapreduce.Job: map 53% reduce 0%
21/03/25 21:09:30 INFO mapreduce.Job: map 77% reduce 0%
21/03/25 21:09:32 INFO mapreduce.Job: map 80% reduce 0%
21/03/25 21:09:33 INFO mapreduce.Job: map 87% reduce 0%
21/03/25 21:09:34 INFO mapreduce.Job: map 97% reduce 0%
```

图 60: 客户端 ecnu04 MapReduce 执行中

```
base-dis@ecnu04:~/hadoop-2.10.1$ ./bin/hdfs dfs -cat output/grep/p*
6    dfs.audit.logger
4    dfs.class
3    dfs.logger
3    dfs.server.namenode.
2    dfs.audit.log.maxbackupindex
2    dfs.period
2    dfs.audit.log.maxfilesize
1    dfs.replication
1    dfs.log
1    dfs.file
1    dfs.datanode.name.dir
1    dfs.servers
1    dfsadmin
1    dfsmetrics.log
1    dfs.namenode.name.dir
```

图 61: 客户端 ecnu04 MapReduce 执行完毕

在客户端提交的任务运行过程中，在从节点可以看到相应启动的进程

```
dase-dis@ecnu03:~$ jps
14594 YarnChild
14579 YarnChild
14550 YarnChild
14615 YarnChild
13016 DataNode
14537 YarnChild
12793 NodeManager
14859 Jps
14571 YarnChild
14619 YarnChild
14540 YarnChild
dase-dis@ecnu03:~$
```

图 62: 从节点使用 jps 查看任务执行过程中的进程状态

随后，在客户端提交 WordCount 示例任务

```
dase-dis@ecnu04:~/hadoop-2.10.1$ ./bin/yarn jar share/hadoop/mapreduce/hadoop-mapreduce-examples-2.10.1.jar wordcount input/pd.train output/wordcount
21/03/25 21:16:50 INFO client.RMProxy: Connecting to ResourceManager at ecnu01/10.24.21.138:8032
21/03/25 21:16:50 INFO input.FileInputFormat: Total input files to process : 1
21/03/25 21:16:50 INFO mapreduce.JobSubmitter: number of splits:17
21/03/25 21:16:50 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1616671487701_0003
21/03/25 21:16:50 INFO conf.Configuration: resource-types.xml not found
21/03/25 21:16:50 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
21/03/25 21:16:50 INFO resource.ResourceUtils: Adding resource type - name = memory-mb, units = Mi, type = COUNTABLE
21/03/25 21:16:50 INFO resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE
21/03/25 21:16:51 INFO impl.YarnClientImpl: Submitted application application_1616671487701_0003
21/03/25 21:16:51 INFO mapreduce.Job: The url to track the job: http://ecnu01:8088/proxy/application_1616671487701_0003/
21/03/25 21:16:51 INFO mapreduce.Job: Running job: job_1616671487701_0003
21/03/25 21:16:56 INFO mapreduce.Job: Job job_1616671487701_0003 running in uber mode : false
21/03/25 21:16:56 INFO mapreduce.Job: map 0% reduce 0%
21/03/25 21:17:16 INFO mapreduce.Job: map 3% reduce 0%
21/03/25 21:17:17 INFO mapreduce.Job: map 8% reduce 0%
21/03/25 21:17:18 INFO mapreduce.Job: map 9% reduce 0%
21/03/25 21:17:21 INFO mapreduce.Job: map 10% reduce 0%
21/03/25 21:17:22 INFO mapreduce.Job: map 11% reduce 0%
21/03/25 21:17:23 INFO mapreduce.Job: map 14% reduce 0%
21/03/25 21:17:24 INFO mapreduce.Job: map 15% reduce 0%
21/03/25 21:17:27 INFO mapreduce.Job: map 16% reduce 0%
21/03/25 21:17:28 INFO mapreduce.Job: map 18% reduce 0%
21/03/25 21:17:29 INFO mapreduce.Job: map 19% reduce 0%
21/03/25 21:17:34 INFO mapreduce.Job: map 20% reduce 0%
21/03/25 21:17:35 INFO mapreduce.Job: map 23% reduce 0%
21/03/25 21:17:36 INFO mapreduce.Job: map 24% reduce 0%
```

图 63: 分布式模式下执行 WordCount 任务

```
dase-dis@ecnu01:~$ jps
9748 NameNode
11285 Jps
9179 ResourceManager
10061 SecondaryNameNode
9470 JobHistoryServer
```

图 64: 分布式模式下执行 WordCount: 主节点状态

```
dase-dis@ecnu01:~$ jps
9748 NameNode
11285 Jps
9179 ResourceManager
10061 SecondaryNameNode
9470 JobHistoryServer
dase-dis@ecnu02:~$ jps
17973 MRAppMaster
18024 Jps
16106 DataNode
15867 NodeManager
dase-dis@ecnu02:~$ jps
17973 MRAppMaster
18198 YarnChild
18135 YarnChild
18281 Jps
16106 DataNode
18155 YarnChild
15867 NodeManager
18188 YarnChild
18206 YarnChild
18175 YarnChild
dase-dis@ecnu02:~$
```

图 65: 分布式模式下执行 WordCount: 从节点状态

```
dase-dis@ecnu04:~$ jps
10689 RunJar
10778 Jps
```

图 66: 分布式模式下执行 WordCount: 客户端状态

然后查看 Map Reduce 程序的运行信息访问 Hadoop Map/Reduce Administration 界面，截图如下：

The screenshot shows the YARN application overview page. On the left, a sidebar lists navigation options: About, Nodes, Node Labels, Applications (with sub-options NEW, NEW\_SAVING, SUBMITTED, ACCEPTED, RUNNING, FINISHED, FAILED, KILLED), Scheduler, and Tools. The main content area displays cluster metrics: Apps Submitted (3), Apps Pending (0), Apps Running (1), and Apps Completed (2). Below this, the 'Cluster Nodes Metrics' section shows Active Nodes (2) and Decommissioning Nodes (0). The 'Scheduler Metrics' section details the Scheduler Type (Capacity Scheduler) and its configuration. A table lists 20 entries of submitted applications, each with columns for ID, User, Name, Application Type, Queue, Application Priority, and Start Time. The three visible entries are:

ID	User	Name	Application Type	Queue	Application Priority	Start Time
application_1616671487701_0003	dase-dis	word count	MAPREDUCE	default	0	Thu Mar 25 21:16:51 +0800 2021
application_1616671487701_0002	dase-dis	grep-sort	MAPREDUCE	default	0	Thu Mar 25 21:09:39 +0800 2021
application_1616671487701_0001	dase-dis	grep-search	MAPREDUCE	default	0	Thu Mar 25 21:09:06 +0800 2021

Showing 1 to 3 of 3 entries

图 67: MapReduce Administration

上述图片同样是访问 8088 端口可得到，可以看到所有的启动 Yarn 服务后提交的应用程序相关信息。

The screenshot shows the YARN job history page. On the left, a sidebar lists Application (About, Jobs) and Tools. The main content area displays the 'Retired Jobs' section with a table of 20 entries. The three visible entries are:

Submit Time	Start Time	Finish Time	Job ID	Name	User	Queue	State
2021.03.25 21:16:51 CST	2021.03.25 21:16:54 CST	2021.03.25 21:18:58 CST	job_1616671487701_0003	word count	dase-dis	default	SUCCEEDED
2021.03.25 21:09:39 CST	2021.03.25 21:09:48 CST	2021.03.25 21:09:50 CST	job_1616671487701_0002	grep-sort	dase-dis	default	SUCCEEDED
2021.03.25 21:09:06 CST	2021.03.25 21:09:11 CST	2021.03.25 21:09:37 CST	job_1616671487701_0001	grep-search	dase-dis	default	SUCCEEDED

Showing 1 to 3 of 3 entries

图 68: MapReduce History

MapReduce 的 JobHistory 可以通过访问端口 19888 查看，里面记录着程序的历史运行记录。

```
dase-dis@ecnu01:~$ cd ~/hadoop-2.10.1/logs
dase-dis@ecnu01:~/hadoop-2.10.1/logs$ ls
hadoop-dase-dis-namenode-ecnu01.log
hadoop-dase-dis-namenode-ecnu01.out
hadoop-dase-dis-namenode-ecnu01.out.1
hadoop-dase-dis-namenode-ecnu01.out.2
hadoop-dase-dis-secondarynamenode-ecnu01.log
hadoop-dase-dis-secondarynamenode-ecnu01.out
hadoop-dase-dis-secondarynamenode-ecnu01.out.1
hadoop-dase-dis-secondarynamenode-ecnu01.out.2
hadoop-dase-dis-secondarynamenode-ecnu01.out.3
hadoop-dase-dis-secondarynamenode-ecnu01.out.4
hadoop-dase-dis-secondarynamenode-ecnu01.out.5
mapred-dase-dis-historyserver-ecnu01.log
mapred-dase-dis-historyserver-ecnu01.out
SecurityAuth-dase-dis.audit
yarn-dase-dis-resourcemanager-ecnu01.log
yarn-dase-dis-resourcemanager-ecnu01.out
```

图 69: MapReduce 任务日志

最后，主节点关闭 MapReduce 服务后，从节点不再有任何相关进程

```
dase-dis@ecnu03:~$ jps
15810 Jps
dase-dis@ecnu03:~$
```

图 70: 从节点不再有任何相关进程

## Part 5

思考题

### Section 1

根据分布式部署的 HDFS 的 Web UI，能看到 HDFS 架构中的哪些角色？每个角色对应所在主机的 IP 地址是多少？

In operation								
Show 25 entries	Search:							
Node	Http Address	Last contact	Last Block Report	Capacity	Blocks	Block pool used	Version	
✓ecnu02:50010 (10.23.244.158:50010)	http://ecnu02:50075	2s	10m	49.09 GB	62	2.05 GB (4.17%)	2.10.1	
✓ecnu03:50010 (10.24.21.163:50010)	http://ecnu03:50075	2s	19m	49.09 GB	62	2.05 GB (4.17%)	2.10.1	
Showing 1 to 2 of 2 entries								
				Previous		1	Next	

图 71: 思考题 1

Answer: 通过 Web UI 中可以查看到 HDFS 两个子节点角色，在 Node 下方可以看到其 IP 地址

### Section 2

根据分布式部署的 Yarn 的 Web UI，能看到 Yarn 架构中的哪些角色？每个角色对应所在主机的 IP 地址是多少？

Answer: 通过 Web UI 中可以查看到 Yarn 架构中两个子节点角色的状况，其 IP 地址似乎并不在 Web UI 中可以查看

## 华东师范大学数据科学与工程学院学生实验报告

经过一些资料查询和同学之间的交流发现，如果在实验一开始不执行 `sethostnamectl`，直接使用 ip 命名的用户做这个实验的话，在 `active nodes` 点进去可以看到当前机器和没有改名的从节点的 ip 地址，为了截一次图一次性开四台云主机太消耗资源，所以就没有复现。

The screenshot shows the Hadoop Web UI interface. On the left, there's a sidebar with navigation links like 'Nodes of the cluster', 'Cluster Metrics', 'Cluster Nodes Metrics', 'Scheduler Metrics', and 'Tools'. The main content area displays 'Cluster Metrics' with tables for Apps Submitted (1), Apps Pending (0), Apps Running (0), Apps Completed (1), Containers Running (0), and Used Resources (<memory:0, vCores:0>). Below this are sections for 'Cluster Nodes Metrics' (Active Nodes: 2, Decommissioning Nodes: 0, Decommissioned Nodes: 0, Lost Nodes: 0, Unhealthy: 0) and 'Scheduler Metrics' (Capacity Scheduler: <name=memory-mb default-unit=Mi type=COUNTABLE>, <name=vcores default-unit= type=COUNTABLE>). The bottom part shows a table of nodes with columns: Node Labels, Rack, Node State, Node Address, Node HTTP Address, Last health-update, Health-report, Containers, and Mem Used. Two nodes are listed: one with IP ecnu02:18639 and another with IP ecnu03:26579.

图 72: 思考题 2

Answer: 通过 Web UI 中可以查看到两个子节点的状况，在 Node 下方可以看到其 IP 地址

### Section 3

能否从 Hadoop Map/Reduce Administration 的 Web UI 中获取某个 MapReduce 应用程序的总执行时间？请简要说明

Answer: 通过 Web UI 中的 Retired Jobs 可以查看，时间为 Start Time 和 Finish Time 之差

Retired Jobs							
Show 20 entries	Submit Time	Start Time	Finish Time	Job ID	Name	User	Queue
	2021.03.29 21:30:27 CST	2021.03.29 21:30:32 CST	2021.03.29 21:30:33 CST	job_1617024132543_0001	WordCount	dase-dis	default
	2021.03.25 21:16:51 CST	2021.03.25 21:16:54 CST	2021.03.25 21:18:58 CST	job_1616671487701_0003	word count	dase-dis	default
	2021.03.25 21:09:39 CST	2021.03.25 21:09:48 CST	2021.03.25 21:09:50 CST	job_1616671487701_0002	grep-sort	dase-dis	default
	2021.03.25 21:09:06 CST	2021.03.25 21:09:11 CST	2021.03.25 21:09:37 CST	job_1616671487701_0001	grep-search	dase-dis	default
Showing 1 to 4 of 4 entries							

图 73: 思考题 3

## Part 6

### 实验总结

最后根据课上老师的建议尝试记时了一下伪分布式部署情况下的时间开销，可以看到比实验报告最初的 9min 快了不少（虚拟机有点卡，请原谅我直接拍了照）

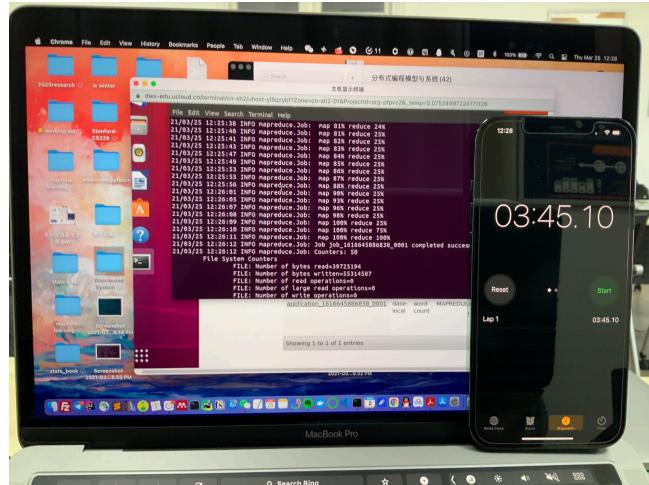


图 74: Pseudo-Distributed Mode

从网上搜集了一些资料，结合实验自己消化吸收后对 Hadoop2.x 的特性以及通用资源管理系统 Yarn 有了基本的认识。

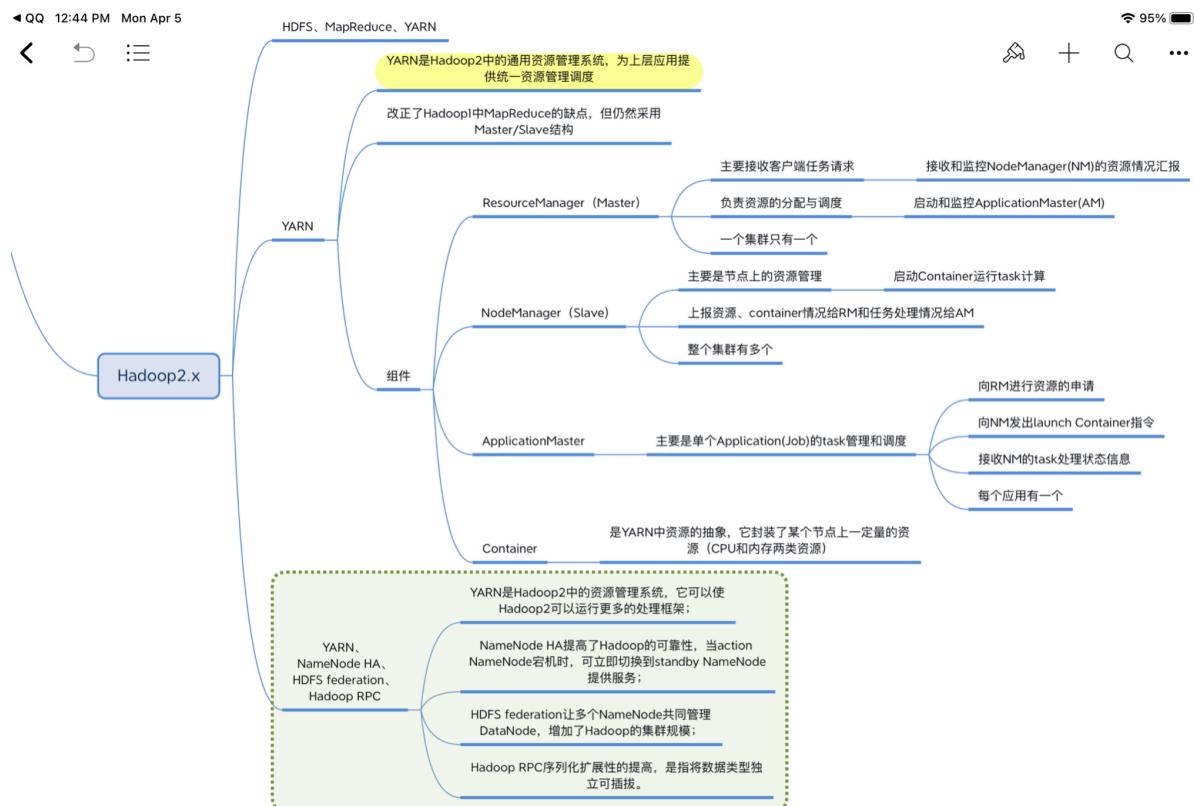


图 75: Hadoop2.x 总结

Hadoop2.0 解决了 Hadoop1.x 版本的一些问题如：

- 解决了 NameNode 单点故障问题。
- 解决 NameNode 内存压力过大难以扩展问题。
- 解决 JobTracker 单点故障问题。
- 解决 JobTracker 访问压力过大问题。
- 解决对 MapReduce 之外的框架支持问题。