## Lab 01: Setup Environment for run Hadoop System

**Group name:** Bigdata boi
**List of members:**
- 19127276 – Nguyen Dang Thi
- 19127592 – Le Minh Tri
- 19127639 – Phan Vi Giai
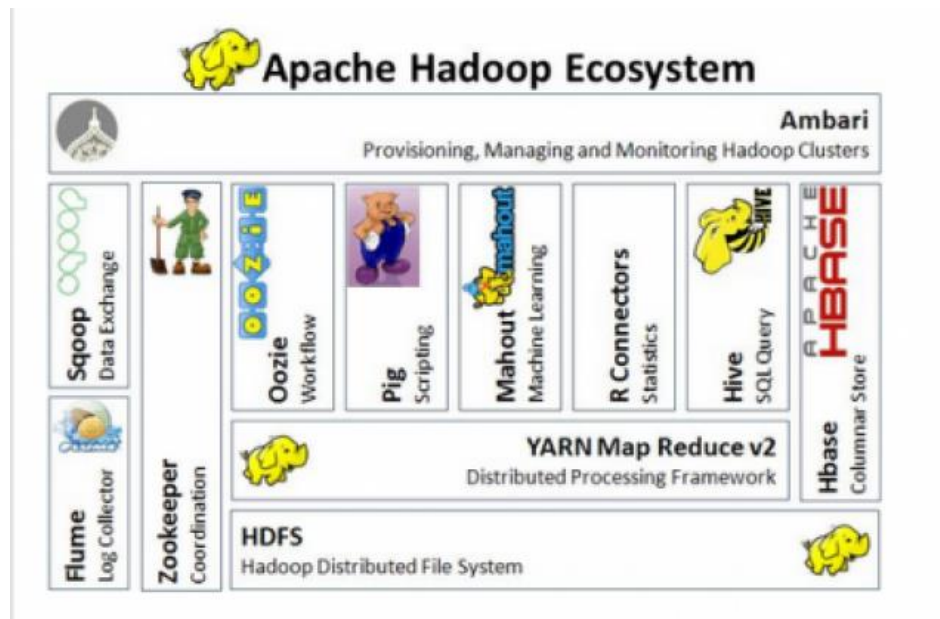- 19127642 – Vo Nhat Huy

# Table content

## I. Introduction to Hadoop Ecosystem

The Apache Hadoop software library is an open-source framework maintained by the ASF that allows for the distributed processing of large data sets across clusters of computers using simple programming models. It is written in Java and optimized for massive amounts of data through distribution.

The Apache Hadoop framework is designed to scale up from single machine to thousands of machines which is integrated a Hadoop cluster, each offering local computation and storage. Further rely on hardware to possess high-availability, the architecture of framework is designed to resilient to failures at the application layer, so that delivering a high-availability service on top of a cluster of computers.

Beside the main purpose of data storage, several services such as processing, indexing and manipulating can be easily integrated on top a Hadoop cluster through Yet Another Resource Negotiator (>= Hadoop 2.0). These services help to solve any problems on Big Data in the biggest business companies.



*Hadoop ecosystem (https://opensource.com/life/14/8/intro-apache-hadoop-big-data)*

## II. Installation

### 1. Setup Hadoop run Pseudo-Distributed operation mode

**Overview:**

*"Hadoop can be installed in 3 different modes: Standalone mode, Pseudo-Distributed mode and Fully-Distributed mode.*

*...*

*Pseudo-distributed mode is also known as a single-node cluster where both NameNode and DataNode will reside on the same machine."*

https://medium0.com/analytics-vidhya/hadoop-single-node-cluster-setup-b11b957681f2

With Pseudo-distributed mode, Hadoop daemon is allowed to run as a single Java process.

**Requirements and versions:**

- Java engine version 8
- Hadoop version 2.10.1

**Follow step by step installation:**

*Note: Reports on installations in this section performed on Linux (Ubuntu 20.04 LTS)*

**Step 1. Install OpenJDK**

Hadoop is written in Java. So, it would require the JRE (Java Runtime Environment) at the very least. However, Hadoop is a program,Ong framework (for running map/reduce jobs or supporting YARN applications across a cluster).

To install OpenJDK, refer here: https://openjdk.java.net/install/

On ubuntu, we can install OpenJDK by following steps:

Update system before initiating a new installation:



Type the following command in your terminal to install OpenJDK 8:



Check current Java version:



**Step 2. Set Up a Non-Root User for Hadoop Environment**

**Install OpenSSH**

**Create Hadoop user**

Use **adduser** to create a new Hadoop user:



**Enable Passwordless SSH for Hadoop User**

## Step 3. Download and Install Hadoop

Visit the Apache Hadoop homepage and select a Hadoop version to install:
https://hadoop.apache.org/releases.html



You can click on the link contains tar.gz file or copy link and run below command to download:



Move to the folder containing the downloaded file and extract it



Remember to replace with the version you want to download accordingly

## Step 4. Setup Hadoop run Pseudo-Distributed operation mode

To setup Hadoop for Pseudo-distributed mode, set up the following file:

- .bashrc
- hadoop-env.sh
- yarn-site.xml
- core-site.xml
- hdfs-site.xml
- mapred-site.xml

### .bashrc

Using nano for configure the .bashrc shell. You can use another text editor to do the same.

Or



Put these lines into the end of .bashrc

```
#Hadoop Related Options
export HADOOP_HOME=<path to installed hadoop folder>
export HADOOP_INSTALL=$HADOOP_HOME
export HADOOP_MAPRED_HOME=$HADOOP_HOME
export HADOOP_COMMON_HOME=$HADOOP_HOME
export HADOOP_HDFS_HOME=$HADOOP_HOME
export YARN_HOME=$HADOOP_HOME
export HADOOP_COMMON_LIB_NATIVE_DIR=$HADOOP_HOME/lib/native export
PATH=$PATH:$HADOOP_HOME/sbin:$HADOOP_HOME/bin
export HADOOP_OPTS="-Djava.library.path=$HADOOP_HOME/lib"
```



*.bashrc*

<u>Note:</u> You have to edit the path to the hadoop folder (HADOOP_HOME) based on your path.

**hadoop-env.sh**

The hadoop-env.sh file is a file that contains some environment variable settings used by Hadoop. It is also a prerequisite for configuring settings related to YARN, DHFS and MapReduce.

In the hadoop-env.sh file, we need to uncomment the **$JAVA_HOME** variable and add the full path to the OpenJDK installation on your system. To check the correct path, you can refer here:
https://www.baeldung.com/find-java-home



The file editing process will go from like this:

```
# Generic settings for HADOOP
###

# Technically, the only required environment variable is JAVA_HOME.
# All others are optional.  However, the defaults are probably not
# preferred.  Many sites configure these options outside of Hadoop,
# such as in /etc/profile.d

# The java implementation to use. By default, this environment
# variable is REQUIRED on ALL platforms except OS X!
# export JAVA_HOME=

# Location of Hadoop.  By default, Hadoop will attempt to determine
# this location based upon its execution path.
# export HADOOP_HOME=

# Location of Hadoop's configuration information.  i.e., where this
# file is living. If this is not defined, Hadoop will attempt to
# locate it based upon its execution path.
#
```

*hadoop-env.sh*

To this:

```
# preferred.  Many sites configure these options outside of Hadoop,
# such as in /etc/profile.d

# The java implementation to use. By default, this environment
# variable is REQUIRED on ALL platforms except OS X!
export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64

# Location of Hadoop.  By default, Hadoop will attempt to determine
# this location based upon its execution path.
# export HADOOP_HOME=
```

*hadoop-env.sh*

## yarn-site.xml

The yarn-site.xml file is used to define YARN related settings. It contains configurations for Node Manager, Resource Manager, Containers and Application Master.

Open the yarn-site.xml file in a text editor, then append the following lines into the end of file:

```
<configuration>
<property>
    <name>yarn.nodemanager.aux-services</name>
    <value>mapreduce_shuffle</value>
</property>
<property>
    <name>yarn.nodemanager.aux-services.mapreduce.shuffle.class</name>
    <value>org.apache.hadoop.mapred.ShuffleHandler</value>
</property>
<property>
    <name>yarn.resourcemanager.hostname</name>
    <value>127.0.0.1</value>
</property>
<property>
    <name>yarn.acl.enable</name>
    <value>0</value>
</property>
<property>
    <name>yarn.nodemanager.env-whitelist</name>
    <value>JAVA_HOME,HADOOP_COMMON_HOME,HADOOP_HDFS_HOME,HADOOP_CONF_DIR,CLASSPAT
    H_PERPEND_DISTCACHE,HADOOP_YARN_HOME,HADOOP_MAPRED_HOME</value>
```

```
</property>
</configuration>
```



*yarn-site.xml*

## core-site.xml

The core-site. xml file informs Hadoop daemon where NameNode runs in the cluster. It contains the configuration settings for Hadoop Core such as I/O settings that are common to HDFS and MapReduce.

Open the core-site.xml in text editor, then add the following configuration to override the default values for the temporary directory and add your HDFS URL to replace the default local file system setting:

```
<configuration>
<property>
<name>hadoop.tmp.dir</name>
<value>file:/home/hadoop/tmp</value>
</property>
<property>
<name>fs.defaultFS</name>
<value>hdfs://master:9000</value>
</property>
</configuration>
```

You can customize the name and value, remember to create a directory in the location you specified for your temporary data.

*core-site.xml*

## hdfs-site.xml

The hdfs-site.xml file contains the configuration settings for HDFS daemons; the NameNode, the Secondary NameNode, and the DataNodes. ... xml to specify default block replication and permission checking on HDFS. The actual number of replications can also be specified when the file is created.

Additionally, the default dfs.replication value of 3 needs to be changed to 1 to match the single node setup.

Open the hdfs-site.xml file for editing:

Add the following configuration to the file and adjust the NameNode and DataNode directories to your custom locations:

```
<configuration>
 <property>
  <name>dfs.replication</name>
  <value>2</value>
 </property>
 <property>
  <name>dfs.namenode.name.dir</name>
  <value>file:/home/hadoop/hdfs/namenode</value>
 </property>
 <property>
  <name>dfs.datanode.data.dir</name>
  <value>file:/home/hadoop/hdfs/datanode</value>
 </property>
 <property>
  <name>dfs.namenode.secondary.http-address</name>
  <value>master:9001</value>
 </property>
 <property>
  <name>dfs.namenode.rpc-address</name>
  <value>master:9000</value>
 </property>
 <property>
  <name>dfs.webhdfs.enabled</name>
  <value>true</value>
 </property>
</configuration>
```

*hdfs-site.xml*

## mapred-site.xml

The mapred-site. xml file contains the configuration settings for MapReduce daemons; the job tracker and the task-trackers.

Add the following configuration to change the default MapReduce framework name value to yarn:

```
<configuration>
<property>
    <name>mapreduce.framework.name</name>
    <value>yarn</value>
</property>
</configuration>
```



*mapred-site.xml*

**Format HDFS NameNode**

Before starting Hadoop service for the fist time, we need to format HDFS NameNode





*format HDFS NameNode*

**Step 5. Start Hadoop Cluster**

Open terminal in sbin folder inside hadoop folder was installed before.



*sbin folder location*

Execute 2 files start-dfs.sh and start-yarn.sh, it may take a few moments.



*Execute start-dfs.sh*

*Execute start-yarn.sh*

Use *jps* command to check if all the daemons are active and running as Java processes:



## 2. Setup Hadoop multi-nodes run on Docker environment

**Overview:**

*Docker is an open-source containerization platform. It enables developers to package applications into containers—standardized executable components combining application source code with the operating system (OS) libraries and dependencies required to run that code in any environment. Containers simplify delivery of distributed applications and have become increasingly popular as organizations shift to cloud-native development and hybrid multicloud environments.*

Based on the idea of Docker running several containers like virtual machines, we take advantage of this to simulate a Hadoop cluster in a single computer. Our virtual Hadoop cluster includes 3 containers corresponding to 1 master and 2 slaves.

To connect these containers (virtual machines), we need docker container management to handle this problem. Docker compose is the key solution – a tool for defining, running and connect multi-container via a virtual network.

**Full configurations for this installation can be found [there](there).**

**Requirements and versions:**

- Docker engine and docker compose:
    - https://docs.docker.com/desktop/windows/install/
    - https://docs.docker.com/compose/install/
- Java engine version 8
- Hadoop version 2.10.1

**Follow step by step installation:**

**Step 1. Install Docker environment**

- First, we need to install docker and docker compose environment through links in the requirement section
- After installing docker, we check whether it is installed by commands:
    - *docker version*
    - *docker compose version*

*docker is installed*



*docker compose is installed*

## Step 2. Setup folder structure



*Folder structure*

- **configurations**: The folder contains all configurations for Hadoop and other services if needed
- **docker-compose.yml**: The file defines our container configurations such as container name, ram, network ip, command, ports, ... for each container.
- **Dockerfile**: The file contains all commands that set up our Hadoop and other services like downloading package, decompress, add configurations to Hadoop, ...
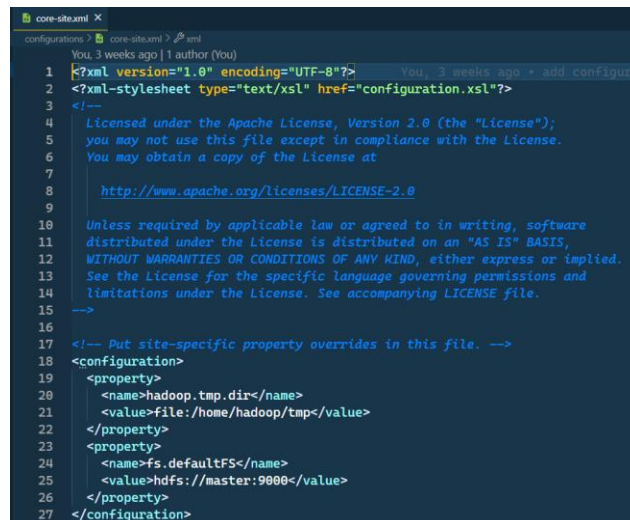
## Step 3. Create files and configurations



*configurations folder*

**configurations** folder contains several files for running Hadoop, that all files we needed to create:

*core-site.xml*
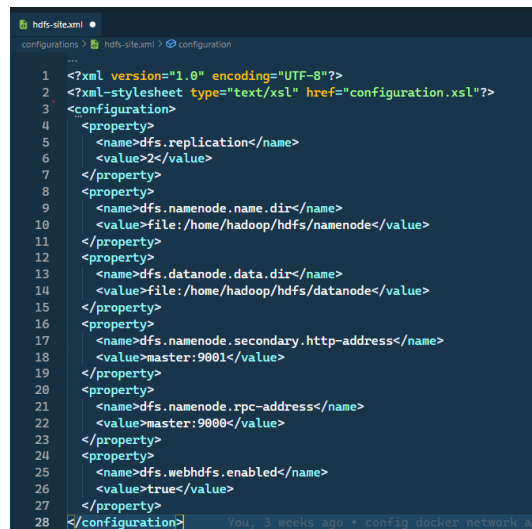
- **core-site.xml:** The file defines temporary directory for Hadoop and master host
  - **hadoop.tmp.dir:** direct to temporary directory which is manually created
  - **fs.defaultFS:** specify an address that allows dfs commands without providing full site name in the command



*hdfs-site.xml*

- **hdfs-site.xml:** The file configures Hadoop DFS properties like replication factors, directory for namenode, directory for datanode, and other properties
  - **dfs.replication:** we set 2 for 2 datanodes
  - **dfs.namenode.name.dir:** working directory for namenode
  - **dfs.datanode.data.dir:** working directory for datanode
  - **dfs.namenode.secondary.http-address:** The secondary namenode http server address and port.
  - **dfs.namenode.rpc-address:** RPC address that handles all client's requests
  - **dfs.webhdfs.enabled:** Enable WebHDFS (REST API) in Namenodes and Datanodes

*mapred-site.xml*

- **mapred-site.xml:** The file use Yarn as Resource Manager for Map-Reduce application on Hadoop
  - **mapreduce.framework.name:** Specify *yarn* value for framework name as using Map-Reduce version 2



*yarn-site.xml*

- **yarn-site.xml:** The file configures some property of Yarn ResourceManager
  - **yarn.nodemanager.aux-services:** use default shuffle and sort of Map-Reduce
  - **yarn.nodemanager.axu-services.mapreduce.shuffle.class:** The auxiliary service class to use as default value
  - **yarn.resourcemanager.address:** The address of the applications manager interface in the ResourceManager
  - **yarn.resourcemanager.scheduler.adress:** The address of the scheduler interface

- o **yarn.resourcemanager.resource-tracker.address:** The address of the resource-tracker interface
- o **yarn.resourcemanager.admin.address:** The address of the admin interface
- o **yarn.resourcemanager.webapp.address:** The http address of the RM web application



*slaves*

- **slaves:** Define name of slaves which are slave1 for datanode 1 and slave2 for datanode 2



*start-dfs.sh*

- **start-dfs.sh:** That is just a default file when we install Hadoop locally. The following command starts demons like starting namenode, secondary namenode (if any), datanodes.

*stop-dfs.sh*

- **stop-dfs.sh:** That is also just a default file when we install Hadoop locally. The following commands stop demons like stopping namenode, secondary namenode (if any), datanodes.



*start-yarn.sh*

- **start-yarn.sh:** The file is starting Yarn as a Resource Manager and Node Manager.



*stop-yarn.sh*

- **stop-yarn.sh:** The file is a default file when installing Hadoop. The following commands are stopping Yarn daemons.

**Step 4. Setup Dockerfile and docker-compose.yml**

We need to create 2 files that are **Dockerfile** and **docker-compose.yml**

First, **Dockerfile** contains all command needed to download package like Hadoop or something you want to integrate with Hadoop system, then setup folder like we installed Hadoop in local mode and add configuration files to Hadoop. All those commands are executed in containers.

```
Dockerfile  ×
Dockerfile > ⑨ FROM
        You, a week ago | 1 author (You)
   1    FROM ubuntu        You, 3 weeks ago • config docker network and deploy
   2    RUN apt update && apt install -y openssh-server openssh-client vim openjdk-8-jdk
```

- We use **ubuntu** operating system then install **openssh-client** for remote login with **SSH** protocol and **openjdk-8-jdk** for both the runtime environment and development kit of Java 8

```
   5    # SSH
   6    RUN ssh-keygen -t rsa -P '' -f ~/.ssh/id_rsa
   7    RUN cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys
   8    RUN chmod 0600 ~/.ssh/authorized_keys
```

- To enable passwordless **SSH** for Hadoop user
    o First, we need to generate an SSH key pair and define the location to store it
    o Copy value to the *authorized_keys* folder
    o Set the permissions for your user with the *chmod* command

```
  11    # CONFIG JAVA ENVINROMENT VARIBLES
  12    ENV JAVA_HOME /usr/lib/jvm/java-8-openjdk-amd64
```

- Configures global Java environment variables

```
  15    # HADOOP
  16    RUN wget https://archive.apache.org/dist/hadoop/common/hadoop-2.10.1/hadoop-2.10.1.tar.gz
  17    RUN tar -xzf hadoop-2.10.1.tar.gz
  18    RUN mv hadoop-2.10.1 usr/local/hadoop
  19    ENV HADOOP_HOME /usr/local/hadoop
  20    ENV CONF_DIR $HADOOP_HOME/etc/hadoop
  21    ENV PATH $HADOOP_HOME/sbin:$PATH
```

- Download Hadoop package
    o First, we download Hadoop package from repository
    o Decompress package
    o Move folder after decompressing to our specific location and rename it for shorthand
    o Configure some Hadoop environment variables such as $HADOOP_HOME, $CONF_DIR, $PATH in order to use these variables in short later

```
24    # HADOOP - CREATE DIRECTORY FOR STORING DOCUMENTS
25    RUN mkdir /home/hadoop /home/hadoop/hdfs
26    RUN mkdir /home/hadoop/tmp /home/hadoop/hdfs/namenode /home/hadoop/hdfs/datanode
27    RUN chmod 777 /home/hadoop/hdfs/namenode
28    RUN chmod 777 /home/hadoop/tmp
29    RUN chmod 777 /home/hadoop/hdfs/datanode
30
31    ADD configurations/start-dfs.sh $HADOOP_HOME/sbin
32    ADD configurations/stop-dfs.sh $HADOOP_HOME/sbin
33    ADD configurations/start-yarn.sh $HADOOP_HOME/sbin
34    ADD configurations/stop-yarn.sh $HADOOP_HOME/sbin
35    ADD configurations/core-site.xml $HADOOP_HOME/etc/hadoop/core-site.xml
36    ADD configurations/hdfs-site.xml $HADOOP_HOME/etc/hadoop/hdfs-site.xml
37    ADD configurations/mapred-site.xml $HADOOP_HOME/etc/hadoop/mapred-site.xml
38    ADD configurations/yarn-site.xml $HADOOP_HOME/etc/hadoop/yarn-site.xml
39    ADD configurations/slaves $HADOOP_HOME/etc/hadoop/slaves
40
41    ENV PATH $HADOOP_HOME/bin:$PATH
```

- Setup directory and configurations
  - o We create working spaces for namenode and datanodes
  - o Give these folders permission to read, change, delete and execute files with **chmod 777** mode
  - o We add all files which are created in Step 3 are files in **configurations folder** to Hadoop configurations

```
65    # FORMAT NAMENODE
66    ARG FORMAT_NAMENODE_COMMAND
67    RUN $FORMAT_NAMENODE_COMMAND
68    EXPOSE 22
```

- Finally, run format namenode command
  - o *hdfs namenode -format*

Second, **docker-compose.yml** file is used to configure our application's services. In this context, three containers are needed that are one master and two slaves. Here are the following steps:

```
3     slave1:
4       build:
5         context: .
6         shm_size: '2gb'
7       container_name: slave1
8       networks:
9         default:
10          ipv4_address: 172.10.0.3
11      extra_hosts:
12        - "master: 172.10.0.2"
13        - "slave2: 172.10.0.4"
14      command: bash -c  "hadoop-daemon.sh --config /usr/local/hadoop/etc/hadoop \
15            start datanode && yarn-daemon.sh --config /usr/local/hadoop/etc/hadoop \
16            start nodemanager && start-slave.sh master:7077 && tail -f /dev/null"
17      hostname: slave1
18      restart: always
```

- Configure for slave1 (datanode1):
  - o build
    - ▪ *context: .* : It's a path to the directory that contains **Dockerfile**. In our case, **Dockerfile** file exists in the current directory (in the same folder of the **docker-compose.yml** file).
    - ▪ *shm_size: '2gb'*: Set the size of the /dev/shm partition for this build's containers.
  - o container_name: Define *slave1* for container's name
  - o networks: Set virtual ip_address for containers to connect others
  - o extra_hosts: Add hostname mappings at build-time. Here are our two other containers such as master and slave2

- command: Client command to start Yarn, Hadoop daemons but not namenode
- hostname: Hostname is that computers on the network know each other and thus communicate between themselves. For *container slave1,* we set hostname value is *slave1*
- restart: Always restart container if started fail

```
19    slave2:
20      build:
21        context: .
22        shm_size: '2gb'
23      container_name: slave2
24      networks:
25        default:
26          ipv4_address: 172.10.0.4
27      extra_hosts:
28        - "master: 172.10.0.2"
29        - "slave1: 172.10.0.3"
30      command: bash -c  "hadoop-daemon.sh --config /usr/local/hadoop/etc/hadoop \
31              start datanode && yarn-daemon.sh --config /usr/local/hadoop/etc/hadoop \
32              start nodemanager && start-slave.sh master:7077 && tail -f /dev/null"
33      hostname: slave2
34      restart: always
```

- Configure for slave2 (datanode2): The same configurations as slave1 (datanode1)

```
35    master:
36      build:
37        context: .
38        shm_size: '2gb'
39        args:
40          FORMAT_NAMENODE_COMMAND: hdfs namenode -format
41      container_name: master
42      networks:
43        default:
44          ipv4_address: 172.10.0.2
45      extra_hosts:
46        - "slave1: 172.10.0.3"
47        - "slave2: 172.10.0.4"
48      command: bash -c  "start-dfs.sh && start-yarn.sh && mr-jobhistory-daemon.sh \
49              start historyserver && start-master.sh && start-hbase.sh && tail -f /dev/null"
50      ports:
51        - 50070:50070
52        - 8088:8088
53        - 8080:8080
54        - 4040:4040
55        - 16010:16010
56      hostname: master
57      restart: always
```

- Configure for master (namenode): some attributes are same configurations of slave
  - build
    - args:  Specify arguments using for command in **Dockerfile**
      - *FORMAT_NAMENODE_COMMAND: hdfs namenode -format*  command use to format namenode before starting Hadoop services for the first time
  - command: Client command to start Yarn, Hadoop daemons but not datanode
  - ports: Activates the container to listen for specified ports from the world outside of the docker and accessible world inside docker
    - Resource Manager WebUI port
    - Namenode WebUI port
    - Datanode port

- networks: Create network for allowing to connect all containers together

## Step 5. Run Hadoop cluster on Docker environment

### Build & Run project

- Create subnets and bridge for *hadoop-network*
  - *docker network create --driver bridge hadoop-network --subnet=172.10.0.0/16*
- Build image and run our containers
  - *docker-compose up*
- Check whether containers are running which are daemons of Hadoop cluster
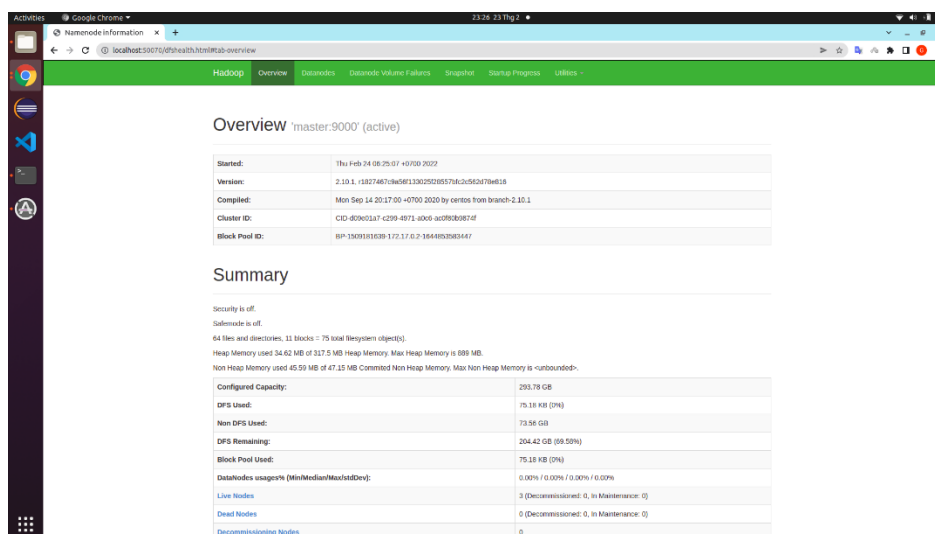  - *docker container ls*

If want to end sessions, run

- *docker-compose down*

### Manipulate with daemons

- Attach node of HDFS cluster
  - *docker exec -it master/slave1/slave2/... /bin/bash*
- Webapp daemons UI
  - For Resource Manager (YARN): http://localhost:8080/
  - For Master Management (Namenode): http://localhost:50070/



*Yarn Resource Manager*

*Namenode*

## III.   Demonstration functionalities

Section III focuses on demonstrating functionalities of Hadoop Distributed Filesystem.

This section has two parts:

- Hadoop command: To illustrate interaction between master and slaves in reading, writing files and other management functions of Hadoop system
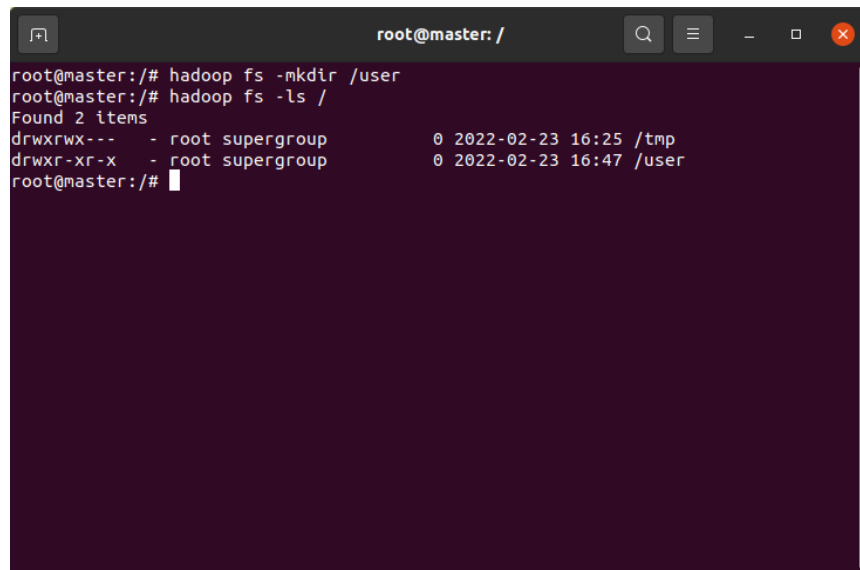- Framework operation: Writing Map-Reduce application that is running on Hadoop DFS

### 1.   Hadoop commands



*All the instrumental hotspots that the JVM is running in the system*

- jps: Java Virtual Machine Process Status Tool use to check all JVM process running in the system

*Create directory /user*

- -mkdir: Command use to create directory in Hadoop DFS



*List file in home location*

*List files in root location*
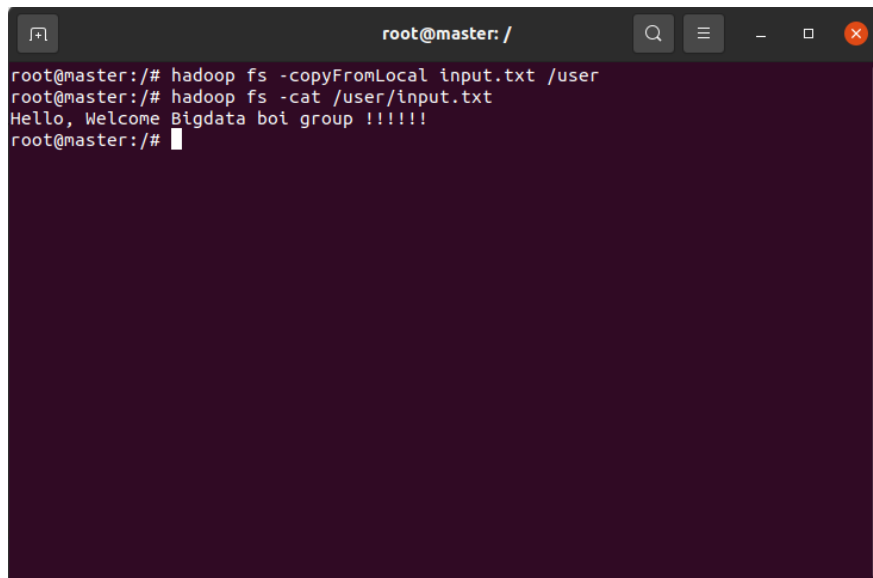
- -ls : List files in Hadoop DFS. This command has an optional parameter:
  - If no param: List files in home location
  - If / : List files in root location
  - If directory: List files in this directory



*Touch empty file input.txt in folder user*

- -touchz: Create empty file in specific directory

*Copy input.txt from master container to /user folder*

- -copyFromLocal: Copy file from local to destination in Hadoop DFS



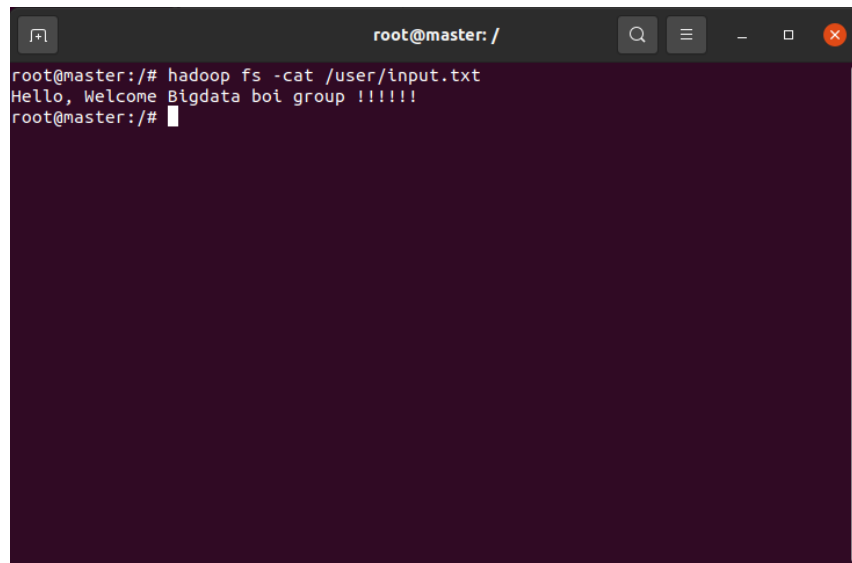*Copy file input.txt to local master container*

- -copyToLocal: Copy file from Hadoop DFS to local machine

*Concatenate input.txt*

- -cat: Concatenate file (read) in location



*Disk usage of root location*

- -du: Show total bytes of directories in specific location of Hadoop DFS

*Healthy report root location*

- fsck – directory: Used to check the health of directory in Hadoop DFS

## 2. Use Map-Reduce on Hadoop DFS with Pseudo-Distributed operation mode

**Overview:**

We take advantage of the distributed-computing framework which is designed to Map phase and Reduce phase calling MapReduce framework.

To demonstrate how MapReduce framework works on top distributed filesystem Hadoop, we play around with simple MapReduce application. The application is called **Sum Operation** that calculates summarization of even numbers and summarization of odd numbers in a very long sequence number stored in text file.

**Requirements and versions:**

- Linux system
- Java engine version 8
- Hadoop version 2.10.1

**Follow step by step implementation:**

**Step 1. Start Hadoop on local:**

- We need to start Hadoop on a local machine with Pseudo-Distributed mode which serves datanode, namenode and yarn daemons in one machine.

**Step 2. MapReduce implementation**

- To run a MapReduce application, we need a Java program for writing Map function and Reduce function to submit to Hadoop for creating Map task(s) and Reduce task(s)

*Our input file*

- Input file has a long sequence number. But to fit with the performance of our machine, input file contains 1000 rows, each row contains one number.

```java
package SumOperation;

import java.io.IOException;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

public class Map extends Mapper<LongWritable, Text, Text, IntWritable> {

    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
        String data[] = value.toString().split(",");
        for (String num : data) {
            int number = Integer.parseInt(num);
            if (number % 2 == 0) {
                context.write(new Text("EVEN"), new IntWritable(number));
            } else {
                context.write(new Text("ODD"), new IntWritable(number));
            }
        }
    }
}
```

*Map function*

- The strategy of map function is:
    o Convert number character to Hadoop internal integer type
    o If number % 2 == 0, label number key with *even*
    o Else, label number key with *odd*
    o We do that give each value has specific 'odd' or 'even' key and send these pairs to Reduce phase

```java
package SumOperation;

import java.io.IOException;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

public class Reduce extends Reducer<Text, IntWritable, Text, IntWritable> {
    public void reduce(Text key, Iterable<IntWritable> values, Context context)
            throws IOException, InterruptedException {
        int sum = 0;
        for (IntWritable value : values) {
            sum += value.get();
        }
        context.write(key, new IntWritable(sum));
    }
}
```

*Reduce function*

- The strategy of reduce function is:
    o With each pair has the same key, we accumulate these value of pairs
    o Return result for 'odd' and 'even' key with values

## Step 3. Build MapReduce application

- Export Hadoop classpath
    o *export HADOOP_CLASSPATH=$(hadoop classpath)*
    o *echo $HADOOP_CLASSPATH*
- Create compiled directory
    o *mkdir container*
- Compile application to JARs packages
    o *javac -cp ${HADOOP_CLASSPATH}:container/:. -d container/ Map.java*

- o *javac -cp ${HADOOP_CLASSPATH}:container/:. -d container/ Reduce.java*
- o *javac -cp ${HADOOP_CLASSPATH}:container/:. -d container/ Main.java*
- Create project JARs file
  - o jar -cvf so.jar -C container/ .

## Step 4. Run MapReduce application

- Create Hadoop input directory for storing our data
  - o *hadoop fs -mkdir inputdir*
- Copy our local data to Hadoop HDFS directory
  - o *hadoop fs -put inputdir/* inputdir*
- Run MapReduce job in Hadoop
  - o *hadoop jar so.jar SumOperation.Main inputdir outputdir*
- View MapReduce result
  - o *hadoop fs -cat outputdir/part-r-00000*

## Demonstration





*MapReduce Job report*

*Yarn Resource Manager manage MR Job*



*MapReduce application result*

## IV.     Self-evaluation

| LAB STATUS | Installation:<br>• Install Pseudo-distributed mode (Single Node Cluster) on local machine: FINISHED<br>• Install Fully-Distributed Mode (Multi-Node Cluster) on Docker environment: FINISHED<br>Test run process:<br>• Demo functionalities of Multi-Node Cluster in docker environments: FINISHED |
|---|---|
| PROBLEM | • Teamworking in remote<br>• Unfamiliar with Docker container, network, …<br>• The way connects each container together is complicated and even more difficult is how to configure Hadoop cluster on those containers<br>• Some Hadoop ecosystem documents are complicated<br>• Inconsistent OS working environment between members causes several unexpected errors<br>• Installations is flaky, it sometimes causes error, however when reinstalling, its run normally |
| SOLUTION | • Take advantage of online working tools to guarantee communication between members<br>• Research document from Docker website<br>• Try many times as well as search for solutions on the internet, YouTube, tutorial, …<br>• Research from many sources about Hadoop ecosystem<br>• Use containerization platform (Docker) to guarantee Hadoop's configuration in team<br>• Hold weekly meetings to review, evaluate each member's work or assist members who are in trouble |

## V.     References

[1] https://hadoop.apache.org/

[2] https://www.cloudera.com/products/open-source/apache-hadoop.html

[3] https://www.docker.com/get-started

[4] https://hadoop.apache.org/docs/stable/hadoop-project-dist/hadoop-common/SingleCluster.html

[5] https://www.geeksforgeeks.org/hdfs-commands/

[6] https://hadoop.apache.org/docs/r1.2.1/mapred_tutorial.html

[7] https://www.baeldung.com/find-java-home