

Установка и настройка **Java**

1. Install OpenJDK (JDK 8):

```
sudo apt-get update
sudo apt-get install openjdk-8-jdk
```

2. Verify installation:

```
java -version
openjdk version "1.8.0_275"
OpenJDK Runtime Environment (build 1.8.0_275-8u275-b01-0ubuntu1~20.04-b01)
OpenJDK 64-Bit Server VM (build 25.275-b01, mixed mode)
```

2. SET JAVA_HOME and JRE_HOME:

```
sudo vi /etc/environment

JAVA_HOME="/usr/lib/jvm/java-8-openjdk-amd64"

JRE_HOME="/usr/lib/jvm/java-8-openjdk-amd64/jre"
```

Настройка пользователя Hadoop

1. Create User:

sudo adduser hadoop
sudo passwd hadoop

2. Switch To User:

sudo su hadoop

3. Switch Back To Root user:

exit

Настройка **SSH** (требуется для компонентов **Hadoop**)

1. Install SSH and PDSH:

```
sudo apt-get install ssh pdsh
```

2. Create Private/Public Keypair for hadoop user (without passphrase):

```
sudo su hadoop
cd
ssh-keygen -t rsa -N "" -f /home/hadoop/.ssh/id_rsa
```

3. Add Public Key To Authorized Keys file (to enable passwordless ssh login)

```
cat /home/hadoop/.ssh/id_rsa.pub >> /home/hadoop/.ssh/authorized_keys
chmod 0600 /home/hadoop/.ssh/authorized_keys
```

Настройка **SSH** (требуется для компонентов **Hadoop**)

4. Check If SSH Is Working

```
hadoop@big-data:~$ ssh localhost
The authenticity of host 'localhost (127.0.0.1)' can't be established.
ECDSA key fingerprint is SHA256:YEUFliBVczkz2rvKWnYU9hB2ix2jnhBqLlbsJQfuBpE.
Are you sure you want to continue connecting (yes/no)? yes
Warning: Permanently added 'localhost' (ECDSA) to the list of known hosts.
Welcome to Ubuntu 18.04.3 LTS (GNU/Linux 4.15.0-1044-gcp x86 64)
 * Management: https://landscape.canonical.com
                  https://ubuntu.com/advantage
 System information as of Sat Oct 12 15:01:56 UTC 2019
 System load: 0.0
 Usage of /: 5.8% of 28.90GB Users logged in:
                     IP address for ens4: 10.156.0.6
 Memory usage: 2%
 Swap usage: 0%
30 packages can be updated.
17 updates are security updates.
Last login: Sat Oct 12 14:49:27 2019 from 80.144.211.195
hadoop@big-data:~$ exit
Connection to localhost closed.
hadoop@big-data:~$
```

Установка **Hadoop**

1. Download Hadoop (v3.1.1):

```
wget https://archive.apache.org/dist/hadoop/common/hadoop-3.1.2/hadoop-3.1.2.tar.gz
```

2. Extract Binaries:

```
tar -xvzf hadoop-3.1.2.tar.gz
```

3. Move Binaries:

mv hadoop-3.1.2 hadoop

1. Set Up **UNIX** Environment Variables

```
vi .bashrc
export HADOOP HOME=/home/hadoop/hadoop
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 PDSH RCMD TYPE=ssh
source .bashrc
```

2. Add **Hadoop** Environment Variables (hadoop-env.sh)

vi /home/hadoop/etc/hadoop/hadoop-env.sh

export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64

3. Set Up CORE Variables (core-site.xml)

4. Set Up HDFS Variables (hdfs-site.xml)

```
vi /home/hadoop/hadoop/etc/hadoop/hdfs-site.xml
```

```
<configuration>
     cproperty>
            <name>dfs.replication</name>
            <value>1
     </property>
     cproperty>
            <name>dfs.name.dir
            <value>file:///home/hadoop/hadoopdata/hdfs/namenode</value>
     </property>
     cproperty>
           <name>dfs.data.dir</name>
           <value>file://home/hadoop/hadoopdata/hdfs/datanode</value>
     </property>
</configuration>
```

5. Set Up MapReduce Variables (mapred-site.xml)

```
vi /home/hadoop/hadoop/etc/hadoop/mapred-site.xml
<configuration>
     cproperty>
           <name>mapreduce.framework.name
           <value>varn
     </property>
     cproperty>
           <name>yarn.app.mapreduce.am.env</name>
           <value>HADOOP MAPRED HOME=${HADOOP HOME}
     </property>
     cproperty>
           <name>mapreduce.map.env</name>
           <value>HADOOP MAPRED HOME=${HADOOP HOME}</value>
     </property>
     cproperty>
           <name>mapreduce.reduce.env</name>
           <value>HADOOP MAPRED HOME=${HADOOP HOME}</value>
     </property>
</configuration>
```

6. Set Up YARN Variables (yarn-site.xml)

```
vi /home/hadoop/hadoop/etc/hadoop/yarn-site.xml
<configuration>
         cproperty>
                   <name>yarn.nodemanager.aux-services
                   <value>mapreduce shuffle</value>
         </property>
         cproperty>
                   <name>yarn.nodemanager.resource.memory-mb</name>
                   <value>16384
         </property>
</configuration>
```

7. Clear HDFS

hdfs namenode -format

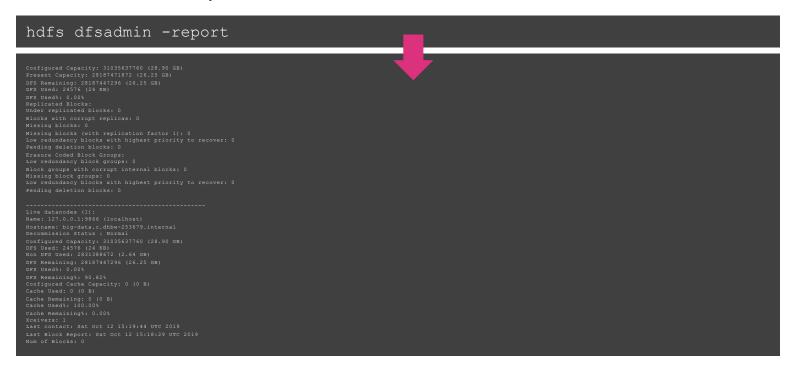
8. Start HDFS:

start-dfs.sh

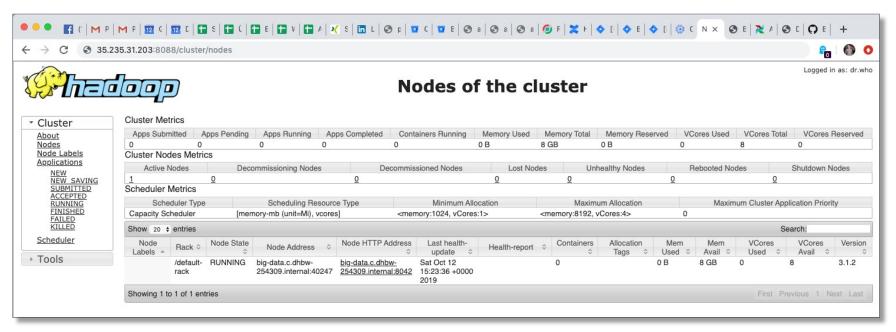
9. Start YARN:

start-yarn.sh

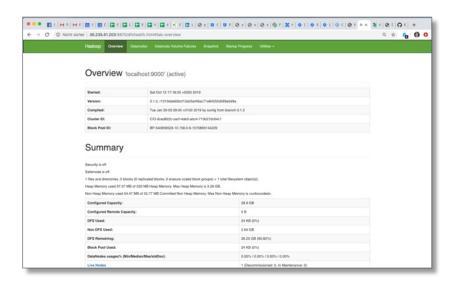
10. Run Admin Status Report

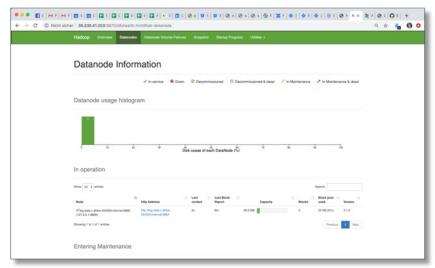


11. Check Ressource Manager Landing Page (http://xxx.xxx.xxx.xxx.xxx.8088/cluster):

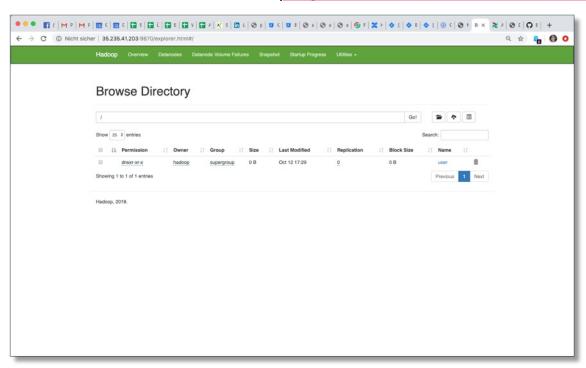


12. Check NameNode Landing and Status Page (http://XXX.XXX.XXX.XXX:9870):





13. Check HDFS File Browser (http://XXX.XXX.XXX.XXX:9870/explorer.html#/)



Работа в HDFS

1. Create User Directory (**on HDFS**):

```
hadoop fs -mkdir /user
hadoop fs -mkdir /user/hadoop
```

2. List Directories (on HDFS):

```
hadoop@big-data:~$ hadoop fs -ls /
Found 1 items
drwxr-xr-x - hadoop supergroup 0 2019-10-12 15:29 /user
hadoop@big-data:~$
```

Работа в **HDFS**

3. Copy File (just a random log file) from local directory to HDFS:

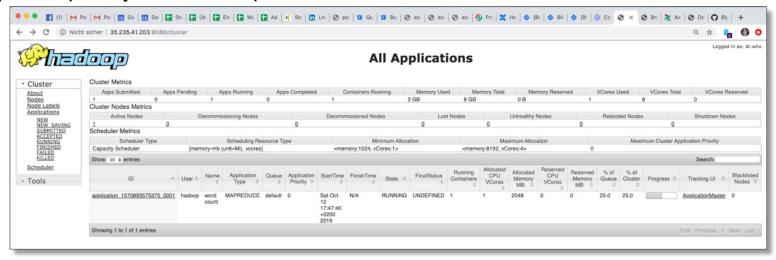
hadoop fs -put /var/log/dpkg.log /user/hadoop/dpkg.log

Запуск Примера MapReduce Job

1. Использование MapReduce WordCount Jar, предоставляемого Hadoop, для подсчета слов в файле

hadoop jar hadoop/share/hadoop/mapreduce/hadoop-mapreduce-examples-3.1.2.jar wordcount /user/hadoop/dpkg.log /user/hadoop/test output

2. Просмотр Запущенного MapReduce Job:



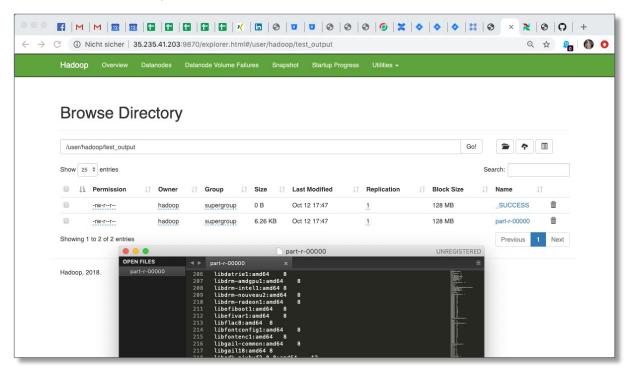
Запуск Примера MapReduce Job

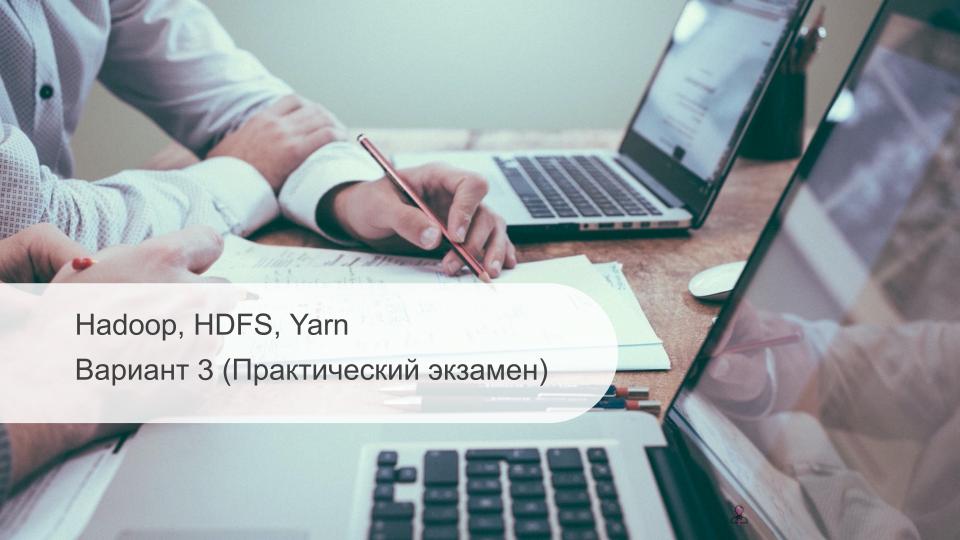
3.Проверить результат Ha Output/Result (via)Ваsh

```
hadoop@big-data:~$ hadoop fs -cat /user/hadoop/test output/part-r-00000
libglx0:amd64 8
libgraphite2-3:amd64 8
libgtk2.0-0:amd64 8
libgtk2.0-bin:amd64 8
libgtk2.0-common:all 9
libharfbuzz0b:amd64 8
libice-dev:amd64 8
libice6:amd64 8
libjbig0:amd64 8
libjpeg-turbo8:amd64 8
libjpeg8:amd64 8
libnss3:amd64 8
libogg0:amd64 8
libpango-1.0-0:amd64 8
libpangocairo-1.0-0:amd64 8
```

Запуск Примера MapReduce Job

4. Проверить результат Ha Output/Result (via Web HDFS File Browser):





Задание

1. Клонировать репозиторий git (чтобы получить пример данных):

git clone https://github.com/BosenkoTM/BigDataWork.git

2

- Скопировать образец файла (/BigDataWork/exercises/winter_semester_2021-2022/01_hadoop/sample_data/Faust_1.txt) из репозитория Git B **HDFS**.
- Запустить MapReduce Jar по умолчанию (hadoop/share/hadoop/mapreduce/hadoop-mapreduce-examples-3.1 .2.jar) для вычисления количества слов в текстовом файле.
- Скопируйте результат MapReduce job обратно в локальную файловую систему ubuntu.
- Используйте и запустите по умолчанию MapReduce Jar (hadoop/share/hadoop/mapreduce/hadoop-mapreduce-examples-3.1 .2.jar), чтобы получить количество вхождений точной строки 'Faust' в текстовом файле.
- Скопируйте результат MapReduce job обратно в локальную файловую систему ubuntu.