# Hadoop分布式环境安装与配置

## 一、搭建虚拟机

本实验将建立三台虚拟机组成的集群，各虚拟机的节点分配如下

| 计算机名称 | 运行节点 |
| --- | --- |
| hadoop01 | NameNode、DataNode、ResourceManager、NodeManager |
| hadoop02 | SecondaryNameNode、DataNode、NodeManager |
| hadoop03 | DataNode、NodeManager |

创建第一个名为hadoop01的虚拟机（略）

启动虚拟机后，输入命令安装需要的工具

```
sudo apt upgrade
sudo apt install openssh-server
sudo apt install net-tools
sudo apt install vim
```

## 二、安装hadoop环境包

1、准备好所有需要的安装包复制到之前设置的共享文件夹VMshare中，包含jdk、hadoop、spark等。建议使用与本实验版本一致的软件包，防止兼容问题。

用以下命令将共享文件夹所有的安装包复制到主目录

```
sudo cp -R /mnt/hgfs/VMShare/ /home/hadoop/
```
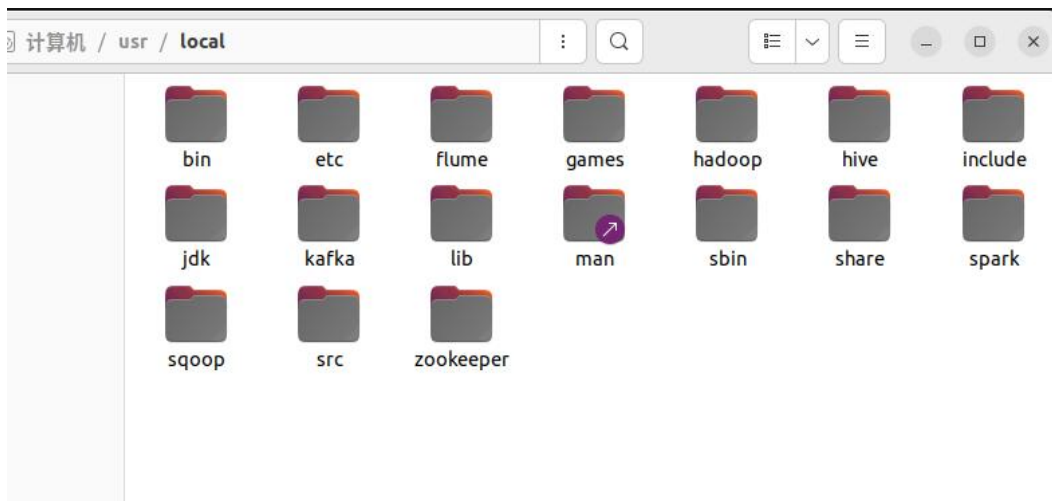
给此目录赋予操作权限

```
sudo chmod -R 777 /home/hadoop/VMShare
```

进入到主目录的VMShare，先安装jdk和hadoop，执行命令解压

```
tar -zxvf jdk-8u261-linux-x64.tar.gz
tar -zxvf hadoop-2.7.7.tar.gz
```

将解压后的文件夹移动到 /usr/local 目录下，同时改成简单的名字

```
sudo mv jdk1.8.0_261/ /usr/local/jdk
sudo mv hadoop-2.7.7/ /usr/local/hadoop
```

### 2、配置环境变量

使用以下命令打开配置文件

```
vi ~/.bashrc
```

在末尾插入以下代码（实际存放位置如不一致需要修改）

```
export JAVA_HOME=/usr/local/jdk
export PATH=$PATH:$JAVA_HOME/bin

export HADOOP_HOME=/usr/local/hadoop
export PATH=$PATH:$HADOOP_HOME/bin:$HADOOP_HOME/sbin
```

保存退出，执行以下代码使配置生效

```
source ~/.bashrc
```

### 3、配置hadoop

进入hadoop配置目录hadoop/etc/hadoop，打开hadoop-env.sh，设置JAVA_HOME的完整存放路径



打开core-site.xml文件,在configuration标签中加入以下代码设置主节点名字，设置一个目录用于存放临时文件

```xml
    <property>
        <name>fs.defaultFS</name>
        <value>hdfs://hadoop01:9000</value>
    </property>
    <property>
        <name>hadoop.tmp.dir</name>
        <value>/usr/local/hadoop/data</value>
    </property>
```

打开hdfs-site.xml文件，往configuration加入设置，设置数据块的副本数量为3，设置hdfs的名字空间元数据、数据块存储位置，根据计划将SecondaryNameNode放在02节点

```xml
    <property>
        <name>dfs.replication</name>
        <value>3</value>
    </property>
    <property>
        <name>dfs.namenode.name.dir</name>
        <value>/usr/local/hadoop/data/dfs/name</value>
    </property>
    <property>
        <name>dfs.datanode.data.dir</name>
        <value>/usr/local/hadoop/data/dfs/data</value>
    </property>
    <property>
        <name>dfs.secondary.http.address</name>
        <value>hadoop02:50090</value>
    </property>
```

将mapred-site.xml.template复制一份并改名为mapred-site.xml，添加设置，mapreduce任务会提交到yarn上

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

打开yarn-site.xml文件，往configuration加入设置，使yarn在01节点运行，并启用mapreduce_shuffle"的辅助服务

```xml
    <property>
        <name>yarn.resourcemanager.hostname</name>
        <value>hadoop01</value>
    </property>
    <property>
        <name>yarn.nodemanager.aux-services</name>
        <value>mapreduce_shuffle</value>
    </property>
```
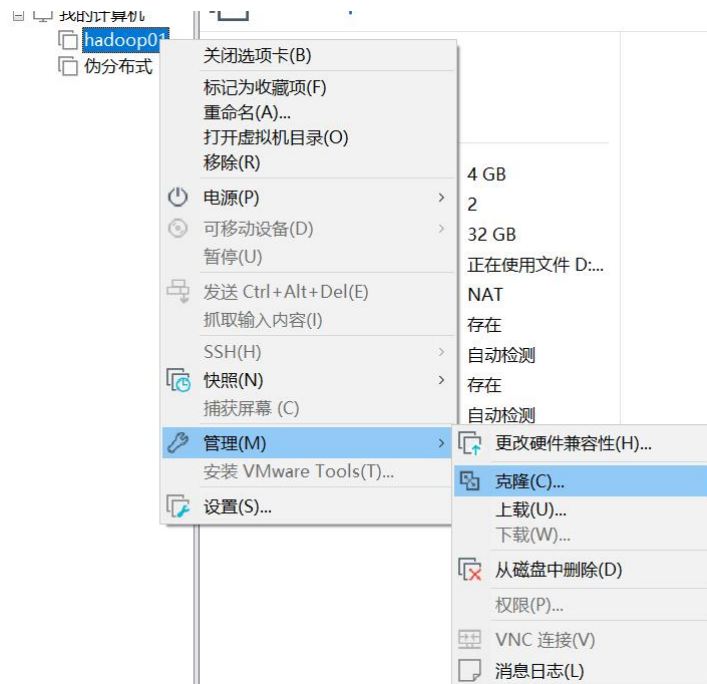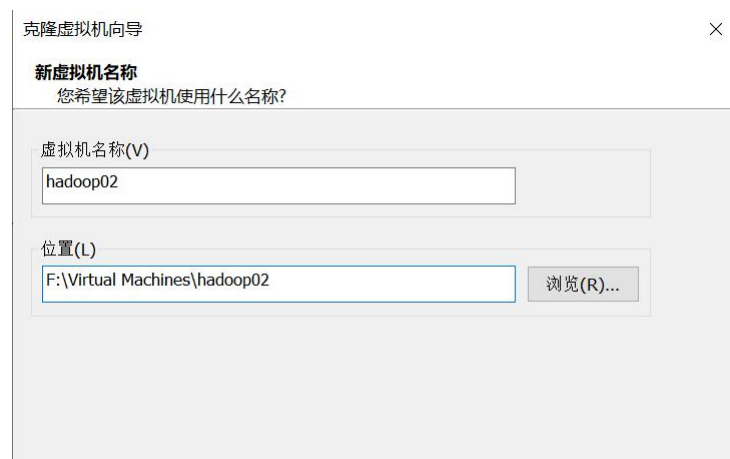
打开slaves文件，编写三台主机名

```
hadoop01
hadoop02
hadoop03
```

## 三、构建集群
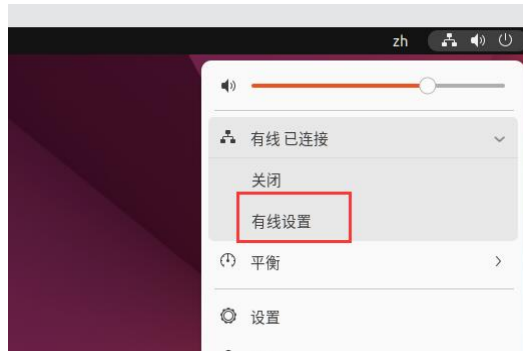
1、克隆虚拟机

将节点01关机，右键——管理——克隆



创建完整克隆



设置节点名称和文件位置



用同样的方法，再将01节点克隆一份为03

2、配置网络

将三个虚拟机全部打开，首先设置01节点的网络，打开有线设置



点开设置按钮，将ip地址复制



切到IPv4页签，选择手动，地址、网关（默认路由）、DNS均与原本一致



**点击应用后，关闭再打开网络生效**

其余两个节点重复同样的操作





3、配置主机IP

回到01节点，打开hosts文件

```
sudo vim /etc/hosts
```

添加以上3台虚拟机的实际ip地址和对应主机名（原本的内容可不删除）

同样的操作在02、03节点也重复一次

4、设置免密登录

回到01节点，打开终端输入命令，中途的对话全部留空直接按回车即可。

```
ssh-keygen -t rsa
```



同样的操作在02、03节点也重复一次。

到01节点拷贝公钥，三个命令都需要输入yes和输密码。

```
ssh-copy-id -i ~/.ssh/id_rsa.pub  hadoop01
ssh-copy-id -i ~/.ssh/id_rsa.pub  hadoop02
ssh-copy-id -i ~/.ssh/id_rsa.pub  hadoop03
```

```
hadoop@master:~$ ssh-copy-id -i ~/.ssh/id_rsa.pub  hadoop01
/usr/bin/ssh-copy-id: INFO: Source of key(s) to be installed: "/home/hadoop/.ssh
/id_rsa.pub"
The authenticity of host 'hadoop01 (192.168.71.134)' can't be established.
ED25519 key fingerprint is SHA256:4BzWbqWwTpAGdhNfnQOWISK/7xDspsZqk604Q5wt85s.
This key is not known by any other names
Are you sure you want to continue connecting (yes/no/[fingerprint])? yes
/usr/bin/ssh-copy-id: INFO: attempting to log in with the new key(s), to filter
out any that are already installed
/usr/bin/ssh-copy-id: INFO: 1 key(s) remain to be installed -- if you are prompt
ed now it is to install the new keys
hadoop@hadoop01's password:

Number of key(s) added: 1

Now try logging into the machine, with:   "ssh 'hadoop01'"
and check to make sure that only the key(s) you wanted were added.

hadoop@master:~$ a
```

同样的操作在02、03节点也重复一次。(每个节点都要拷贝3个公钥)。

任意两个节点间使用ssh命令远程连接,均能直接出现以下结果,而不需要输密码,则设置成功

> ssh 对方主机名

```
hadoop@sz:~$ ssh hadoop02
Welcome to Ubuntu 22.04.5 LTS (GNU/Linux 6.8.0-87-generic x86_64)

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

扩展安全维护(ESM) Applications 未启用。

0 更新可以立即应用。

启用 ESM Apps 来获取未来的额外安全更新
请参见 https://ubuntu.com/esm 或者运行: sudo pro status

New release '24.04.3 LTS' available.
Run 'do-release-upgrade' to upgrade to it.
```
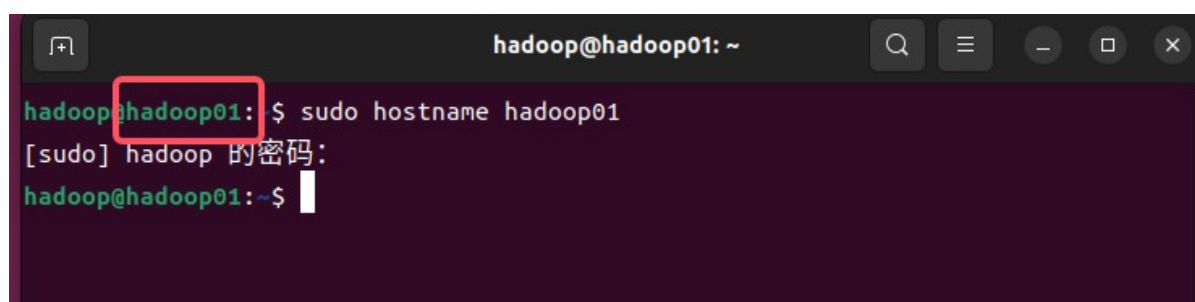
设置主机名,使用以下命令,将三个节点的主机名分别设为对应的hadoop01、hadoop02、hadoop03

> sudo hostname 主机名

设置后,关闭终端重新打开,光标处应显示新的主机名



```
hadoop@hadoop01:~$ sudo hostname hadoop01
[sudo] hadoop 的密码:
hadoop@hadoop01:~$
```

## 四、启动测试

格式化、启动和停止都只需在01节点执行

1、格式化hdfs（首次使用）

```
hdfs namenode –format
```

2、启动hadoop，在输出日志中可以看到哪个节点启动了哪个进程

```
start-all.sh
```

```
hadoop@hadoop01:~$ start-all.sh
This script is Deprecated. Instead use start-dfs.sh and start-yarn.sh
Starting namenodes on [hadoop01]
hadoop01: starting namenode, logging to /usr/local/hadoop/logs/hadoop-hadoop-namenode-hadoop01.out
hadoop02: starting datanode, logging to /usr/local/hadoop/logs/hadoop-hadoop-datanode-hadoop02.out
hadoop03: starting datanode, logging to /usr/local/hadoop/logs/hadoop-hadoop-datanode-hadoop03.out
hadoop01: starting datanode, logging to /usr/local/hadoop/logs/hadoop-hadoop-datanode-hadoop01.out
Starting secondary namenodes [hadoop02]
hadoop02: starting secondarynamenode, logging to /usr/local/hadoop/logs/hadoop-hadoop-secondarynamenode-h
adoop02.out
starting yarn daemons
starting resourcemanager, logging to /usr/local/hadoop/logs/yarn-hadoop-resourcemanager-hadoop01.out
hadoop03: starting nodemanager, logging to /usr/local/hadoop/logs/yarn-hadoop-nodemanager-hadoop03.out
hadoop02: starting nodemanager, logging to /usr/local/hadoop/logs/yarn-hadoop-nodemanager-hadoop02.out
hadoop01: starting nodemanager, logging to /usr/local/hadoop/logs/yarn-hadoop-nodemanager-hadoop01.out
hadoop@hadoop01:~$
```

查分别查看3个节点的进程

```
jps
```

01

```
hadoop@hadoop01:~$ jps
9985 NodeManager
10115 Jps
9690 ResourceManager
9322 NameNode
9471 DataNode
hadoop@hadoop01:~$
```

02

```
hadoop@hadoop02:~$ jps
11410 Jps
11074 DataNode
11290 NodeManager
11197 SecondaryNameNode
hadoop@hadoop02:~$
```
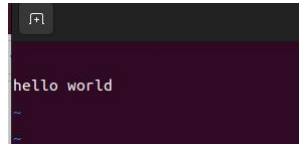
03

```
hadoop@hadoop03:~$ jps
6144 NodeManager
6018 DataNode
6264 Jps
hadoop@hadoop03:~$
```

与开头的规划表对应，部署成功。


4、测试上传文件

在本地创建一个文本文件，内容自定
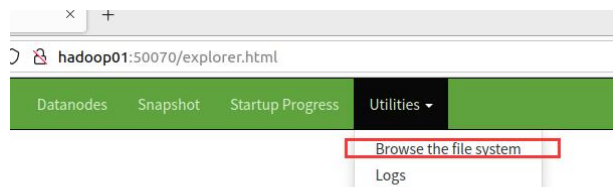
```
vi test.txt
```



在hdfs上创建一个目录

```
hdfs dfs -mkdir /home
```

上传文件到hdfs

```
hdfs dfs -put test.txt /home
```

可通过浏览器访问hdfs的管理页面，地址：http://hadoop01:50070/

点击Utilities——Browse file system可以查看文件系统中的文件



查看刚才上传的文件信息



要结束集群，在01节点执行

```
stop-all.sh
```