

Big Data Cluster

JieLiu (Deployment plan)

1 集群规划

1.1 总体方案

整个计算集群都是基于分布式文件系统HDFS之上，YARN用来实现集群资源的管理与调度，MapReduce用于并行计算。Hive提供了SQL标准的数据存储，管理服务，Pig是一个分析大型数据集的平台。Spark是一个内存分布式计算框架，HBase是分布式数据库，用来管理和存储大规模的数据。

1.2 主机规划

Table 1: 主机规划				
主机名	IP	用户	角色	
cu01	192.168.1.21	hadoop	master	
cu02	192.168.1.22	hadoop	slave01	
cu03	192.168.1.28	hadoop	slave02	
cu04	192.168.1.29	hadoop	slave03	
cu05	192.168.1.30	hadoop	slave04	

2 Hadoop安装配置

2.1 主机名配置

修改集群中每个主机的hosts文件

```
1 #master
2 192.168.1.21 cu01
3 #slaves
4 192.168.1.22 cu02
5 192.168.1.28 cu03
6 192.168.1.29 cu04
7 192.168.1.30 cu05
```

2.2 创建hadoop用户

```
1 root$useradd hadoop
2 root$passwd hadoop
```

2.3 master无密码登录slaves主机配置

生成ssh-key密钥在控制主机上

```
1 ssh-keygen -b 4096
```

拷贝密钥到其他节点

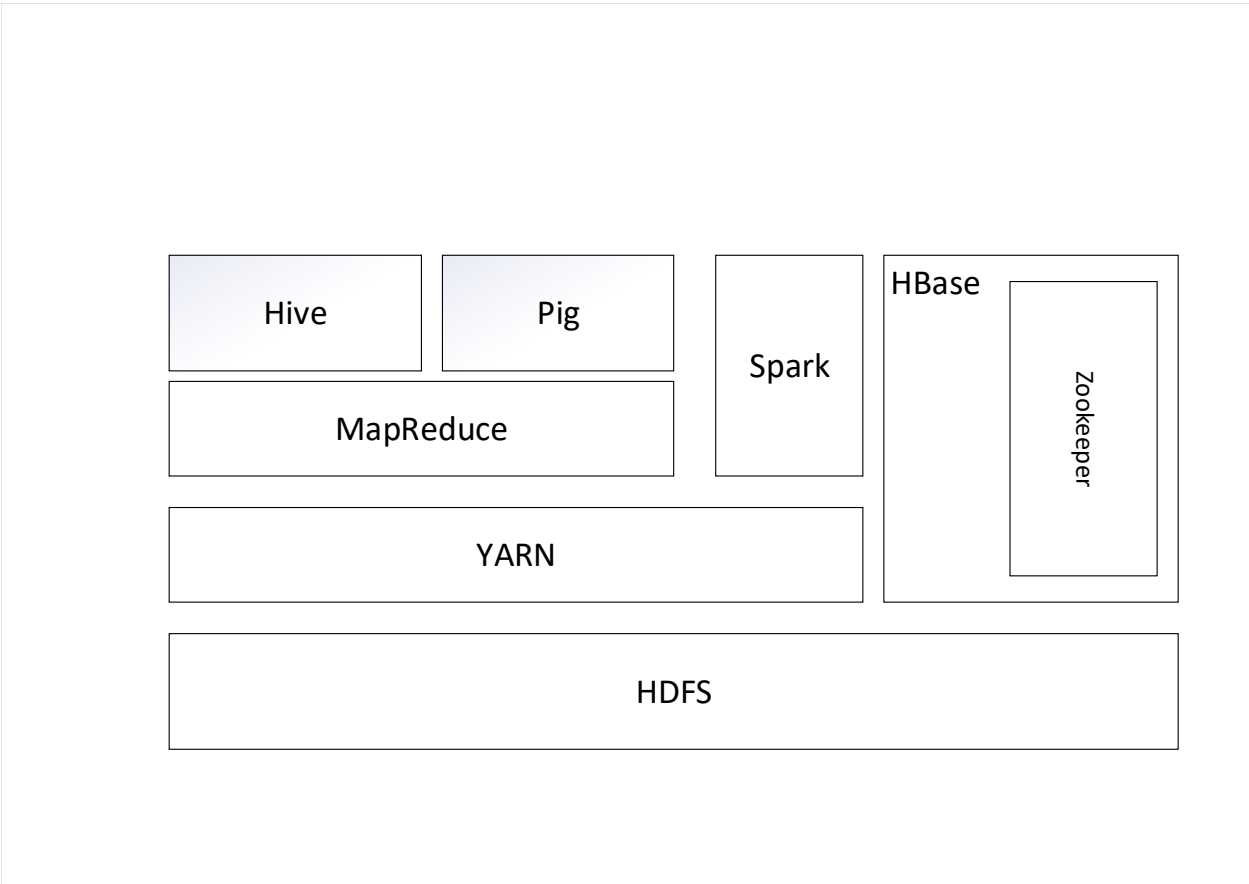


Figure 1: 软件架构

```
1 ssh-copy-id -i $HOME/.ssh/id_rsa.pub hadoop@cu01
2 ssh-copy-id -i $HOME/.ssh/id_rsa.pub hadoop@cu02
3 ssh-copy-id -i $HOME/.ssh/id_rsa.pub hadoop@cu03
4 ssh-copy-id -i $HOME/.ssh/id_rsa.pub hadoop@cu04
5 ssh-copy-id -i $HOME/.ssh/id_rsa.pub hadoop@cu05
```

2.4 下载解压安装包，设置环境变量

```
1 cd ~
2 wget http://good.ncu.edu.cn/mirrors/hadoop-2.8.1.tar.gz
3 tar -xzf hadoop-2.8.1.tar.gz
4 mv hadoop-2.8.1 hadoop
```

配置环境变量(/etc/profile文件)

```
1 export HADOOP_CONF_DIR=/home/hadoop/hadoop/etc/hadoop
2 export LD_LIBRARY_PATH=/home/hadoop/hadoop/lib/native:$LD_LIBRARY_PATH
3 export PATH=/home/hadoop/hadoop/bin:/home/hadoop/hadoop/sbin:$PATH
```

2.5 安装配置java环境

```
1 cd ~
2 wget http://good.ncu.edu.cn/mirrors/jdk-8u161-linux-x64.tar.gz
3 tar -zxvf jdk-8u161-linux-x64.tar.gz
```

配置环境变量

```
1 export JAVA_HOME=/home/hadoop/jdk1.8.0_161
2 export PATH=$PATH:$JAVA_HOME/bin
3 export CLASSPATH=$JAVA_HOME/jre/lib/ext:$JAVA_HOME/lib/tools.jar
```

2.6 master节点配置

修改 /hadoop/etc/hadoop/hadoop-env.sh

```
1 export JAVA_HOME=/home/hadoop/jdk1.8.0_161
```

2.7 配置NameNode位置

修改 /hadoop/etc/hadoop/core-site.xml

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
3
4 <!-- Put site-specific property overrides in this file. -->
5
6 <configuration>
7     <property>
8         <name>fs.default.name</name>
9         <value>hdfs://192.168.1.21:9000</value>
10    </property>
11 </configuration>
```

2.8 配置HDFS路径

修改 /hadoop/etc/hadoop/hdfs-site.xml

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
3
4 <!-- Put site-specific property overrides in this file. -->
5
6 <configuration>
7
8     <property>
9         <name>dfs.namenode.name.dir</name>
10        <value>/home/hadoop/data/nameNode</value>
11    </property>
12
13    <property>
14        <name>dfs.datanode.data.dir</name>
15        <value>/home/hadoop/data/dataNode</value>
16    </property>
17
18    <property>
19        <name>dfs.replication</name>
20        <value>3</value>
21    </property>
22
23 </configuration>
```

2.9 配置YARN为作业调度器

```
1 cd ~/hadoop/etc/hadoop
2 mv mapred-site.xml.template mapred-site.xml
```

修改文件内容

```
1 <?xml version="1.0"?>
2 <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
3
4 <!-- Put site-specific property overrides in this file. -->
5
6 <configuration>
7     <property>
8         <name>mapreduce.framework.name</name>
9         <value>yarn</value>
10    </property>
11
12    <property>
13        <name>yarn.app.mapreduce.am.resource.mb</name>
14        <value>512</value>
15    </property>
16
17    <property>
18        <name>mapreduce.map.memory.mb</name>
19        <value>256</value>
20    </property>
21
22    <property>
23        <name>mapreduce.reduce.memory.mb</name>
```

```
24         <value>256</value>
25     </property>
26
27 </configuration>
```

2.10 配置YARN

修改 /hadoop/etc/hadoop/yarn-site.xml

```
1 <?xml version="1.0"?>
2
3 <configuration>
4
5 <!-- Site specific YARN configuration properties -->
6     <property>
7         <name>yarn.acl.enable</name>
8         <value>0</value>
9     </property>
10
11     <property>
12         <name>yarn.resourcemanager.hostname</name>
13         <value>cu01</value>
14     </property>
15
16     <property>
17         <name>yarn.nodemanager.aux-services</name>
18         <value>mapreduce_shuffle</value>
19     </property>
20
21     <property>
22         <name>yarn.nodemanager.resource.memory-mb</name>
23         <value>1536</value>
24     </property>
25
26     <property>
27         <name>yarn.scheduler.maximum-allocation-mb</name>
28         <value>1536</value>
29     </property>
30
31     <property>
32         <name>yarn.scheduler.minimum-allocation-mb</name>
33         <value>128</value>
34     </property>
35
36     <property>
37         <name>yarn.nodemanager.vmem-check-enabled</name>
38         <value>false</value>
39     </property>
40
41 </configuration>
```

2.11 配置slaves

修改 /hadoop/etc/hadoop/slaves

```
1 cu02
2 cu03
```

```
3 cu04
4 cu05
```

2.12 拷贝文件到每一个slave节点

```
1 cd /home/hadoop/
2 scp -r hadoop cu02:/home/hadoop
3 scp -r hadoop cu03:/home/hadoop
4 scp -r hadoop cu04:/home/hadoop
5 scp -r hadoop cu05:/home/hadoop
```

2.13 配置slave节点的环境变量

同master环境变量配置

2.14 运行测试

```
1 hdfs namenode -format
2 hadoop/sbin/start-all.sh
```

jps查看状态

3 HBase安装配置

3.1 下载解压安装包，设置环境变量

```
1 cd ~
2 wget http://good.ncu.edu.cn/mirrors/hbase-1.2.6-bin.tar.gz
3 tar -xzf hbase-1.2.6-bin.tar.gz
4 mv hbase-1.2.6 hbase
```

配置环境变量(/etc/profile文件)

```
1 export HBASE_HOME=/home/hadoop/hbase
2 export PATH=$PATH:$HBASE_HOME/bin
3 export CLASSPATH=$CLASSPATH:/home/hadoop/hbase/lib/*:.
```

3.2 配置hbase环境

修改hbase/conf下hbase-env.sh

```
1 export JAVA_HOME=/home/hadoop/jdk1.8.0_161
2 export HBASE_MANAGES_ZK=true
```

3.3 配置hbase-site.xml

```
1 <?xml version="1.0"?>
2 <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
3 <configuration>
4
5 <property>
```

```
6 <name>hbase.rootdir</name>
7 <value>hdfs://192.168.1.21:9000/hbase</value>
8 </property>
9 <property>
10 <name>hbase.cluster.distributed</name>
11 <value>true</value>
12 </property>
13 <property>
14 <name>hbase.zookeeper.property.dataDir</name>
15 <value>/home/hadoop/zookeeper</value>
16 </property>
17 <property>
18 <name>hbase.zookeeper.quorum</name>
19 <value>cu01,cu02,cu03,cu04,cu05</value>
20 </property>
21 <property>
22 <name>hbase.zookeeper.property.clientPort</name>
23 <value>2181</value>
24 </property>
25
26 </configuration>
```

3.4 配置从服务器

修改文件conf/regionservers

```
1 cu02
2 cu03
3 cu04
4 cu05
```

3.5 创建zookeeper目录

在每个节点与hbase同级目录创建zookeeper文件夹

3.6 region节点配置

拷贝hbase文件到region节点

```
1 scp -r hbase cu02:/home/hadoop
2 scp -r hbase cu03:/home/hadoop
3 scp -r hbase cu04:/home/hadoop
4 scp -r hbase cu05:/home/hadoop
```

3.7 region节点环境变量配置

同master节点配置

3.8 启动测试hbase

```
1 $HBASE_HOME/bin/start-hbase.sh
2 hbase shell
```

4 Hive安装配置

4.1 下载解压安装包，设置环境变量

```
1 cd ~
2 wget http://good.ncu.edu.cn/mirrors/apache-hive-1.2.2-bin.tar.gz
3 tar -xzf apache-hive-1.2.2-bin.tar.gz
4 mv apache-hive-1.2.2-bin.tar.gz hive
```

配置环境变量(/etc/profile文件)

```
1 export HIVE_HOME=/home/hadoop/hive
2 export PATH=$HIVE_HOME/bin:$PATH
```

4.2 安装mysql用于存储元数据

```
1
2 wget http://repo.mysql.com/mysql-community-release-el7-5.noarch.rpm
3 sudo rpm -ivh mysql-community-release-el7-5.noarch.rpm
4 yum update
5
6 sudo yum install mysql-server
7 sudo systemctl start mysqld
8
9 sudo mysql_secure_installation
```

登录mysql创建hive用户,并创建hive数据库

```
1 mysql>CREATE USER 'hive' IDENTIFIED BY 'hive';
2 mysql>GRANT ALL PRIVILEGES ON *.* TO 'hive'@'cu01' WITH GRANT OPTION;
3 mysql>flush privileges;
4
5 mysql>create database hive;
```

4.3 配置hive-site.xml

找到hive-default.xml.template,cp一份为hive-default.xml

```
1
2 <?xml version="1.0"?>
3 <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
4
5 <configuration>
6     <property>
7         <name>javax.jdo.option.ConnectionURL</name>
8         <value>jdbc:mysql://192.168.1.21:3306/hive?createDatabaseIfNotExist=true</value>
9         <description>JDBC connect string for a JDBC metastore</description>
10    </property>
11    <property>
12        <name>javax.jdo.option.ConnectionDriverName</name>
13        <value>com.mysql.jdbc.Driver</value>
14        <description>Driver class name for a JDBC metastore</description>
15    </property>
16
17    <property>
18        <name>javax.jdo.option.ConnectionUserName</name>
19        <value>root</value>
```



```
20     <description>username to use against metastore database</description>
21     </property>
22
23     <property>
24     <name>javax.jdo.option.ConnectionPassword</name>
25     <value>work@good308</value>
26     <description>password to use against metastore database</description>
27     </property>
28
29 </configuration>
```

4.4 JDBC依赖配置

```
1 wget http://cdn.mysql.com/Downloads/Connector-J/mysql-connector-java-5.1.36.tar.gz
2 #
3 tar -zxvf mysql-connector-java-5.1.36.tar.gz
4 cp mysql-connector-java-5.1.33-bin.jar apache-hive-1.2.1-bin/lib/
```

4.5 hive客户端配置

```
1 <configuration>
2     <property>
3         <name>hive.metastore.uris</name>
4         <value>thrift://192.168.1.21:9083</value>
5     </property>
6 </configuration>
```

4.6 启动并测试hive

```
1 #start metastore service
2 hive --service metastore &
3 #start test hive
4 hive
5
6 hive> show databases;
7 OK
8 default
9 src
10 Time taken: 1.332 seconds, Fetched: 2 row(s)
11 hive> use src;
12 OK
13 Time taken: 0.037 seconds
14 hive> create table test1(id int);
15 OK
16 Time taken: 0.572 seconds
17 hive> show tables;
18 OK
19 abc
20 test
21 test1
22 Time taken: 0.057 seconds, Fetched: 3 row(s)
23 hive>
```

5 安装配置Spark

5.1 下载解压安装包，设置环境变量

```
1 cd ~
2 wget https://good.ncu.edu.cn/mirrors/spark-2.2.0-bin-hadoop2.7.tgz
3 tar -xvf spark-2.2.0-bin-hadoop2.7.tgz
4 mv spark-2.2.0-bin-hadoop2.7 spark
```

```
1 pathmunge /home/hadoop/spark/bin
```

配置环境变量(/etc/profile文件)

```
1 export HADOOP_CONF_DIR=/home/hadoop/hadoop/etc/hadoop
2 export SPARK_HOME=/home/hadoop/spark
3 export LD_LIBRARY_PATH=/home/hadoop/hadoop/lib/native:$LD_LIBRARY_PATH
```

5.2 配置YARN管理spark

```
1 #
2 # Licensed to the Apache Software Foundation (ASF) under one or more
3 # contributor license agreements. See the NOTICE file distributed with
4 # this work for additional information regarding copyright ownership.
5 # The ASF licenses this file to You under the Apache License, Version 2.0
6 # (the "License"); you may not use this file except in compliance with
7 # the License. You may obtain a copy of the License at
8 #
9 # http://www.apache.org/licenses/LICENSE-2.0
10 #
11 # Unless required by applicable law or agreed to in writing, software
12 # distributed under the License is distributed on an "AS IS" BASIS,
13 # WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
14 # See the License for the specific language governing permissions and
15 # limitations under the License.
16 #
17
18 # Default system properties included when running spark-submit.
19 # This is useful for setting default environmental settings.
20
21 # Example:
22 spark.master yarn
23 spark.yarn.am.memory 512m
24 spark.executor.memory 512m
25
26 spark.eventLog.enabled true
27 spark.eventLog.dir hdfs://cu01:9000/spark-logs
28 # spark.serializer org.apache.spark.serializer.KryoSerializer
29 spark.driver.memory 512m
30 spark.history.provider org.apache.spark.deploy.history.FsHistoryProvider
31 spark.history.fs.logDirectory hdfs://cu01:9000/spark-logs
32 spark.history.fs.update.interval 10s
33 spark.history.ui.port 18080
34
35 # spark.executor.extraJavaOptions -XX:+PrintGCDetails -Dkey=value -Dnumbers="one two three"
```

6 安装配置pig

6.1 下载解压安装包，设置环境变量

```
1 cd ~
2 wget https://good.ncu.edu.cn/mirrors/pig-0.16.0.tar.gz
3 tar -xvf pig-0.16.0.tar.gz
4 mv pig-0.16.0.tar.gz pig
```

配置环境变量(/etc/profile文件)

```
1 export PIG_HOME=/home/hadoop/pig
2 export PATH=$PATH:/home/hadoop/pig/bin
3 export PIG_CLASSPATH=$HADOOP_HOME/conf
```
