# SequoiaDB Spark Yarn 部署及案例演示

# 1、 背景

由于 MRv1 在扩展性、可靠性、资源利用率和多框架等方面存在明显的不足，在 Hadoop MRv2 中引入了资源管理和调度系统 YARN。YARN 是 Hadoop MRv2 计算机框架中构建的一个独立的、通用的资源管理系统，可为上层应用提供统一的资源管理和调度，它的引入为集群在利用率、资源统一管理和数据共享等方面带来了巨大好处。主要体现在以下几个方面：

（1）资源利用率大大提高。一种计算框架一个集群，往往会由于应用程序数量和资源需求的不均衡性，使得在某段时间有些计算框架集群资源紧张，而另外一些集群资源空闲。共享集群模式则通过多种框架共享资源，使得集群中的资源得到更加充分的利用；

（2）运维成本大大降低。共享集群模式使得少数管理员就可以完成多个框架的统一管理；

（3）共享集群的模式也让多种框架共享数据和硬件资源更为方便。

# 2、 产品介绍

巨杉数据库 SequoiaDB 是一款分布式非关系型文档数据库，可以被用来存取海量非关系型的数据，其底层主要基于分布式，高可用，高性能与动态数据类型设计，它兼顾了关系型数据库中众多的优秀设计：如索引、动态查询和更新等，同时以文档记录为基础更好地处理了动态灵活的数据类型。并且为了用户能够使用常见的分布式计算框架，SequoiaDB 可以和常见分布式计算框架如 Spark、Hadoop、HBase 进行整合。本文主要讲解 SequoiaDB 与 Spark、YARN 的整合以及通过一个案例来演示 MapReduce 分析存储在 SequoiaDB 中的业务数据。

# 3、 环境搭建

## 3.1、 服务器分布

| 服务器 | 服务名称 |
|---|---|
| 192.168.1.46 | |
| 192.168.1.47 | |
| 192.168.1.48 | NameNode、DataNode、 |

## 3.2、 软件配置

操作系统：RedHat6.5
JDK 版本：1.7.0_80 64 位
Scala 版本：
Hadoop 版本：2.7.2
Spark 版本：2.0
SequoiaDB 版本：2.0


## 3.3、 安装步骤

1、JDK 安装
```
tar  -xvf  jdk-7u45-linux-x64.tar.gz   -C /usr/local
cd  /usr/local
ln  -s  jdk1.7.0_45  jdk
```

配置环境变量
```
vim  ~/.bash_profile
export  JAVA_HOME=/usr/local/jdk
export  CLASS_PATH=$JAVA_HOME/lib:$JAVA_HOME/jre/lib
export  PATH=$PATH:$JAVA_HOME/bin
source  /etc/profile
```

2、Scala 安装
```
tar  -xvf  scala-2.11.8.tgz   -C /usr/local
cd  /usr/local
ln  -s  scala-2.11.8 scala
```

配置环境变量
```
vim  ~/.bash_profile
export SCALA_HOME=/usr/local/scala
export PATH=$PATH:$SCALA_HOME/bin
```

3、修改主机 hosts 文件配置
在每台主机上修改 host 文件
```
vim /etc/hosts
192.168.1.46 node01
192.168.1.47 node02
192.168.1.48  master
```

4、    SSH 免密钥登录
在 master 节点中执行 ssh-keygen 按回车键

```
cat  ~/.ssh/id_rsa.pub  >>  ~/.ssh/authorized_keys
```

将 master 节点中的授权文件 authorized_keys 传输到 slave 节点中
```
scp ~/.ssh/id_rsa.pub root@master:~/.ssh/
```

在 slave 节点中执行
```
cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys
```

在 slave 节点中验证 SSH 免密钥登录
```
ssh master
```

5、Hadoop 集群安装
拷贝 hadoop 文件 hadoop-2.7.2.tar.gz 到/opt 目录中
解压 hadoop 安装包
```
tar  -xvf hadoop-2.7.2.tar.gz
mv hadoop-2.7.2 /opt/cloud/hadoop
```

创建 hadoop 数据存储及临时目录
```
mkdir  -p /opt/hadoop/data
mkdir  -p /opt/hadoop/tmp
```

配置 Hadoop  jdk 环境变量
```
vim hadoop-env.sh
export JAVA_HOME=/usr/local/jdk
```

编辑 core.xml 文件
```
<configuration>
    <property>
        <name>fs.defaultFS</name>
        <value>hdfs://master:9000</value>
    </property>
    <property>
        <name>hadoop.tmp.dir</name>
        <value>/opt/data/tmp</value>
    </property>
    <property>
        <name>io.file.buffer.size</name>
        <value>4096</value>
    </property>
</configuration>
```

编辑 mapred-site.xml
```
<configuration>
```

```xml
    <property>
        <name>mapreduce.framework.name</name>
        <value>yarn</value>
    </property>
    <property>
        <name>mapreduce.jobtracker.http.address</name>
        <value> master:50030</value>
    </property>
    <property>
        <name>mapreduce.jobhistory.address</name>
        <value> master:10020</value>
    </property>
    <property>
        <name>mapreduce.jobhistory.webapp.address</name>
        <value>master:19888</value>
    </property>
</configuration>
```

编辑 hdfs-site.xml
```xml
<configuration>
    <property>
        <name>dfs.nameservices</name>
        <value>master</value>
    </property>
    <property>
        <name>dfs.namenode.secondary.http-address</name>
        <value> master:50090</value>
    </property>
    <property>
        <name>dfs.namenode.name.dir</name>
        <value>file:///opt/hadoop/data/name</value>
    </property>
    <property>
        <name>dfs.datanode.data.dir</name>
        <value>file:///opt/hadoop/data</value>
    </property>
    <property>
        <name>dfs.replication</name>
        <value>3</value>
    </property>
    <property>
        <name>dfs.webhdfs.enabled</name>
        <value>true</value>
```

```
        </property>
</configuration>

编辑 yarn-site.xml
<configuration>
    <property>
        <name>yarn.nodemanager.aux-services</name>
        <value>mapreduce_shuffle</value>
    </property>
    <property>
        <name>yarn.resourcemanager.address</name>
        <value> master:8032</value>
    </property>
    <property>
        <name>yarn.resourcemanager.scheduler.address</name>
        <value> master:8030</value>
    </property>
    <property>
        <name>yarn.resourcemanager.resource-tracker.address</name>
        <value> master:8031</value>
    </property>
    <property>
        <name>yarn.resourcemanager.admin.address</name>
        <value> master:8033</value>
    </property>
    <property>
        <name>yarn.resourcemanager.webapp.address</name>
        <value>master:8088</value>
    </property>
    <property>
        <name>yarn.nodemanager.resource.memory-mb</name>
        <value>12288</value>
    </property>
  <property>
        <name>yarn.nodemanager.log-dirs</name>
        <value>/opt/hadoop/tmp/userlogs</value>
    </property>
</configuration>
```

启动 Hadoop
首次启动集群时，做如下操作
进入到/opt/cloud/hadoop/bin 目录中执行./hdfs namenode –format 格式化

hdfs 文件系统

进入到/opt/cloud/hadoop/sbin 目录中执行./start-all.sh 启动 hadoop 集群

6、安装 Spark 集群
　　拷贝 Spark 安装包到/opt 目录中，解压
　　tar －xvf spark-2.0.0-bin-hadoop2.7.tgz
　　mv spark-2.0.0-bin-hadoop2.7 /opt/cloud/spark

　　编辑 spark-env.sh
　　vim spark-env.sh
　　JAVA_HOME="/usr/jdk1.7"
　　SPARK_DRIVER_MEMORY="1g"
　　SPARK_EXECUTOR_CORES=1
　　SPARK_EXECUTOR_MEMORY="512m"
　　SPARK_MASTER_PORT="7077"
　　SPARK_MASTER_WEBUI_PORT="8070"
　　SPARK_CLASSPATH="/opt/cloud/spark/jars/sequoiadb.jar:/opt/cloud/spark/jars/spark-sequoiadb_2.11-2.6.0.jar"
　　SPARK_MASTER_IP="node03"
　　SPARK_WORKER_MEMORY="712m"
　　SPARK_WORKER_CORES=1
　　SPARK_WORKER_INSTANCES=1
　　SPARK_WORKER_DIR="/opt/data/spark/work"
　　SPARK_LOCAL_DIRS="/opt/data/spark/tmp"
　　HADOOP_HOME="/opt/cloud/hadoop"
　　HADOOP_CONF_DIR=$HADOOP_HOME/etc/hadoop

　　编辑  slaves
　　node02
　　node03

　　启动 spark 集群
　　进入到目录/opt/cloud/spark/sbin 目录中
　　./start-all.sh
　　Spark 成功启动后截图如下：

The following is a browser window showing Spark Master page.

URL bar: 192.168.1.48:8070

**Spark** 2.0.1  **Spark Master at spark://node03:7077**

**URL:** spark://node03:7077
**REST URL:** spark://node03:6066 *(cluster mode)*
**Alive Workers:** 2
**Cores in use:** 2 Total, 0 Used
**Memory in use:** 1424.0 MB Total, 0.0 B Used
**Applications:** 0 Running, 0 Completed
**Drivers:** 0 Running, 0 Completed
**Status:** ALIVE

**Workers**

| Worker Id | Address | State | Cores | Memory |
|---|---|---|---|---|
| worker-20161116001004-192.168.1.48-59316 | 192.168.1.48:59316 | ALIVE | 1 (0 Used) | 712.0 MB (0.0 B Used) |
| worker-20161118215804-192.168.1.47-40327 | 192.168.1.47:40327 | ALIVE | 1 (0 Used) | 712.0 MB (0.0 B Used) |

**Running Applications**

| Application ID | Name | Cores | Memory per Node | Submitted Time | User | State | Duration |
|---|---|---|---|---|---|---|---|

**Completed Applications**

| Application ID | Name | Cores | Memory per Node | Submitted Time | User | State | Duration |
|---|---|---|---|---|---|---|---|

7、Spark Yarn 连接 SequoiaDB

在 SequoiaDB 中创建集合空间、集合
db.createCS('poc');
db.poc.createCL('test');

```
help() for help, Ctrl+c or quit to exit
> db = new Sdb();
localhost:11810
Takes 0.6467s.
> db.createCS('poc');
localhost:11810.poc
Takes 0.26202s.
> db.poc.createCL('test');
localhost:11810.poc.test
Takes 19.459596s.
> db.poc.test.insert({id:1,name:'测试1'});
Takes 0.35498s.
> db.poc.test.insert({id:2,name:'测试2'});
Takes 0.17728s.
> db.poc.test.find();
{
  "_id": {
    "$oid": "582f0b47041f736581000000"
  },
  "id": 1,
  "name": "测试1"
}
{
  "_id": {
    "$oid": "582f0b4d041f736581000001"
  },
  "id": 2,
  "name": "测试2"
}
Return 2 row(s).
Takes 0.42897s.
```

进入到 spark 安装目录 bin 中，执行./spark-sql –master yarn 启动 spark sql 交互界面



创建表，映射到上述 poc 集合空间中 test 集合

```
CREATE TABLE `test` (`id` INT, `name` STRING)
USING com.sequoiadb.spark
OPTIONS (
    `collection` 'test',
    `host` 'node02:11810,node03:11810',
    `serialization.format` '1',
    `collectionspace` 'poc'
);
```

查询表 test 数据，执行：

```
Select * from test;
```



进入到 yarn 管理页面查看 spark 任务

# 5、 案例演示

为了配合司法部门的执法和银行内部的风险监管，部分商业银行对于存取款业务定制了相关预警方案，本案例以个人存取款业务高频交易来讲述 MapReduce 如何分析 SequoiaDB 中的个人交易明细数据。

具体场景为：分析同一实体柜员办理，1 小时内同一账户连续 3 笔以上支取类金额的交易账户及明细。

本演示案例采用 Hadoop Map Reduce 实现，开发语言为 Java 语言。整个测试程序分为两个部分 Map 算法和 Reduce 算法。演示程序中 Map 算法负责将同一个账号的所有对应交易明细归并在一起并输出给 Reduce 端，Reduce 端根据 Map 算法的结果运算具体的业务场景，最后将运算结果写入到 SequoiaDB 中。

具体架构如下：



Reduce 端具体算法流程如下：

开始

SequoiaDB Hadoop Connector

Map端对同账号对应的交易明细进行归并

对同账号对应的交易明细进行遍历转换成BSONObject并保存到List集合中

对List集合中的BSONObject按照时间字段进行降序排序

遍历List集合并获取交易明细相关字段信息

判断List集合遍历是否结束 —— 是

否

判断交易代码是否满足要求

是

将交易明细BSONObject对象放入临时Map对象tempMap中

遍历List集合并获取交易明细相关字段信息

判断List集合遍历是否结束 —— 是

否

判断是否是同一交易柜员办理

否

是

判断交易时间是否是一个小时内

是

将交易明细BSONObject对象放入临时Map对象tempMap中

判断tempMap的长度是否大于等于3

是

将交易明细BSONObject加入到结果Map对象result中

清空tempMap元素

遍历List集合并获取交易明细相关字段信息

判断List集合遍历是否结束 —— 是

否

将交易明细BSONObject对象放入临时Map对象tempMap中

遍历List集合并获取交易明细相关字段信息

判断List集合遍历是否结束 —— 是

否

判断是否是同一交易柜员办理

否

是

判断交易时间是否是半个小时内

是

将交易明细BSONObject对象放入临时Map对象tempMap中

判断tempMap的长度是否大于等于2和交易金额是否大于等于10000000

是

将交易明细BSONObject加入到结果Map对象result中

清空tempMap元素

将result和result2中的BSONObject保存到SequoiaDB中

结束

Map 端算法代码如下：

```java
static class TMapper extends Mapper<Object, BSONWritable,Text,BSONWritable>{
    @Override
    protected void map(Object key, BSONWritable value, Context context)
            throws IOException, InterruptedException {
                BSONObject obj = value.getBson();
        String acct_no=(String) obj.get("ACCT_NO");
```

```
                context.write(new Text(acct_no), value);
        }
    }
```

Reduce 端算法代码如下：

```java
static class TReducer extends Reducer<Text,BSONWritable,NullWritable,NullWritable>{
        private static String pattern = "yyyy-MM-dd HH:mm:ss";
        private DateFormat df = new SimpleDateFormat(pattern);
        private static int tradeNum1 = 3;
        private static int tradeTime1 = 3600;

        private static int tradeNum2 = 2;
        private static int tradeTime2 = 1800;
        private static int tradeAll = 100000;

        private Sequoiadb sdb = null;
        private CollectionSpace cs = null;
        private DBCollection cl_1 = null;
        private DBCollection cl_2 = null;
        private static String CS_NAME="";
        private static String CL_NAME_1="";
        private static String CL_NAME_2="";

        public TReducer(){
            if (null == sdb) {
                sdb = ConnectionPool.getInstance().getConnection();
            }

            if (sdb.isCollectionSpaceExist(CS_NAME)) {
                cs = sdb.getCollectionSpace(CS_NAME);
            } else {

                throw new BaseException("集合空间" + CS_NAME + "不存在！");

            }

            if (null == cs) {

                throw new BaseException("集合空间不能为null！");
```

```java
        } else {
            this.cl_1 = cs.getCollection(CL_NAME_1);

        }
        if (null == cs) {

            throw new BaseException("集合空间不能为null！");

        } else {
            this.cl_2 = cs.getCollection(CL_NAME_2);

        }
    }

    @Override
    protected void reduce(Text key, Iterable<BSONWritable> values,
            Context context)
            throws IOException, InterruptedException{
        Iterator<BSONWritable> iterator=values.iterator();
        long sum=0;
        List<BSONWritable> oldList = new ArrayList<BSONWritable>();

        while(iterator.hasNext()){
          BSONWritable bsonWritable = iterator.next();
          oldList.add(bsonWritable);
        }

        //对values进行排序，排序字段为TRN_TIME（交易时间）

        Collections.sort(oldList, new Comparator<BSONWritable>() {
            @Override
            public int compare(BSONWritable o1, BSONWritable o2) {
                String trn_time1 = (String)o1.getBson().get("TRN_TIME");
                String trn_time2 = (String)o2.getBson().get("TRN_TIME");
                return trn_time2.compareTo(trn_time1);
            }
        });

         Map<String,BSONWritable> result = new HashMap<St
```

```java
ring,BSONWritable>();
        if(oldList != null && oldList.size() > 0){
            //记录同一账户满足条件的笔数
            Map<String,BSONWritable> tempMap = new HashMap<String,BSONWritable>();
            for(int i=0;i<oldList.size()-1;i++){
                BSONWritable bSONWritable1 = oldList.get(i);
                //交易代码
                String trn_cd = (String)bSONWritable1.getBson().get("TRN_CD");
                if(trn_cd.equals("000045") || trn_cd.equals("001045")
                    || trn_cd.equals("021031") || trn_cd.equals("020031")
                    || trn_cd.equals("001060") || trn_cd.equals("000060")){
                    //交易柜员
                    String tran_teller_no1 = (String)bSONWritable1.getBson().get("TRAN_TELLER_NO");
                    //流水号
                    String jrnl_no = (String)bSONWritable1.getBson().get("JRNL_NO");
                    //交易日期
                    String trn_date1 = (String)bSONWritable1.getBson().get("TRN_DATE");
                    //交易时间
                    String trn_time1 = (String)bSONWritable1.getBson().get("TRN_TIME");
                    Date bigDate = null;
                    try {
                        bigDate = df.parse(trn_date1+" "+trn_time1);
                    } catch (ParseException e) {
                        e.printStackTrace();
                    }
                    tempMap.put(jrnl_no,bSONWritable1);
                    for(int j=i+1;j<oldList.size();j++){
```

```java
                        BSONWritable bSONWritable2 = oldList.
get(j);

                        //交易代码

                        String trn_cd1 = (String)bSONWritable
2.getBson().get("TRN_CD");
                        if(trn_cd1.equals("000045") || trn_cd
1.equals("001045")
                                || trn_cd1.equals("021031") ||
trn_cd1.equals("020031")
                                || trn_cd1.equals("001060") ||
trn_cd1.equals("000060")){
                            //交易柜员

                            String tran_teller_no2 = (String)b
SONWritable2.getBson().get("TRAN_TELLER_NO");
                            //流水号

                            String jrnl_no2 = (String)bSONWrit
able2.getBson().get("JRNL_NO");
                            //交易日期

                            String trn_date2 = (String)bSONWri
table2.getBson().get("TRN_DATE");
                            //交易时间

                            String trn_time2 = (String)bSONWri
table2.getBson().get("TRN_TIME");
                            Date smallDate = null;
                            try {
                                smallDate = df.parse(trn_date1+
" "+trn_time1);
                            } catch (ParseException e) {
                                e.printStackTrace();
                            }
                            //判断是否是同一实体{交易柜员}办理

                            if(!tran_teller_no1.equals(tran_te
ller_no2)){

                                continue;
                            }
                            //判断{交易日期}{交易时间}是否是[1小时]
内
```

```java
                    if((bigDate.getTime()-smallDate.ge
tTime())/1000 > tradeTime1){
                            break;
                    }
                    tempMap.put(jrnl_no2,bSONWritable
2);

                }else{ //end if TRN_CD1.equals("00004
5")

                    continue;
                }
            }//end for

            if(tempMap.size() >= tradeNum1){
                result.putAll(tempMap);
                tempMap.clear();
            }
        }else{
            continue;
        }//end if ||
    }//end for
}

Map<String,BSONWritable> result2 = new HashMap<S
tring,BSONWritable>();

List<BSONObject> cl_1_list = new ArrayList<BSONO
bject>();

//结果写入sdb

Iterator iter1 = result.keySet().iterator();
while(iter1.hasNext()){
    String keyValue = (String)iter1.next();
    BSONWritable resultValue = result.get(keyValu
e);
    cl_1_list.add(resultValue.getBson());
    cl_1.insert(resultValue.getBson());
}
cl_1.bulkInsert(cl_1_list, DBCollection.FLG_INS
ERT_CONTONDUP);
cl_1_list = null;
List<BSONObject> cl_2_list = new ArrayList<BSONO
bject>();
context.write(null,null);
```

```
        }
    }
```