

docker搭建hadoop完全分布式

计算机科学与技术17-1 陈巍

0. 说明

在已经构建好的单机伪分布式基础上搭建hadoop完全分布式。

【注意】 提前创建好docker网络

```
1 | sudo docker network create --subnet=172.24.0.0/24 will-network
```

1. 集群规划

这里搭建一个 5 个节点的 Hadoop 集群，其中四台主机均部署 `DataNode` 和 `NodeManager` 服务，但其中 `hadoop02` 上部署 `SecondaryNameNode` 服务和 `NTP服务器`，`hadoop01` 上部署 `NameNode`、`ResourceManager` 和 `JobHistoryServer` 服务。

```
master = hadoop01 => NameNode && ResourceManager && JobHistoryServer
```

```
slave = {
```

```
hadoop02; => DataNode && NodeManager && SecondaryNamenode && NTP-Server
```

```
hadoop03; => DataNode && NodeManager
```

```
hadoop04; => DataNode && NodeManager
```

```
hadoop05; => DataNode && NodeManager
```

```
}
```

```
1 # 节点ip分配:
2 hadoop01=172.24.0.2
3 hadoop02=172.24.0.3
4 hadoop03=172.24.0.4
5 hadoop04=172.24.0.5
6 hadoop05=172.24.0.6
```

2. 创建并启动容器

使用一个「单机伪分布式镜像」启动多个镜像。

2.1 手动启动方式

```
1 # 下面会用到systemctl 管理服务（配置时间同步），需要加上参数 --
  privileged 来增加权并且不能使用默认的bash，换成 init（能够使用
  systemctl 命令）
2 # hadoop01
3 # 50070 hdfs
4 # 8088 yarn
5 # 19888 JobHistoryServer
6 sudo docker run -d --name hadoop01 --hostname hadoop01
  --net will-network --ip 172.24.0.2 -P -p 50070:50070 -p
  8088:8088 -p 19888:19888 --privileged 2d7d22776b99
  /usr/sbin/init
7
8 # hadoop02
9 # 50090 SecondaryNameNode
10 sudo docker run -d --name hadoop02 --hostname hadoop02
  --net will-network --ip 172.24.0.3 -P -p 50090:50090 --
  privileged 2d7d22776b99 /usr/sbin/init
11
12 # hadoop03
13 sudo docker run -d --name hadoop03 --hostname hadoop03
  --net will-network --ip 172.24.0.4 -P --privileged
  2d7d22776b99 /usr/sbin/init
14
15 # hadoop04
```

```

16 sudo docker run -d --name hadoop04 --hostname hadoop04
   --net will-network --ip 172.24.0.5 -P --privileged
   2d7d22776b99 /usr/sbin/init
17
18 # hadoop05
19 sudo docker run -d --name hadoop05 --hostname hadoop05
   --net will-network --ip 172.24.0.6 -P --privileged
   2d7d22776b99 /usr/sbin/init

```

2.2 Shell脚本启动方式

编写脚本 `hadoop-cluster.sh`，批量启动容器

```

1 #!/bin/bash
2 # description:Batch start Containers Script (批量启动容
   器)
3 # author:will
4
5 # 预判启动容器的数量：不小于1 && 提示信息
6 if [ $# -ne 1 ];then
7     echo "You need to start several containers
   explicitly (批量启动容器的方法如下) ."
8     echo "Some like './hadoop-cluster.sh 3' or 'sh
   hadoop-cluster.sh 3'"
9     exit 1
10 fi
11
12 #要启动的容器数量
13 NUM_CONTAINERS=$1
14 #自定义网络名称
15 NETWORK_NAME=will-network
16 #镜像ID
17 IMAGE_ID=2d7d22776b99
18 #前缀
19 PREFIX="0"
20 for (( i=1;i<=$NUM_CONTAINERS;i++ ))
21     do

```

```

22         if [ $i -eq 1 ];then
23             sudo docker run -d --name
hadoop$PREFIX$i --hostname hadoop$PREFIX$i --net
${NETWORK_NAME} --ip 172.24.0.${i+1} -P -p 50070:50070
-p 8088:8088 -p 19888:19888 --privileged $IMAGE_ID
/usr/sbin/init
24         elif [ $i -eq 2 ];then
25             sudo docker run -d --name
hadoop$PREFIX$i --hostname hadoop$PREFIX$i --net
${NETWORK_NAME} --ip 172.24.0.${i+1} -P -p 50090:50090
--privileged $IMAGE_ID /usr/sbin/init
26         else
27             sudo docker run -d --name
hadoop$PREFIX$i --hostname hadoop$PREFIX$i --net
${NETWORK_NAME} --ip 172.24.0.${i+1} -P --privileged
$IMAGE_ID /usr/sbin/init
28         fi
29     done
30 echo "$NUM_CONTAINERS containers started!"
31 echo "=====
32 sudo docker ps | grep hadoop
33 echo "=====IPv4=====
34 sudo docker inspect $(sudo docker ps -q) | grep -i ipv4
35 echo "=====

```

使用shell脚本启动：

```

[will@master ~]$ sh ./Desktop/hadoop-cluster.sh 3
2840cee259340a22e910af7e1cbd35be060d84be92989dab5238e60dd74b6d9c
7485d39d5df39950b8f17e7dbf3f2dc489c029283925d586f0b4b1d9ea124bdf
4bd249d461cb7b3969d58d6dc0617d534c38b5d1f122edb0a89c14f98ebfe938
3 containers started!
=====
4bd249d461cb      2d7d22776b99      "/usr/sbin/init"   4 seconds ago      Up Less than a second   0.0.0.0: 32772->22/tcp
                                hadoop03
7485d39d5df3      2d7d22776b99      "/usr/sbin/init"   8 seconds ago      Up 4 seconds            0.0.0.0: 50090->50090/tcp, 0.0.0.
0: 32771->22/tcp
                                hadoop02
2840cee25934      2d7d22776b99      "/usr/sbin/init"   8 seconds ago      Up 7 seconds            0.0.0.0: 8088->8088/tcp, 0.0.0.0:
19888->19888/tcp, 0.0.0.0: 50070->50070/tcp, 0.0.0.0: 32770->22/tcp
                                hadoop01
=====IPv4=====
                                "IPv4Address": "172.24.0.4"
                                "IPv4Address": "172.24.0.3"
                                "IPv4Address": "172.24.0.2"

```

【受虚拟机性能所限制，只能同时启动三个容器】

3. 进入容器进行集群配置

3.1 进入容器

```
1 sudo docker exec -ti hadoop01 /bin/bash
2
3 sudo docker exec -ti hadoop02 /bin/bash
4
5 sudo docker exec -ti hadoop03 /bin/bash
6
7 sudo docker exec -ti hadoop04 /bin/bash
8
9 sudo docker exec -ti hadoop05 /bin/bash
10
11 # 每进入一个之后，还要进入 will用户
12 su will
```

3.2 配置映射

配置ip地址和主机名映射（所有节点都要配置）

```
1 172.24.02 hadoop01
2 172.24.03 hadoop02
3 172.24.04 hadoop03
4 172.24.05 hadoop04
5 172.24.06 hadoop05
```

```
1 # 修改
2 vim /etc/hosts
3
4 # 删除最后一行
5
6 # 插入
7 172.24.02 hadoop01
8 172.24.03 hadoop02
9 172.24.04 hadoop03
10 172.24.05 hadoop04
11 172.24.06 hadoop05
```

3.3 配置免密登录

3.3.1 生成密钥

在每一台主机上使用 `ssh-keygen` 命令生成公钥和私钥

```
1 ssh-keygen -t rsa -N '' -C 'willl'
2 # -N ''      指定：不需要密码
3 # -N和-C可写，也可以不写
4
5 # 然后一路「回车」
```

3.3.2 复制公钥

每一个节点都要把公钥给其他节点（每个容器中都执行一遍下面的内容）

```
1 ssh-copy-id -i ~/.ssh/id_rsa.pub hadoop01
2 ssh-copy-id -i ~/.ssh/id_rsa.pub hadoop02
3 ssh-copy-id -i ~/.ssh/id_rsa.pub hadoop03
4 ssh-copy-id -i ~/.ssh/id_rsa.pub hadoop04
5 ssh-copy-id -i ~/.ssh/id_rsa.pub hadoop05
```

3.3.3 验证免密登录

```
1 ssh hadoop01
2 ssh hadoop02
3 ssh hadoop03
4 ssh hadoop04
5 ssh hadoop05
```

3.4 同步集群时间（hadoop02）

思路：

hadoop02开启ntp服务；其他节点同步hadoop02

3.4.1 检查ntp包是否安装

```
1 rpm -qa | grep ntp
2
3 # 如果没有安装的话, 执行 sudo yum -y install ntp
```

hadoop02的操作:

3.4.2 设置时间文件

```
1 sudo vim /etc/ntp.conf
```

```
1 #修改一 (设置本地网络上的主机不受限制)
2 #restrict 192.168.1.0 mask 255.255.255.0 nomodify
   notrap
3 restrict 172.24.0.0 mask 255.255.255.0 nomodify notrap
4 #修改二 (添加默认的一个内部时钟数据, 使用它为局域网用户提供服务)
5 server 127.127.1.0
6 fudge 127.127.1.0 stratum 10
7 #修改三 (设置为不采用公共的服务器)
8 server 0.centos.pool.ntp.org iburst
9 server 1.centos.pool.ntp.org iburst
10 server 2.centos.pool.ntp.org iburst
11 server 3.centos.pool.ntp.org iburst
12 #server 0.centos.pool.ntp.org iburst
13 #server 1.centos.pool.ntp.org iburst
14 #server 2.centos.pool.ntp.org iburst
15 #server 3.centos.pool.ntp.org iburst
```

3.4.3 设置BIOS与系统时间同步

```
1 sudo vim /etc/sysconfig/ntpd
2
3 # 增加下面的内容
4 OPTIONS="-u ntp:ntp -p /var/run/ntpd.pid -g"
5 SYNC_HWCLOCK=yes
```

3.4.4 启动ntp服务并测试

```
1  sudo systemctl start ntpd
2
3  systemctl status ntpd
4
5  sudo systemctl enable ntpd.service
6
7  # 用来查看ntp状态（下同）
8  ntpstat
9
10 sudo ntpdq -p
```

3.4.5 关闭非时间服务器节点的ntpd服务

```
1  sudo systemctl stop ntpd
```

手动同步（所有节点）

```
1  # 手动设置时区
2  sudo timedatectl set-timezone Asia/Shanghai
```

其他节点（非hadoop02）的操作

定时同步

```
1  #非时间服务器都安装crond服务
2  sudo yum -y install vixie-cron crontabs
3  #非时间服务器节点都编写定时同步时间任务
4  sudo vim /etc/crontab
5  编写定时任务如下：
6  */1 * * * * /usr/sbin/ntpdate hadoop02
7
8  #加载任务,使之生效
9  sudo crontab /etc/crontab
```

```
1  sudo ntpdate -u ntp.api.bz
```


4. 修改主节点配置文件

hadoop01中的操作;

在hadoop目录下。

```
1 mkdir /opt/software/hadoop-2.7.3/tmp
2
3 mkdir -p /opt/software/hadoop-2.7.3/dfs/namenode_data
4
5 mkdir -p /opt/software/hadoop-2.7.3/dfs/datanode_data
6
7 mkdir -p /opt/software/hadoop-2.7.3/checkpoint/dfs/cname
```

4.1 hadoop-env.sh

```
1 #25行 export JAVA_HOME
2 export JAVA_HOME=/opt/moudle/jdk1.8.0_251
3 #33行 export HADOOP_CONF_DIR
4 export HADOOP_CONF_DIR=/opt/software/hadoop-
  2.7.3/etc/hadoop
```

4.2 core-site.xml

```
1 <configuration>
2     <property>
3         <!--用来指定hdfs的主节点，namenode的地址-->
4         <name>fs.defaultFS</name>
5         <value>hdfs://hadoop01:9000</value>
6     </property>
7     <property>
8         <!--用来指定hadoop运行时产生文件的存放目录-->
9         <name>hadoop.tmp.dir</name>
10        <value>/opt/software/hadoop-
  2.7.3/tmp</value>
11    </property>
12    <property>
```

```

13         <!--设置缓存大小, 默认4kb-->
14         <name>io.file.buffer.size</name>
15         <value>4096</value>
16     </property>
17 </configuration>

```

4.3 hdfs-site.xml

```

1 <configuration>
2     <property>
3         <!--数据块默认大小128M-->
4         <name>dfs.block.size</name>
5         <value>134217728</value>
6     </property>
7     <property>
8         <!--副本数量, 不配置的话默认为3-->
9         <name>dfs.replication</name>
10        <value>3</value>
11    </property>
12    <property>
13        <!--定点检查-->
14        <name>fs.checkpoint.dir</name>
15        <value>/opt/software/hadoop-
16        2.7.3/checkpoint/dfs/cname</value>
17    </property>
18    <property>
19        <!--namenode节点数据 (元数据) 的存放位置-->
20        <name>dfs.name.dir</name>
21        <value>/opt/software/hadoop-
22        2.7.3/dfs/namenode_data</value>
23    </property>
24    <property>
25        <!--datanode节点数据 (元数据) 的存放位置-->
26        <name>dfs.data.dir</name>
27        <value>/opt/software/hadoop-
28        2.7.3/dfs/datanode_data</value>
29    </property>

```

```

27         <property>
28             <!--指定secondarynamenode的web地址-->
29             <name>dfs.namenode.secondary.http-
address</name>
30             <value>hadoop02:50090</value>
31         </property>
32         <property>
33             <!--hdfs文件操作权限,false为不验证-->
34             <name>dfs.permissions</name>
35             <value>false</value>
36         </property>
37 </configuration>

```

4.4 mapreduce-site.xml

在\${HADOOP_HOME}/etc/hadoop的目录下，只有一个mapred-site.xml.template文件，复制一个进行更改。

```

1  cp mapred-site.xml.template mapred-site.xml
2
3  <configuration>
4      <property>
5          <!--指定mapreduce运行在yarn上-->
6          <name>mapreduce.framework.name</name>
7          <value>yarn</value>
8      </property>
9      <property>
10         <!--配置任务历史服务器地址-->
11         <name>mapreduce.jobhistory.address</name>
12         <value>hadoop01:10020</value>
13     </property>
14     <property>
15         <!--配置任务历史服务器web-UI地址-->
16
17         <name>mapreduce.jobhistory.webapp.address</name>
18         <value>hadoop01:19888</value>
19     </property>

```

4.5 yarn-site.xml

```
1 <configuration>
2     <property>
3         <!--指定yarn的老大resourcemanager的地址-->
4         <name>yarn.resourcemanager.hostname</name>
5         <value>hadoop01</value>
6     </property>
7     <property>
8         <name>yarn.resourcemanager.address</name>
9         <value>hadoop01:8032</value>
10    </property>
11    <property>
12
13        <name>yarn.resourcemanager.webapp.address</name>
14        <value>hadoop01:8088</value>
15    </property>
16    <property>
17
18        <name>yarn.resourcemanager.scheduler.address</name>
19        <value>hadoop01:8030</value>
20    </property>
21    <property>
22        <name>yarn.resourcemanager.resource-
23        tracker.address</name>
24        <value>hadoop01:8031</value>
25    </property>
26    <property>
27
28        <name>yarn.resourcemanager.admin.address</name>
29        <value>hadoop01:8033</value>
30    </property>
31    <property>
32
33        <!--NodeManager获取数据的方式-->
34        <name>yarn.nodemanager.aux-services</name>
```

```
30         <value>mapreduce_shuffle</value>
31     </property>
32     <property>
33         <!--开启日志聚集功能-->
34         <name>yarn.log-aggregation-enable</name>
35         <value>true</value>
36     </property>
37     <property>
38         <!--配置日志保留7天-->
39         <name>yarn.log-aggregation.retain-
seconds</name>
40         <value>604800</value>
41     </property>
42 </configuration>
```

4.6 master

```
1 # 在 hadoop目录下
2 vim master
3 # 添加内容
4 hadoop01
```

4.7 slaves

```
1 # 在 hadoop目录下
2 vim slave
3 # 添加内容
4 hadoop02
5 hadoop03
6 hadoop04
7 hadoop05
```

4.8 分发配置文件

hadoop01分发到其他机器上去

手动

```
1 sudo scp -r /opt/software/hadoop-2.7.3/
  will@hadoop02:/opt/software/
2 sudo scp -r /opt/software/hadoop-2.7.3/
  will@hadoop03:/opt/software/
3 sudo scp -r /opt/software/hadoop-2.7.3/
  will@hadoop04:/opt/software/
4 sudo scp -r /opt/software/hadoop-2.7.3/
  will@hadoop05:/opt/software/
```

shell脚本 scp-config.sh

```
1 #!/bin/bash
2 #description: 节点间复制文件
3 #author: will
4
5 #首先判断参数是否存在
6 args=$#
7 if [ args -eq 0 ];then
8     echo "no args"
9     exit 1
10 fi
11 #获取文件名称
12 p1=$1
13 fname=$(basename $p1)
14 echo fname=$fname
15 #获取上级目录到绝对路径
16 pdir=$(cd $(dirname $p1);pwd -P)
17 echo pdir=$pdir
18 #获取当前用户名称
19 user=$(whoami)
20 #循环分发
21 # 【注意，如果节点数量更多】当前是，hadoop02 ~ hadoop05节点
22 for(( host=2;host<6;host++ ));do
23     echo "-----hadoop0$host-----"
24     scp -r $pdir/$fname $user@hadoop0$host:$pdir
25 done
26 echo "-----All done -----"
```

```
1 # 使用shell脚本
2 sh scp-config.sh /opt/software/hadoop-2.7.3/
```

4.9 初始化

hadoop01

```
1 hdfs namenode -format
```

5 启动集群

在hadoop01上启动集

群，hadoop02、hadoop03、hadoop04、hadoop05上相关服务也会被启动

```
1 # hadoop01上操作
2 start-dfs.sh
3
4 start-yarn.sh
5
6 mr-jobhistory-daemon.sh start historyserver
```

每一台服务器上jps查看即可

参考资料及笔记

[hadoop搭建完全分布式（文案）](#)

[hadoop搭建完全分布式（视频）](#)

[hadoop基本原理（读写 && mapreduce）](#)

NameNode：是Master节点（主节点），可以看作是分布式文件系统的管理者，主要负责管理文件系统的命名空间、集群配置信息和存储块的复制等（HDFS的元数据信息）。NameNode会将文件系统的Meta-data存储在内存中，这些信息主要包括了文件信息、每一个文件对应的文件块的信息和每一个文件块在DataNode的信息等（HDFS的元数据信息的持久化）。

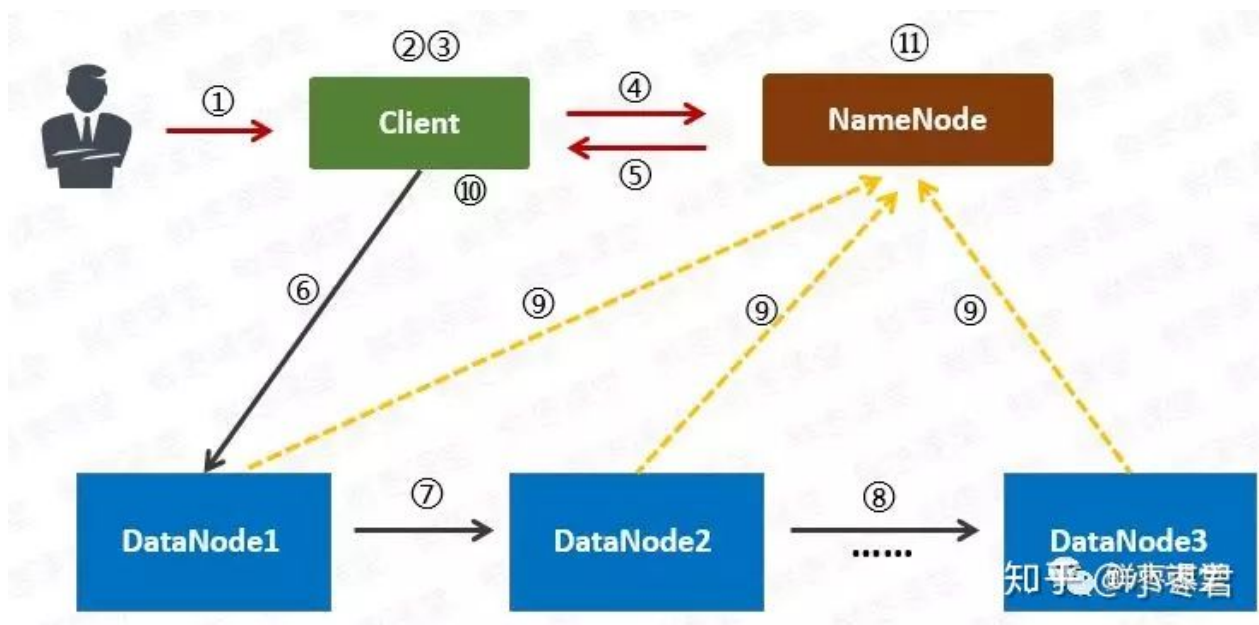
DataNode：是Slave节点（从节点），是文件存储的基本单元，它将Block存储在本地文件系统中，保存了Block的Meta-data，同时周期性地将所有存在的Block信息发送给NameNode。

Client：切分文件；访问HDFS；与NameNode交互，获得文件位置信息；与DataNode交互，读取和写入数据。

还有一个**Block（块）**的概念：Block是HDFS中的基本读写单元；HDFS中的文件都是被切割为block（块）进行存储的；这些块被复制到多个DataNode中；块的大小（通常为64MB）和复制的块数量在创建文件时由Client决定。

我们来简单看看HDFS的读写流程。

首先是写入流程：



1 用户向Client（客户机）提出请求。例如，需要写入200MB的数据。

2 Client制定计划：将数据按照64MB为块，进行切割；所有的块都保存三份。

3 Client将大文件切分成块（block）。

4 针对第一个块，Client告诉NameNode（主控节点），请帮助我，将64MB的块复制三份。

5 NameNode告诉Client三个DataNode（数据节点）的地址，并且将它们根据到Client的距离，进行了排序。

6 Client把数据和清单发给第一个DataNode。

7 第一个DataNode将数据复制给第二个DataNode。

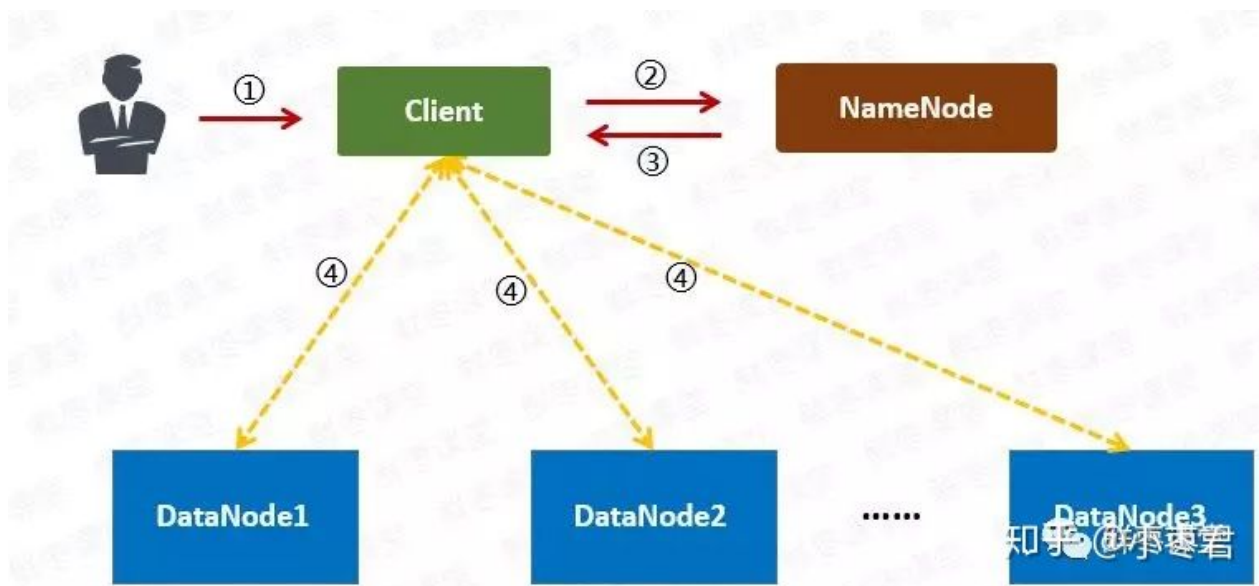
8 第二个DataNode将数据复制给第三个DataNode。

9 如果某一个块的所有数据都已写入，就会向NameNode反馈已完成。

10 对第二个Block，也进行相同的操作。

11 所有Block都完成后，关闭文件。NameNode会将数据持久化到磁盘上。

读取流程：



1 用户向Client提出读取请求。

2 Client向NameNode请求这个文件的所有信息。

3 NameNode将给Client这个文件的块列表，以及存储各个块的数据节点清单（按照和客户端的距离排序）。

4 Client从距离最近的数据节点下载所需的块。

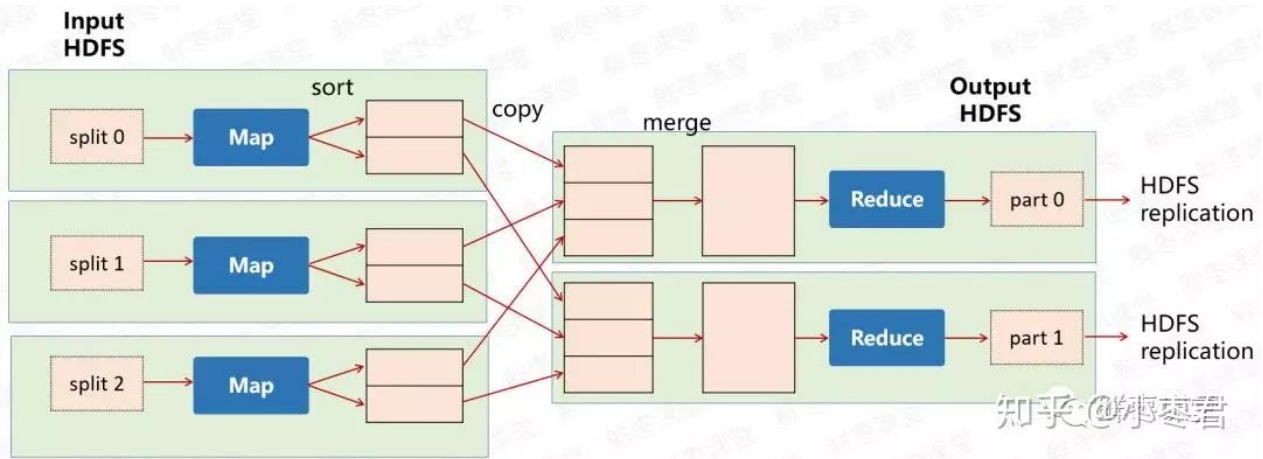
（注意：以上只是简化的描述，实际过程会更加复杂。）

再来看**MapReduce**。

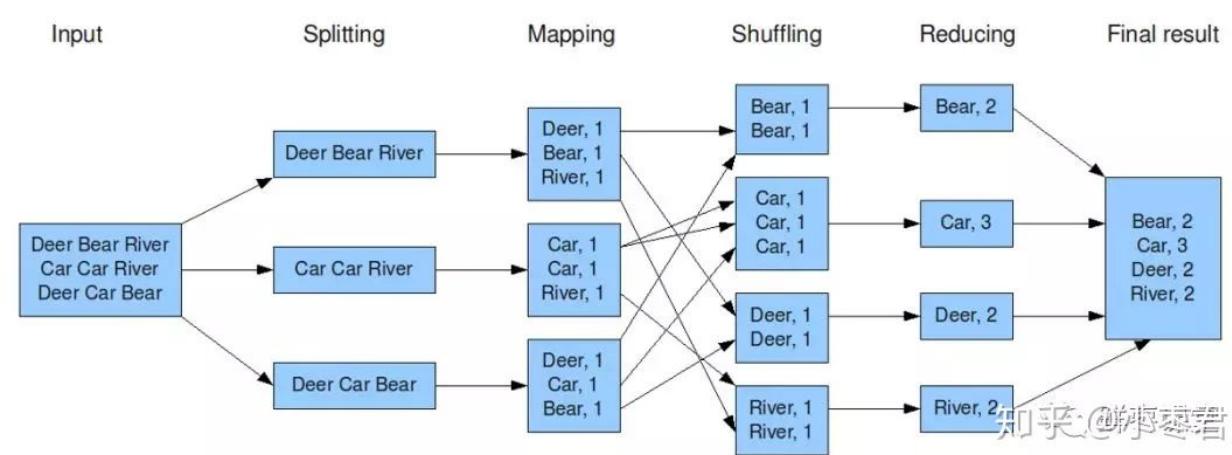
MapReduce其实是一种编程模型。这个模型的核心步骤主要分两部分：**Map**（映射）和**Reduce**（归约）。

当你向MapReduce框架提交一个计算作业时，它会首先把计算作业拆分成若干个**Map任务**，然后分配到不同的节点上去执行，每一个Map任务处理输入数据中的一部分，当Map任务完成后，它会生成一些中间文件，这些中间文件将会作为**Reduce任务**的输入数据。Reduce任务的主要目标就是

把前面若干个Map的输出汇总到一起并输出。



例子：单词统计（WordCount）



1 Hadoop将输入数据切成若干个分片，并将每个split（分割）交给一个map task（Map任务）处理。

2 Mapping之后，相当于得出这个task里面，每个词以及它出现的次数。

3 shuffle（拖移）将相同的词放在一起，并对它们进行排序，分成若干个分片。

4 根据这些分片，进行reduce（归约）。

5 统计出reduce task的结果，输出到文件。

节点解释

ResourceManager：资源管理，**yarn**集群的主节点。功能上：与客户端交流，处理客户端的请求；管理NodeManager，接收来自NodeManager的资源汇报信息，并向NodeManager下达管理指令。

JobHistoryServer：**yarn**查看已经完成的任务的历史日志记录的服务。功能上：查看

NodeManager：节点管理，**yarn**集群的从节点。功能上：任务的计算。

SecondaryNameNode：镜像备份；日志与镜像的定期合并。

NTP-Server：时间服务器。功能上：保证集群内所有主机的时间同步。