

# Hadoop优化&新特性

## 1 Hadoop数据压缩

### 1.1 概述

- 概述

压缩技术能够有效减少底层存储系统（HDFS）读写字节数。压缩提高了网络带宽和磁盘空间的效率。在运行MR程序时，I/O操作、网络数据传输、Shuffle和Merge要花大量的时间，尤其是数据规模很大和工作负载密集的情况下，因此，使用数据压缩显得非常重要。

鉴于磁盘I/O和网络带宽是Hadoop的宝贵资源，数据压缩对于节省资源、最小化磁盘I/O和网络传输非常有帮助。可以在任意MapReduce阶段启用压缩。不过，尽管压缩与解压操作的CPU开销不高，其性能的提升和资源的节省并非没有代价。

- 压缩策略和原则

压缩是提高Hadoop运行效率的一种优化策略。

通过对Mapper、Reducer运行过程的数据进行压缩，以减少磁盘IO，提高MR程序运行速度。

注意：采用压缩技术减少了磁盘IO，但同时增加了CPU运算负担。所以，压缩特性运用得当能提高性能，但运用不当也可能降低性能。

**压缩基本原则：**（1）运算密集型的job，少用压缩。（2）IO密集型的job，多用压缩

### 1.2 MR支持的压缩编码

压缩格式	hadoop自带?	算法	文件扩展名	是否可切分	换成压缩格式后，原来的程序是否需要修改
DEFLATE	是，直接使用	DEFLATE	.deflate	否	和文本处理一样，不需要修改
Gzip	是，直接使用	DEFLATE	.gz	否	和文本处理一样，不需要修改
bzip2	是，直接使用	bzip2	.bz2	是	和文本处理一样，不需要修改
LZO	否，需要安装	LZO	.lzo	是	需要建索引，还需要指定输入格式
Snappy	是，直接使用	Snappy	.snappy	否	和文本处理一样，不需要修改

为了支持多种压缩/解压缩算法，Hadoop引入了编码/解码器，如下表所示。

压缩格式	对应的编码/解码器
DEFLATE	org.apache.hadoop.io.compress.DefaultCodec
gzip	org.apache.hadoop.io.compress.GzipCodec
bzip2	org.apache.hadoop.io.compress.BZip2Codec
LZO	com.hadoop.compression.lzo.LzopCodec
Snappy	org.apache.hadoop.io.compress.SnappyCodec

压缩性能的比较

压缩算法	原始文件大小	压缩文件大小	压缩速度	解压速度
gzip	8.3GB	1.8GB	17.5MB/s	58MB/s
bzip2	8.3GB	1.1GB	2.4MB/s	9.5MB/s
LZO	8.3GB	2.9GB	49.3MB/s	74.6MB/s

## 1.3 压缩方式选择

### 1.3.1 Gzip压缩

优点：压缩率比较高，而且压缩/解压速度也比较快；Hadoop本身支持，在应用中处理Gzip格式的文件就和直接处理文本一样；大部分Linux系统都自带Gzip命令，使用方便。

缺点：不支持Split。

应用场景：当每个文件压缩之后在130M以内的（1个块大小内），都可以考虑用Gzip压缩格式。例如说一天或者一个小时的日志压缩成一个Gzip文件。

### 1.3.2 Bzip2压缩

优点：支持Split；具有很高的压缩率，比Gzip压缩率都高；Hadoop本身自带，使用方便

缺点：压缩/解压速度慢。

应用场景：适合对速度要求不高，但需要较高的压缩率的时候；或者输出之后的数据比较大，处理之后的数据需要压缩存档减少磁盘空间并且以后数据用得比较少的情况；或者对单个很大的文本文件想压缩减少存储空间，同时又需要支持Split，而且兼容之前的应用程序的情况。

### 1.3.3 Lzo压缩

优点：压缩/解压速度也比较快，合理的压缩率；支持Split，是Hadoop中最流行的压缩格式；可以在Linux系统下安装lzop命令，使用方便。

缺点：压缩率比Gzip要低一些；Hadoop本身不支持，需要安装；在应用中对Lzo格式的文件需要做一些特殊处理（为了支持Split需要建索引，还需要指定InputFormat为Lzo格式）。

应用场景：一个很大的文本文件，压缩之后还大于200M以上的可以考虑，而且单个文件越大，Lzo优点越越明显。

1.3.4 Snappy压缩

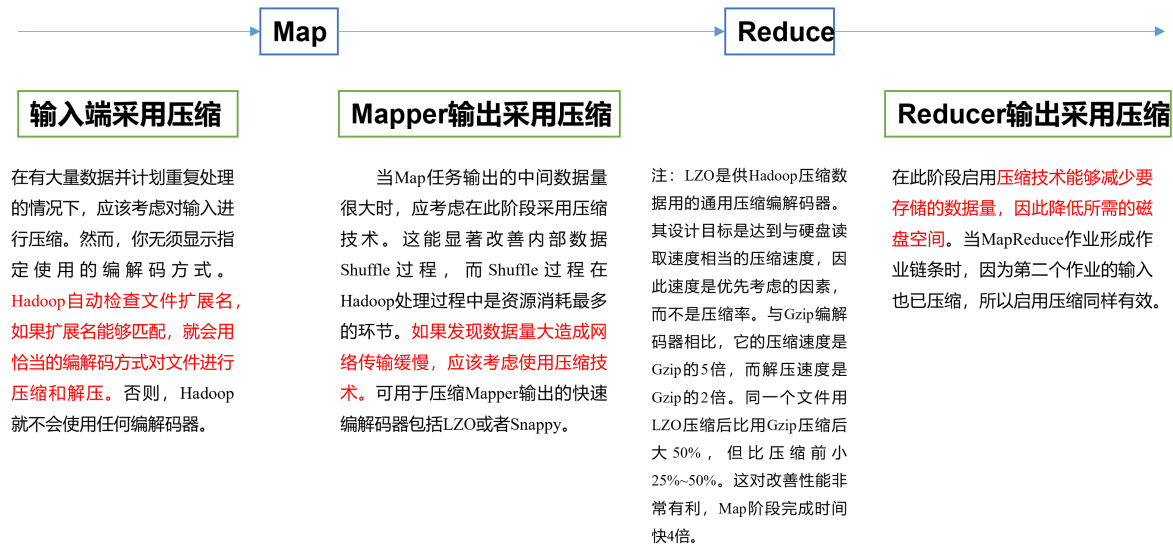
优点：高速压缩速度和合理的压缩率。

缺点：不支持Split；压缩率比Gzip要低；Hadoop本身不支持，需要安装。

应用场景：当MapReduce作业的Map输出的数据比较大的时候，作为Map到Reduce的中间数据的压缩格式；或者作为一个MapReduce作业的输出和另外一个MapReduce作业的输入。

1.4 压缩位置选择

压缩可以在MapReduce作用的任意阶段启用。



1.5 压缩参数配置

要在Hadoop中启用压缩，可以配置如下参数：

参数	默认值	阶段	建议
io.compression.codecs (在core-site.xml中配置)	无，这个需要在命令行输入hadoop checknative查看	输入压缩	Hadoop使用文件扩展名判断是否支持某种编解码器
mapreduce.map.output.compress (在mapred-site.xml中配置)	false	mapper输出	这个参数设为true启用压缩
mapreduce.map.output.compress.codec (在mapred-site.xml中配置)	org.apache.hadoop.io.compress.DefaultCodec	mapper输出	企业多使用LZO或Snappy编解码器在此阶段压缩数据
mapreduce.output.fileoutputformat.compress (在mapred-site.xml中配置)	false	reducer输出	这个参数设为true启用压缩
mapreduce.output.fileoutputformat.compress.codec (在mapred-site.xml中配置)	org.apache.hadoop.io.compress.DefaultCodec	reducer输出	使用标准工具或者编解码器，如gzip和bzip2
mapreduce.output.fileoutputformat.compress.type (在mapred-site.xml中配置)	RECORD	reducer输出	SequenceFile输出使用的压缩类型：NONE和BLOCK

1.6 压缩实操案例

1.6.1 数据流的压缩和解压缩

CompressionCodec有两个方法可以用于轻松地压缩或解压缩数据。

要想对正在被写入一个输出流的数据进行压缩，我们可以使用createOutputStream(OutputStreamout)方法创建一个CompressionOutputStream，将其以压缩格式写入底层的流。

相反，要想对从输入流读取而来的数据进行解压缩，则调用createInputStream(InputStream in)函数，从而获得一个CompressionInputStream，从而从底层的流读取未压缩的数据。

测试一下如下压缩方式：

DEFLATE	org.apache.hadoop.io.compress.DefaultCodec
gzip	org.apache.hadoop.io.compress.GzipCodec
bzip2	org.apache.hadoop.io.compress.BZip2Codec

```
package com.xu1an.mr.compress;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IOUtils;
import org.apache.hadoop.io.compress.CompressionCodec;
import org.apache.hadoop.io.compress.CompressionCodecFactory;
import org.apache.hadoop.io.compress.CompressionInputStream;
import org.apache.hadoop.io.compress.CompressionOutputStream;
import org.apache.hadoop.util.ReflectionUtils;
import org.junit.Test;

import java.io.*;

/**
 * 对文件进行压缩和解压缩操作
 */
public class CompressTest {

    // 准备原始文件的路径
    String srcFile = "F:\\in\\jianai\\ja.txt.deflate";
    // 准备压缩后的文件路径
    String destFile = "F:\\in\\jianai\\ja.txt";

    /**
     * 压缩：通过一个能够具备压缩功能输出流将文件写出
     */
    @Test
    public void testCompress() throws IOException, ClassNotFoundException {

        //声明一个输入流
        FileInputStream in = new FileInputStream(new File(srcFile));
        Configuration conf = new Configuration();
        String classPath = "org.apache.hadoop.io.compress.DefaultCodec";
        Class<?> codecClass = Class.forName(classPath);
        // 获取一个编解码器（压缩工具对象）
        CompressionCodec codec = (CompressionCodec)
        ReflectionUtils.newInstance(codecClass, conf);
        // 声明一个输出流，将文件写出
        FileOutputStream out = new FileOutputStream(new File(destFile +
        codec.getDefaultExtension()));
        // 把普通的输出流让 CompressionCodec 包装一下
        CompressionOutputStream outputStream = codec.createOutputStream(out);

        // 读写数据
        IOUtils.copyBytes(in, outputStream, conf);
    }
}
```

```

        // 关闭流
        IOUtils.closeStream(in);
        IOUtils.closeStream(outputStream);

    }

    /**
     * 解压缩：通过一个能够具备解压缩功能输入流将文件写出
     */
    @Test
    public void testCodecCompress() throws IOException, ClassNotFoundException {

        //声明一个输入流
        FileInputStream in = new FileInputStream(new File(srcFile));
        Configuration conf = new Configuration();
        // 获取一个编解码器（解压缩工具对象）
        CompressionCodec codec =
            new CompressionCodecFactory(conf).getCodec(new Path(srcFile));
        // 把普通的输入流让 CompressionCodec 包装一下
        CompressionInputStream inputStream = codec.createInputStream(in);
        // 声明一个输出流，将文件写出
        FileOutputStream out = new FileOutputStream(new File(destFile));

        // 读写数据
        IOUtils.copyBytes(inputStream, out, conf);

        // 关闭流
        IOUtils.closeStream(inputStream);
        IOUtils.closeStream(out);

    }
}

```

## 1.6.2 Map输出端采用压缩

即使你的MapReduce的输入输出文件都是未压缩的文件，你仍然可以对Map任务的中间结果输出做压缩，因为它要写在硬盘并且通过网络传输到Reduce节点，对其压缩可以提高很多性能，这些工作只要设置两个属性即可，我们来看下代码怎么设置。

1) 给大家提供的Hadoop源码支持的压缩格式有：BZip2Codec、DefaultCodec

//见reduce输出端采用压缩案例

2) Mapper保持不变

3) Reducer保持不变

### 1.6.3 Reduce输出端采用压缩

基于WordCount案例处理。

#### 1) 修改驱动

```
package com.xulan.mr.wordcount;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.CombineFileInputFormat;
import org.apache.hadoop.mapreduce.lib.input.CombineTextInputFormat;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

import java.io.IOException;

/**
 * Created with IntelliJ IDEA.
 *
 * @Author: Xu1Aan
 * @Date: 2022/03/15/20:36
 * @Description:
 * MR程序的驱动类，主要用于提交MR任务
 */
public class WordCountDriver {
    public static void main(String[] args) throws IOException,
        InterruptedException, ClassNotFoundException {

        //声明配置对象
        Configuration conf = new Configuration();
        //指定当前job提交队列名称
        //conf.set("mapreduce.job.queueName", "hello");
        //1、设置在Mapper端输出进行压缩
        conf.set("mapreduce.map.output.compress", "true");
        //或者 conf.setBoolean("mapreduce.map.output.compress", true);

        //设置编解码器（默认的）

        conf.set("mapreduce.map.output.compress.codec", "org.apache.hadoop.io.compress.D
efaultCodec");
        //声明Job对象
        Job job = Job.getInstance(conf);
        //指定当前Job的驱动类
        job.setJarByClass(WordCountDriver.class);
        //指定当前Job的Mapper和Reducer
        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);
        //指定Map端输出数据key的类型和输出数据value的类型
        job.setMapOutputKeyClass(Text.class);
        job.setMapOutputValueClass(IntWritable.class);
        //指定最终输出结果的key的类型和value的类型
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
    }
}
```

```
//指定CombineFileInputFormat中切片的最大值
//      CombineFileInputFormat.setMaxInputSplitSize(job,4194304 * 5);
//      //指定InputFormat的实现，默认为FileInputFormat
//      job.setInputFormatClass(CombineTextInputFormat.class);

//指定输入数据的目录 和 输出数据的目录
FileInputFormat.setInputPaths(job,new Path("E:\\learning\\04_java\\02_大
数据资料\\00_hadoop\\资料\\07_测试数据\\jianai"));
FileOutputFormat.setOutputPath(job, new Path("E:\\learning\\04_java\\02_
大数据资料\\00_hadoop\\out\\data4"));
//      FileInputFormat.setInputPaths(job, new Path(args[0]));
//      FileOutputFormat.setOutputPath(job, new Path(args[1]));

//提交job
job.waitForCompletion(true);
}
}
```

2) Mapper和Reducer保持不变

---

## 2 Hadoop企业优化

### 2.1 MapReduce 跑的慢的原因

MapReduce 程序效率的瓶颈在于两点：

1. 计算机性能  
CPU、内存、磁盘健康、网络
2. I/O 操作优化
  - (1) 数据倾斜
  - (2) Map和Reduce数设置不合理
  - (3) Map运行时间太长，导致Reduce等待过久
  - (4) 小文件过多
  - (5) 大量的不可切片的超大压缩文件
  - (6) Spill次数过多
  - (7) Merge次数过多等。

### 2.2 MapReduce优化方法

MapReduce优化方法主要从六个方面考虑：数据输入、Map阶段、Reduce阶段、IO传输、数据倾斜问题和常用的调优参数。

#### 2.2.1 数据输入

(1) 合并小文件：在执行MR任务前将小文件进行合并，大量的小文件会产生大量的Map任务，增大Map任务装载次数，而任务的装载比较耗时，从而导致MR运行较慢。

(2) 采用CombineTextInputFormat来作为输入，解决输入端大量小文件场景。

## 2.2.2 Map阶段

(1) 减少溢写 (Spill) 次数: 通过调整`mapreduce.task.io.sort.mb`及`mapreduce.map.sort.spill.percent`参数值, 增大触发Spill的内存上限, 减少Spill次数, 从而减少磁盘IO。

(2) 减少合并 (Merge) 次数: 通过调整`mapreduce.task.io.sort.factor`参数, 增大Merge的文件数目, 减少Merge的次数, 从而缩短MR处理时间。

(3) 在Map之后, 不影响业务逻辑前提下, 先进行Combine处理, 减少 I/O。

## 2.2.3 Reduce阶段

(1) 合理设置Map和Reduce数: 两个都不能设置太少, 也不能设置太多。太少, 会导致Task等待, 延长处理时间; 太多, 会导致Map、Reduce任务间竞争资源, 造成处理超时等错误。

(2) 设置Map、Reduce共存: 调整`mapreduce.job.reduce.slowstart.completedmaps`参数, 使Map运行到一定程度后, Reduce也开始运行, 减少Reduce的等待时间。

(3) 规避使用Reduce: 因为Reduce在用于连接数据集的时候将会产生大量的网络消耗。

(4) 合理设置Reduce端的Buffer: 默认情况下, 数据达到一个阈值的时候, Buffer中的数据就会写入磁盘, 然后Reduce会从磁盘中获得所有的数据。也就是说, Buffer和Reduce是没有直接关联的, 中间多次写磁盘->读磁盘的过程, 既然有这个弊端, 那么就可以通过参数来配置, 使得Buffer中的一部分数据可以直接输送到Reduce, 从而减少IO开销: `mapreduce.reduce.input.buffer.percent`, 默认为0.0。当值大于0的时候, 会保留指定比例的内存读Buffer中的数据直接拿给Reduce使用。这样一来, 设置Buffer需要内存, 读取数据需要内存, Reduce计算也要内存, 所以要根据作业的运行情况进行调整。

## 2.2.4 I/O传输

(1) 采用数据压缩的方式, 减少网络IO的时间。安装Snappy和LZO压缩编码器。

(2) 使用SequenceFile二进制文件。

## 2.2.5 数据倾斜问题

### 1. 数据倾斜现象

数据频率倾斜——某一个区域的数据量要远远大于其他区域。

数据大小倾斜——部分记录的大小远远大于平均值。

### 2. 减少数据倾斜的方法

方法1: 抽样和范围分区

可以通过对原始数据进行抽样得到的结果集来预设分区边界值。

方法2: 自定义分区

基于输出键的背景知识进行自定义分区。例如, 如果Map输出键的单词来源于一本书。且其中某几个专业词汇较多。那么就可以自定义分区将这这些专业词汇发送给固定的一部分Reduce实例。而将其他的都发送给剩余的Reduce实例。

方法3: Combiner

使用Combiner可以大量地减小数据倾斜。在可能的情况下, Combine的目的就是聚合并精简数据。

方法4: 采用Map Join, 尽量避免Reduce Join。



## 2.3 常用的调优参数

### (1) 资源相关参数

(a) 以下参数是在用户自己的MR应用程序中配置就可以生效 (mapred-default.xml)

配置参数	参数说明
mapreduce.map.memory.mb	一个MapTask可使用的资源上限 (单位:MB)，默认为1024。如果MapTask实际使用的资源量超过该值，则会被强制杀死。
mapreduce.reduce.memory.mb	一个ReduceTask可使用的资源上限 (单位:MB)，默认为1024。如果ReduceTask实际使用的资源量超过该值，则会被强制杀死。
mapreduce.map.cpu.vcores	每个MapTask可使用的最多cpu core数目，默认值: 1
mapreduce.reduce.cpu.vcores	每个ReduceTask可使用的最多cpu core数目，默认值: 1
mapreduce.reduce.shuffle.parallelcopies	每个Reduce去Map中取数据的并行数。默认值是5
mapreduce.reduce.shuffle.merge.percent	Buffer中的数据达到多少比例开始写入磁盘。默认值0.66
mapreduce.reduce.shuffle.input.buffer.percent	Buffer大小占Reduce可用内存的比例。默认值0.7
mapreduce.reduce.input.buffer.percent	指定多少比例的内存用来存放Buffer中的数据，默认值是0.0

(b) 应该在YARN启动之前就配置在服务器的配置文件中才能生效 (yarn-default.xml)

配置参数	参数说明
yarn.scheduler.minimum-allocation-mb	给应用程序Container分配的最小内存，默认值: 1024
yarn.scheduler.maximum-allocation-mb	给应用程序Container分配的最大内存，默认值: 8192
yarn.scheduler.minimum-allocation-vcores	每个Container申请的最小CPU核数，默认值: 1
yarn.scheduler.maximum-allocation-vcores	每个Container申请的最大CPU核数，默认值: 32
yarn.nodemanager.resource.memory-mb	给Containers分配的最大物理内存，默认值: 8192

(c) Shuffle性能优化的关键参数，应在YARN启动之前就配置好 (mapred-default.xml)

配置参数	参数说明
mapreduce.task.io.sort.mb	Shuffle的环形缓冲区大小，默认100m
mapreduce.map.sort.spill.percent	环形缓冲区溢出的阈值，默认80%

## 2) 容错相关参数（MapReduce性能优化）

配置参数	参数说明
mapreduce.map.maxattempts	每个Map Task最大重试次数，一旦重试次数超过该值，则认为Map Task运行失败，默认值：4。
mapreduce.reduce.maxattempts	每个Reduce Task最大重试次数，一旦重试次数超过该值，则认为Map Task运行失败，默认值：4。
mapreduce.task.timeout	Task超时时间，经常需要设置的一个参数，该参数表达的意思为：如果一个Task在一定时间内没有任何进入，即不会读取新的数据，也没有输出数据，则认为该Task处于Block状态，可能是卡住了，也许永远会卡住，为了防止因为用户程序永远Block住不退出，则强制设置了一个该超时时间（单位毫秒），默认是600000（10分钟）。如果你的程序对每条输入数据的处理时间过长（比如会访问数据库，通过网络拉取数据等），建议将该参数调大，该参数过小常出现的错误提示是： “AttemptID:attempt_14267829456721_123456_m_000224_0 Timed out after 300 secsContainer killed by the ApplicationMaster.”。

## 2.4 Hadoop小文件优化方法

### 2.4.1 Hadoop小文件弊端

HDFS上每个文件都要在NameNode上创建对应的元数据，这个元数据的大小约为150byte，这样当小文件比较多的时候，就会产生很多的元数据文件，一方面会大量占用NameNode的内存空间，另一方面就是元数据文件过多，使得寻址索引速度变慢。

小文件过多，在进行MR计算时，会生成过多切片，需要启动过多的MapTask。每个MapTask处理的数据量小，导致MapTask的处理时间比启动时间还小，白白消耗资源。

### 2.4.2 Hadoop小文件解决方案

#### 1. 小文件优化的方向：

- （1）在数据采集的时候，就将小文件或大批数据合成大文件再上传HDFS。
- （2）在业务处理之前，在HDFS上使用MapReduce程序对小文件进行合并。
- （3）在MapReduce处理时，可采用CombineTextInputFormat提高效率。
- （4）开启uber模式，实现jvm重用

#### 2. Hadoop Archive

是一个高效的将小文件放入HDFS块中的文件存档工具，能够将多个小文件打包成一个HAR文件，从而达到减少NameNode的内存使用

#### 3. SequenceFile

SequenceFile是由一系列的二进制k/v组成，如果为key为文件名，value为文件内容，可将大批小文件合并成一个大文件

#### 4. CombineTextInputFormat

CombineTextInputFormat用于将多个小文件在切片过程中生成一个单独的切片或者少量的切片。

5. 开启uber模式，实现jvm重用。默认情况下，每个Task任务都需要启动一个jvm来运行，如果Task任务计算的数据量很小，我们可以让同一个Job的多个Task运行在一个jvm中，不必为每个Task都开启一个jvm。

开启uber模式，在mapred-site.xml中添加如下配置

```
<!-- 开启uber模式 -->
<property>
  <name>mapreduce.job.ubertask.enable</name>
  <value>true</value>
</property>

<!-- uber模式中最大的mapTask数量，可向下修改 -->
<property>
  <name>mapreduce.job.ubertask.maxmaps</name>
  <value>9</value>
</property>

<!-- uber模式中最大的reduce数量，可向下修改 -->
<property>
  <name>mapreduce.job.ubertask.maxreduces</name>
  <value>1</value>
</property>

<!-- uber模式中最大的输入数据量，默认使用dfs.blocksize 的值，可向下修改 -->
<property>
  <name>mapreduce.job.ubertask.maxbytes</name>
  <value></value>
</property>
```

## 3 Hadoop新特性

### 3.1 Hadoop2.x新特性

#### 3.1.1 集群间数据拷贝

##### 1) scp实现两个远程主机之间的文件复制

```
scp -r hello.txt root@hadoop103:/user/xu1an/hello.txt // 推 push
scp -r root@hadoop103:/user/xu1an/hello.txt hello.txt // 拉 pull
scp -r root@hadoop103:/user/xu1an/hello.txt root@hadoop104:/user/xu1an //是通过
本地主机中转实现两个主机的文件复制；如果在两个远程主机之间ssh没有配置的情况下可以使用该方式。
```

##### 2) 采用distcp命令实现两个Hadoop集群之间的递归数据复制

```
bin/hadoop distcp hdfs://hadoop102:9820/user/xu1an/hello.txt
hdfs://hadoop105:9820/user/xu1an/hello.txt
```

#### 3.1.2 小文件存档

##### 1、HDFS存储小文件弊端

每个文件均按块存储，每个块的元数据存储(NameNode的内存中)，因此HDFS存储小文件会非常低效。因为大量的小文件会耗尽NameNode中的大部分内存。但注意，存储小文件所需要的磁盘容量和数据块的大小无关。例如，一个1MB的文件设置为128MB的块存储，实际使用的是1MB的磁盘空间，而不是128MB。

## 2、解决存储小文件办法之一

HDFS存档文件或HAR文件，是一个更高效的文件存档工具，它将文件存入HDFS块，在减少NameNode内存使用的同时，允许对文件进行透明的访问。具体说来，HDFS存档文件对内还是一个一个独立文件，对NameNode而言却是一个整体，减少了NameNode的内存。



### 1) 案例实操

(1) 需要启动YARN进程

```
start-yarn.sh
```

(2) 归档文件

把/user/xu1an/input目录里面的所有文件归档成一个叫input.har的归档文件，并把归档后文件存储到/user/xu1an/output路径下。

```
hadoop archive -archiveName input.har -p /user/xu1an/input /user/xu1an/output
```

(3) 查看归档

```
hadoop fs -ls /user/xu1an/output/input.har
hadoop fs -ls har:///user/xu1an/output/input.har
```

(4) 解归档文件

```
hadoop fs -cp har:/// user/xu1an/output/input.har/* /user/xu1an
```

## 3.1.3 回收站

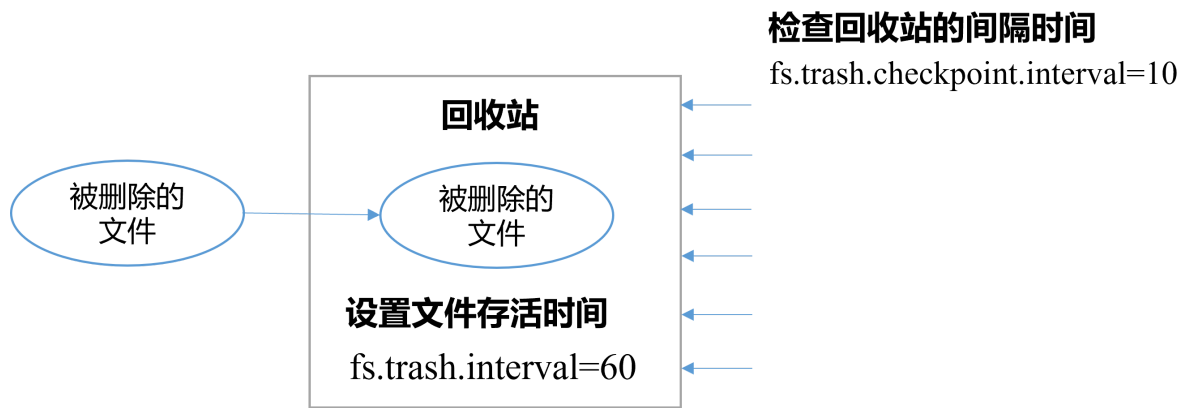
开启回收站功能，可以将删除的文件在不超时的情况下，恢复原数据，起到防止误删除、备份等作用。

1) 回收站参数设置及工作机制

一、开启回收站功能参数说明：

- 1、默认值fs.trash.interval=0，0表示禁用回收站;其他值表示设置文件的存活时间。
- 2、默认值fs.trash.checkpoint.interval=0，检查回收站的间隔时间。如果该值为0，则该值设置和fs.trash.interval的参数值相等。
- 3、要求fs.trash.checkpoint.interval<=fs.trash.interval。

二、回收站工作机制：



## 2) 启用回收站

修改core-site.xml，配置垃圾回收时间为1分钟。

```
<property>
  <name>fs.trash.interval</name>
  <value>1</value>
</property>
<property>
  <name>fs.trash.checkpoint.interval</name>
  <value>1</value>
</property>
```

## 3) 查看回收站

回收站目录在hdfs集群中的路径：/user/xu1an/.Trash/....

4) 通过程序删除的文件不会经过回收站，需要调用moveToTrash()才进入回收站

```
Trash trash = New Trash(conf);
trash.moveToTrash(path);
```

5) 通过网页上直接删除的文件也不会走回收站。

6) 只有在命令行利用hadoop fs -rm命令删除的文件才会走回收站。

```
[xu1an@hadoop102 hadoop-3.1.3]$ hadoop fs -rm -r /user/xu1an/input
2021-07-14 16:13:42,643 INFO fs.TrashPolicyDefault: Moved:
'hd fs://hadoop102:9820/user/xu1an/input' to trash at:
hd fs://hadoop102:9820/user/xu1an/.Trash/Current/user/xu1an/input
```

## 7) 恢复回收站数据

```
[xu1an@hadoop102 hadoop-3.1.3]$ hadoop fs -mv
/user/xu1an/.Trash/Current/user/xu1an/input /user/xu1an/input
```

# 3.2 Hadoop3.x新特性

### 3.2.1 多NN的HA架构

HDFS NameNode高可用性的初始实现为单个活动NameNode和单个备用NameNode，将edits复制到三个JournalNode。该体系结构能够容忍系统中一个NN或一个JN的故障。

但是，某些部署需要更高层次的容错能力。Hadoop3.x允许用户运行多个备用NameNode。例如，通过配置三个NameNode和五个JournalNode，群集能够容忍两个节点而不是一个节点的故障。

### 3.2.2 纠删码

HDFS中的默认3副本方案在存储空间和其他资源（例如，网络带宽）中具有200%的开销。但是，对于I/O活动相对较低暖和冷数据集，在正常操作期间很少访问其他块副本，但仍会消耗与第一个副本相同的资源量。

纠删码（Erasure Coding）能够在不到50%的数据冗余情况下提供和3副本相同的容错能力，因此，使用纠删码作为副本机制的改进是自然而然的。

查看集群支持的纠删码策略：`hdfs ec -listPolicies`

---

## 4 Hadoop HA高可用

### 4.1 HA概述

(1) 所谓HA（High Availability），即高可用（7\*24小时不中断服务）。

(2) 实现高可用最关键的策略是消除单点故障。HA严格来说应该分成各个组件的HA机制：HDFS的HA和YARN的HA。

(3) Hadoop2.0之前，在HDFS集群中NameNode存在单点故障（SPOF）。

(4) NameNode主要在以下两个方面影响HDFS集群

- NameNode机器发生意外，如宕机，集群将无法使用，直到管理员重启
- NameNode机器需要升级，包括软件、硬件升级，此时集群也将无法使用

HDFS HA功能通过配置Active/Standby两个NameNodes实现在集群中对NameNode的热备来解决上述问题。如果出现故障，如机器崩溃或机器需要升级维护，这时可通过此种方式将NameNode很快的切换到另外一台机器。

### 4.2 HDFS-HA工作机制

通过多个NameNode消除单点故障

#### 4.2.1 HDFS-HA工作要点

##### 1) 元数据管理方式需要改变

内存中各自保存一份元数据；

Edits日志只有Active状态的NameNode节点可以做写操作；

所有的NameNode都可以读取Edits；

共享的Edits放在一个共享存储中管理（qjournal和NFS两个主流实现）；

##### 2) 需要一个状态管理功能模块

实现了一个zkfailover，常驻在每一个namenode所在的节点，每一个zkfailover负责监控自己所在NameNode节点，利用zk进行状态标识，当需要进行状态切换时，由zkfailover来负责切换，切换时需要防止brain split现象的发生。

3) 必须保证两个NameNode之间能够ssh无密码登录

4) 隔离 (Fence) , 即同一时刻仅仅有一个NameNode对外提供服务

## 4.2.2 HDFS-HA自动故障转移工作机制

自动故障转移为HDFS部署增加了两个新组件: ZooKeeper和ZKFailoverController (ZKFC) 进程, 如图3-20所示。ZooKeeper是维护少量协调数据, 通知客户端这些数据的改变和监视客户端故障的高可用服务。HA的自动故障转移依赖于ZooKeeper的以下功能:

### 1. 故障检测

集群中的每个NameNode在ZooKeeper中维护了一个会话, 如果机器崩溃, ZooKeeper中的会话将终止, ZooKeeper通知另一个NameNode需要触发故障转移。

### 2. 现役NameNode选择

ZooKeeper提供了一个简单的机制用于唯一的选择一个节点为active状态。如果目前现役NameNode崩溃, 另一个节点可能从ZooKeeper获得特殊的排外锁以表明它应该成为现役NameNode。

ZKFC是自动故障转移中的另一个新组件, 是ZooKeeper的客户端, 也监视和管理NameNode的状态。每个运行NameNode的主机也运行了一个ZKFC进程, ZKFC负责:

#### 1) 健康监测

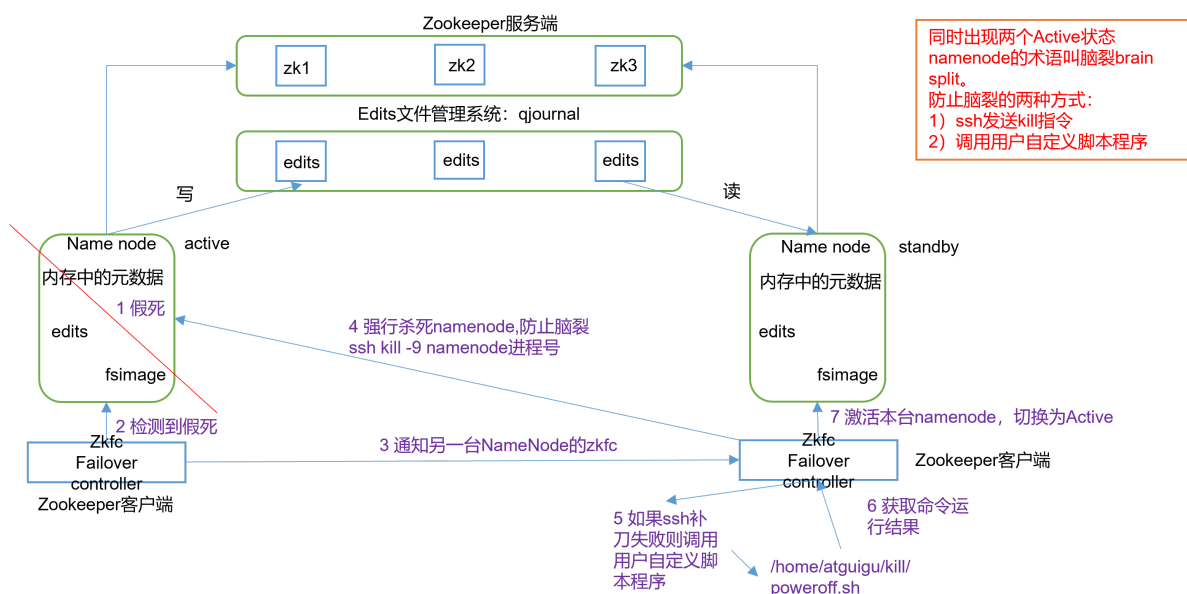
ZKFC使用一个健康检查命令定期地ping与之在相同主机的NameNode, 只要该NameNode及时地回复健康状态, ZKFC认为该节点是健康的。如果该节点崩溃, 冻结或进入不健康状态, 健康监测器标识该节点为非健康的。

#### 2) ZooKeeper会话管理

当本地NameNode是健康的, ZKFC保持一个在ZooKeeper中打开的会话。如果本地NameNode处于active状态, ZKFC也保持一个特殊的znode锁, 该锁使用了ZooKeeper对短暂节点的支持, 如果会话终止, 锁节点将自动删除。

#### 3) 基于ZooKeeper的选择

如果本地NameNode是健康的, 且ZKFC发现没有其它的节点当前持有znode锁, 它将为自己获取该锁。如果成功, 则它已经赢得了选择, 并负责运行故障转移进程以使它的本地NameNode为Active。



## 4.3 HDFS-HA集群配置

### 4.3.1 环境准备

- (1) 修改IP
- (2) 修改主机名及主机名和IP地址的映射
- (3) 关闭防火墙
- (4) ssh免密登录
- (5) 安装JDK，配置环境变量等

### 4.3.2 规划集群

hadoop102	hadoop103	hadoop104
NameNode	NameNode	NameNode
ZKFC	ZKFC	ZKFC
JournalNode	JournalNode	JournalNode
DataNode	DataNode	DataNode
ZK	ZK	ZK
	ResourceManager	
NodeManager	NodeManager	NodeManager

### 4.3.3 配置Zookeeper集群

#### 1) 集群规划

在hadoop102、hadoop103和hadoop104三个节点上部署Zookeeper。

#### 2) 解压安装

- (1) 解压Zookeeper安装包到/opt/module/目录下

```
tar -zxvf zookeeper-3.5.7.tar.gz -C /opt/module/
```

- (2) 在/opt/module/zookeeper-3.5.7/这个目录下创建zkData

```
mkdir -p zkData
```

- (3) 重命名/opt/module/zookeeper-3.4.14/conf这个目录下的zoo\_sample.cfg为zoo.cfg

```
mv zoo_sample.cfg zoo.cfg
```

#### 3) 配置zoo.cfg文件

- (1) 具体配置

```
dataDir=/opt/module/zookeeper-3.5.7/zkData
```



增加如下配置

```
server.2=hadoop102:2888:3888  
server.3=hadoop103:2888:3888  
server.4=hadoop104:2888:3888
```

## (2) 配置参数解读

Server.A=B:C:D。

A是一个数字，表示这个是第几号服务器；

B是这个服务器的IP地址；

C是这个服务器与集群中的Leader服务器交换信息的端口；

D是万一集群中的Leader服务器挂了，需要一个端口来重新进行选举，选出一个新的Leader，而这个端口就是用来执行选举时服务器相互通信的端口。

集群模式下配置一个文件myid，这个文件在dataDir目录下，这个文件里面有一个数据就是A的值，Zookeeper启动时读取此文件，拿到里面的数据与zoo.cfg里面的配置信息比较从而判断到底是哪个server。

## 4) 集群操作

(1) 在/opt/module/zookeeper-3.5.7/zkData目录下创建一个myid的文件

```
[xu1an@hadoop102 zkData]$ touch myid
```

添加myid文件，注意一定要在linux里面创建，在notepad++里面很可能乱码

(2) 编辑myid文件

```
[xu1an@hadoop102 zkData]$ vi myid
```

在文件中添加与server对应的编号：如2

(3) 拷贝配置好的zookeeper到其他机器上

```
[xu1an@hadoop102 module]$ scp -r zookeeper-3.5.7/ xu1an@hadoop103:/opt/module/  
[xu1an@hadoop102 module]$ scp -r zookeeper-3.5.7/ xu1an@hadoop104:/opt/module/
```

并分别修改myid文件中内容为3、4

(4) 分别启动zookeeper

```
[xu1an@hadoop102 zookeeper-3.5.7]$ bin/zkServer.sh start  
[xu1an@hadoop103 zookeeper-3.5.7]$ bin/zkServer.sh start  
[xu1an@hadoop104 zookeeper-3.5.7]$ bin/zkServer.sh start
```

(5) 查看状态

```
[xu1an@hadoop104 zookeeper-3.5.7]$ bin/zkServer.sh status
JMX enabled by default
Using config: /opt/module/zookeeper-3.5.7/bin/../conf/zoo.cfg
Mode: follower
[xu1an@hadoop104 zookeeper-3.5.7]$ bin/zkServer.sh status
JMX enabled by default
Using config: /opt/module/zookeeper-3.5.7/bin/../conf/zoo.cfg
Mode: leader
[xu1an@hadoop104 zookeeper-3.5.7]$ bin/zkServer.sh status
JMX enabled by default
Using config: /opt/module/zookeeper-3.5.7/bin/../conf/zoo.cfg
Mode: follower
```

#### 4.3.4 配置HDFS-HA集群

1) 官方地址: <http://hadoop.apache.org/>

2) 在opt目录下创建一个ha文件夹

```
[xu1an@hadoop102 ~]$ cd /opt
[xu1an@hadoop102 opt]$ sudo mkdir ha
[xu1an@hadoop102 opt]$ sudo chown xu1an:xu1an /opt/ha
```

3) 将/opt/module/下的 hadoop-3.1.3拷贝到/opt/ha目录下 (记得删除data 和 log目录)

```
[xu1an@hadoop102 opt]$ cp -r /opt/module/hadoop-3.1.3 /opt/ha/
```

4) 配置hadoop-env.sh

```
export JAVA_HOME=/opt/module/jdk1.8.0_212
```

5) 配置core-site.xml

```
<configuration>
<!-- 把多个NameNode的地址组装成一个集群mycluster -->
  <property>
    <name>fs.defaultFS</name>
    <value>hdfs://mycluster</value>
  </property>
<!-- 指定hadoop运行时产生文件的存储目录 -->
  <property>
    <name>hadoop.tmp.dir</name>
    <value>/opt/ha/hadoop-3.1.3/data</value>
  </property>
</configuration>
```

6) 配置hdfs-site.xml

```
<configuration>
<!-- NameNode数据存储目录 -->
  <property>
    <name>dfs.namenode.name.dir</name>
    <value>file://${hadoop.tmp.dir}/name</value>
  </property>
```

```
<!-- DataNode数据存储目录 -->
<property>
  <name>dfs.datanode.data.dir</name>
  <value>file://${hadoop.tmp.dir}/data</value>
</property>
<!-- JournalNode数据存储目录 -->
<property>
  <name>dfs.journalnode.edits.dir</name>
  <value>${hadoop.tmp.dir}/jn</value>
</property>
<!-- 完全分布式集群名称 -->
<property>
  <name>dfs.nameservices</name>
  <value>mycluster</value>
</property>
<!-- 集群中NameNode节点都有哪些 -->
<property>
  <name>dfs.ha.namenodes.mycluster</name>
  <value>nn1,nn2,nn3</value>
</property>
<!-- NameNode的RPC通信地址 -->
<property>
  <name>dfs.namenode.rpc-address.mycluster.nn1</name>
  <value>hadoop102:8020</value>
</property>
<property>
  <name>dfs.namenode.rpc-address.mycluster.nn2</name>
  <value>hadoop103:8020</value>
</property>
<property>
  <name>dfs.namenode.rpc-address.mycluster.nn3</name>
  <value>hadoop104:8020</value>
</property>
<!-- NameNode的http通信地址 -->
<property>
  <name>dfs.namenode.http-address.mycluster.nn1</name>
  <value>hadoop102:9870</value>
</property>
<property>
  <name>dfs.namenode.http-address.mycluster.nn2</name>
  <value>hadoop103:9870</value>
</property>
<property>
  <name>dfs.namenode.http-address.mycluster.nn3</name>
  <value>hadoop104:9870</value>
</property>
<!-- 指定NameNode元数据在JournalNode上的存放位置 -->
<property>
  <name>dfs.namenode.shared.edits.dir</name>
  <value>qjournal://hadoop102:8485;hadoop103:8485;hadoop104:8485/mycluster</value>
</property>
<!-- 访问代理类: client用于确定哪个NameNode为Active -->
<property>
  <name>dfs.client.failover.proxy.provider.mycluster</name>

  <value>org.apache.hadoop.hdfs.server.namenode.ha.ConfiguredFailoverProxyProvide
r</value>
</property>
```

```
<!-- 配置隔离机制，即同一时刻只能有一台服务器对外响应 -->
<property>
  <name>dfs.ha.fencing.methods</name>
  <value>sshfence</value>
</property>
<!-- 使用隔离机制时需要ssh密钥登录-->
<property>
  <name>dfs.ha.fencing.ssh.private-key-files</name>
  <value>/home/xu1an/.ssh/id_rsa</value>
</property>
</configuration>
```

7) 分发配置好的hadoop环境到其他节点

### 4.3.5 启动HDFS-HA集群

1) 将HADOOP\_HOME环境变量更改到HA目录

```
[xu1an@hadoop102 ~]$ sudo vim /etc/profile.d/my_env.sh
```

将HADOOP\_HOME部分改为如下

```
##HADOOP_HOME
export HADOOP_HOME=/opt/ha/hadoop-3.1.3
export PATH=$PATH:$HADOOP_HOME/bin
export PATH=$PATH:$HADOOP_HOME/sbin
```

2) 在各个JournalNode节点上，输入以下命令启动journalnode服务

```
[xu1an@hadoop102 ~]$ hdfs --daemon start journalnode
[xu1an@hadoop103 ~]$ hdfs --daemon start journalnode
[xu1an@hadoop104 ~]$ hdfs --daemon start journalnode
```

3) 在[nn1]上，对其进行格式化，并启动

```
[xu1an@hadoop102 ~]$ hdfs namenode -format
[xu1an@hadoop102 ~]$ hdfs --daemon start namenode
```

4) 在[nn2]和[nn3]上，同步nn1的元数据信息

```
[atguigu@hadoop103 ~]$ hdfs namenode -bootstrapStandby
[atguigu@hadoop104 ~]$ hdfs namenode -bootstrapStandby
```

5) 启动[nn2]和[nn3]

```
[atguigu@hadoop103 ~]$ hdfs --daemon start namenode
[atguigu@hadoop104 ~]$ hdfs --daemon start namenode
```

6) 查看web页面显示

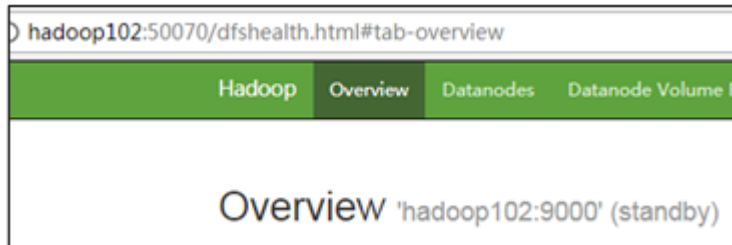


图 hadoop102(standby)

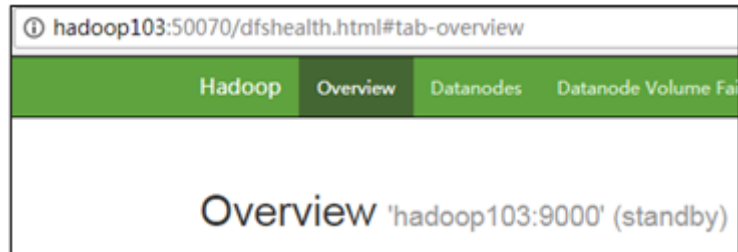


图 hadoop103(standby)

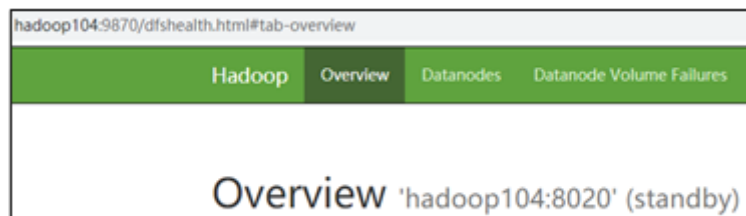


图 hadoop104(standby)

7) 在所有节点上, 启动datanode

```
[atguigu@hadoop102 ~]$ hdfs --daemon start datanode
[atguigu@hadoop103 ~]$ hdfs --daemon start datanode
[atguigu@hadoop104 ~]$ hdfs --daemon start datanode
```

8) 将[nn1]切换为Active

```
[atguigu@hadoop102 ~]$ hdfs haadmin -transitionToActive nn1
```

9) 查看是否Active

```
[atguigu@hadoop102 ~]$ hdfs haadmin -getServiceState nn1
```

### 4.3.6 配置HDFS-HA自动故障转移

#### 1) 具体配置

(1) 在hdfs-site.xml中增加

```
<!-- 启用nn故障自动转移 -->
<property>
  <name>dfs.ha.automatic-failover.enabled</name>
  <value>true</value>
</property>
```

(2) 在core-site.xml文件中增加

```
<!-- 指定zkfc要连接的zkServer地址 -->
<property>
  <name>ha.zookeeper.quorum</name>
  <value>hadoop102:2181,hadoop103:2181,hadoop104:2181</value>
</property>
```

### (3) 修改后分发配置文件

```
[atguigu@hadoop102 etc]$ pwd
/opt/ha/hadoop-3.1.3/etc
[atguigu@hadoop102 etc]$ xsync hadoop/
```

## 2) 启动

### (1) 关闭所有HDFS服务:

```
[atguigu@hadoop102 ~]$ stop-dfs.sh
```

### (2) 启动Zookeeper集群:

```
[atguigu@hadoop102 ~]$ zkServer.sh start
[atguigu@hadoop103 ~]$ zkServer.sh start
[atguigu@hadoop104 ~]$ zkServer.sh start
```

### (3) 启动Zookeeper以后, 然后再初始化HA在Zookeeper中状态:

```
[atguigu@hadoop102 ~]$ hdfs zkfc -formatzk
```

### (4) 启动HDFS服务:

```
[atguigu@hadoop102 ~]$ start-dfs.sh
```

### (5) 可以去zkCli.sh客户端查看Namenode选举锁节点内容:

```
[zk: localhost:2181(CONNECTED) 7] get -s /hadoop-
ha/mycluster/ActiveStandbyElectorLock

    myclusternn2      hadoop103 <>(<>
cZxid = 0x10000000b
ctime = Tue Jul 14 17:00:13 CST 2020
mZxid = 0x10000000b
mtime = Tue Jul 14 17:00:13 CST 2020
pZxid = 0x10000000b
cversion = 0
dataVersion = 0
aclVersion = 0
ephemeralOwner = 0x40000da2eb70000
dataLength = 33
numChildren = 0
```

## 3) 验证

### (1) 将Active NameNode进程kill, 查看网页端三台Namenode的状态变化

```
[atguigu@hadoop102 ~]$ kill -9 namenode的进程id
```

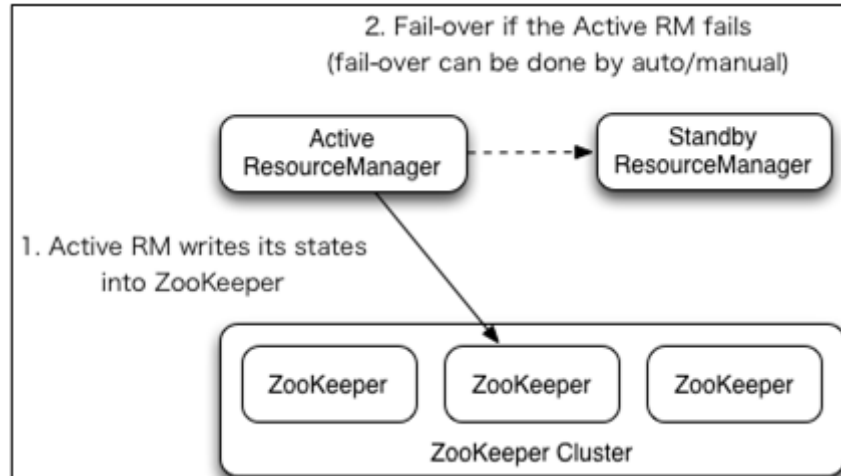
## 4.4 YARN-HA配置

### 4.4.1 YARN-HA工作机制

#### 1) 官方文档:

<http://hadoop.apache.org/docs/r3.1.3/hadoop-yarn/hadoop-yarn-site/ResourceManagerHA.html>

#### 2) YARN-HA工作机制



### 4.4.2 配置YARN-HA集群

#### 1) 环境准备

- (1) 修改IP
- (2) 修改主机名及主机名和IP地址的映射
- (3) 关闭防火墙
- (4) ssh免密登录
- (5) 安装JDK, 配置环境变量等
- (6) 配置Zookeeper集群

#### 2) 规划集群

hadoop102	hadoop103	hadoop104
NameNode	NameNode	NameNode
JournalNode	JournalNode	JournalNode
DataNode	DataNode	DataNode
ZK	ZK	ZK
ResourceManager	ResourceManager	
NodeManager	NodeManager	NodeManager

#### 3) 具体配置

(1) yarn-site.xml

```
<configuration>

  <property>
    <name>yarn.nodemanager.aux-services</name>
    <value>mapreduce_shuffle</value>
  </property>

  <!-- 启用resourcemanager ha -->
  <property>
    <name>yarn.resourcemanager.ha.enabled</name>
    <value>true</value>
  </property>

  <!-- 声明两台resourcemanager的地址 -->
  <property>
    <name>yarn.resourcemanager.cluster-id</name>
    <value>cluster-yarn1</value>
  </property>
  <!--指定resourcemanager的逻辑列表-->
  <property>
    <name>yarn.resourcemanager.ha.rm-ids</name>
    <value>rm1,rm2</value>
  </property>
  <!-- ===== rm1的配置 ===== -->
  <!-- 指定rm1的主机名 -->
  <property>
    <name>yarn.resourcemanager.hostname.rm1</name>
    <value>hadoop102</value>
  </property>
  <!-- 指定rm1的web端地址 -->
  <property>
    <name>yarn.resourcemanager.webapp.address.rm1</name>
    <value>hadoop102:8088</value>
  </property>
  <!-- 指定rm1的内部通信地址 -->
  <property>
    <name>yarn.resourcemanager.address.rm1</name>
    <value>hadoop102:8032</value>
  </property>
  <!-- 指定AM向rm1申请资源的地址 -->
  <property>
    <name>yarn.resourcemanager.scheduler.address.rm1</name>
    <value>hadoop102:8030</value>
  </property>
  <!-- 指定供NM连接的地址 -->
  <property>
    <name>yarn.resourcemanager.resource-tracker.address.rm1</name>
    <value>hadoop102:8031</value>
  </property>
  <!-- ===== rm2的配置 ===== -->
  <!-- 指定rm2的主机名 -->
  <property>
    <name>yarn.resourcemanager.hostname.rm2</name>
    <value>hadoop103</value>
  </property>
```



```

<property>
  <name>yarn.resourcemanager.webapp.address.rm2</name>
  <value>hadoop103:8088</value>
</property>
<property>
  <name>yarn.resourcemanager.address.rm2</name>
  <value>hadoop103:8032</value>
</property>
<property>
  <name>yarn.resourcemanager.scheduler.address.rm2</name>
  <value>hadoop103:8030</value>
</property>
<property>
  <name>yarn.resourcemanager.resource-tracker.address.rm2</name>
  <value>hadoop103:8031</value>
</property>

  <!-- 指定zookeeper集群的地址 -->
  <property>
    <name>yarn.resourcemanager.zk-address</name>
    <value>hadoop102:2181,hadoop103:2181,hadoop104:2181</value>
  </property>

  <!-- 启用自动恢复 -->
  <property>
    <name>yarn.resourcemanager.recovery.enabled</name>
    <value>true</value>
  </property>

  <!-- 指定resourcemanager的状态信息存储在zookeeper集群 -->
  <property>
    <name>yarn.resourcemanager.store.class</name>
    <value>org.apache.hadoop.yarn.server.resourcemanager.recovery.ZKRMStateStore</value>
  </property>
  <!-- 环境变量的继承 -->
  <property>
    <name>yarn.nodemanager.env-whitelist</name>

    <value>JAVA_HOME,HADOOP_COMMON_HOME,HADOOP_HDFS_HOME,HADOOP_CONF_DIR,CLASSPATH_
    PREPEND_DISTCACHE,HADOOP_YARN_HOME,HADOOP_MAPRED_HOME</value>
  </property>
</configuration>

```

(2) 同步更新其他节点的配置信息，分发配置文件

```
[atguigu@hadoop102 etc]$ xsync hadoop/
```

#### 4) 启动hdfs

```
[atguigu@hadoop102 ~]$ start-dfs.sh
```

#### 5) 启动YARN

(1) 在hadoop102或者hadoop103中执行：

```
[atguigu@hadoop102 ~]$ start-yarn.sh
```

(2) 查看服务状态

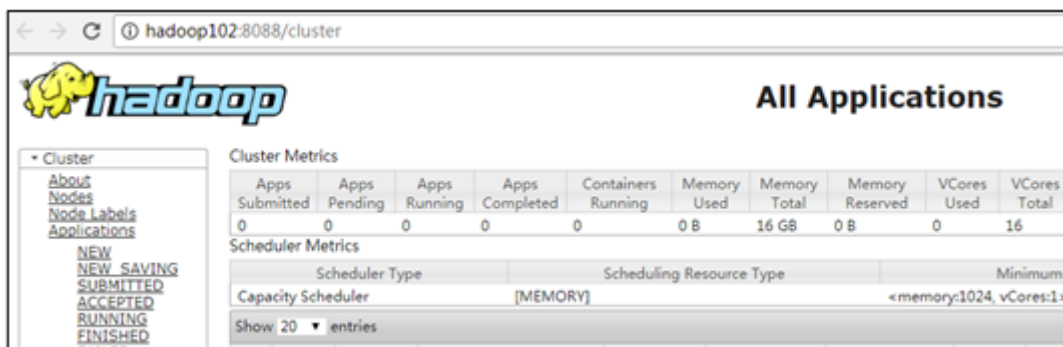
```
[atguigu@hadoop102 ~]$ yarn rmadmin -getServiceState rm1
```

(3) 可以去zkCli.sh客户端查看ResourceManager选举锁节点内容:

```
[atguigu@hadoop102 ~]$ zkCli.sh
[zk: localhost:2181(CONNECTED) 16] get -s /yarn-leader-election/cluster-
yarn1/ActiveStandbyElectorLock

cluster-yarn1rm1
czxid = 0x100000022
ctime = Tue Jul 14 17:06:44 CST 2020
mzxid = 0x100000022
mtime = Tue Jul 14 17:06:44 CST 2020
pzxid = 0x100000022
cversion = 0
dataVersion = 0
aclVersion = 0
ephemeralOwner = 0x30000da33080005
dataLength = 20
numChildren = 0
```

(4) web端查看hadoop102:8088和hadoop103:8088的YARN的状态, 和NameNode对比, 查看区别



Cluster Metrics										
Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Memory Used	Memory Total	Memory Reserved	VCores Used	VCores Total	
0	0	0	0	0	0 B	16 GB	0 B	0	16	

Scheduler Metrics		
Scheduler Type	Scheduling Resource Type	Minimum
Capacity Scheduler	[MEMORY]	<memory:1024, vCores:1>

## 4.5 HDFS Federation架构设计

### 4.5.1 NameNode架构的局限性

#### 1) Namespace (命名空间) 的限制

由于NameNode在内存中存储所有的元数据 (metadata), 因此单个NameNode所能存储的对象 (文件+块) 数目受到NameNode所在JVM的heap size的限制。50G的heap能够存储20亿 (200million) 个对象, 这20亿个对象支持4000个DataNode, 12PB的存储 (假设文件平均大小为40MB)。随着数据的飞速增长, 存储的需求也随之增长。单个DataNode从4T增长到36T, 集群的尺寸增长到8000个DataNode。存储的需求从12PB增长到大于100PB。

#### 2) 隔离问题

由于HDFS仅有一个NameNode, 无法隔离各个程序, 因此HDFS上的一个实验程序就很有可能影响整个HDFS上运行的程序。

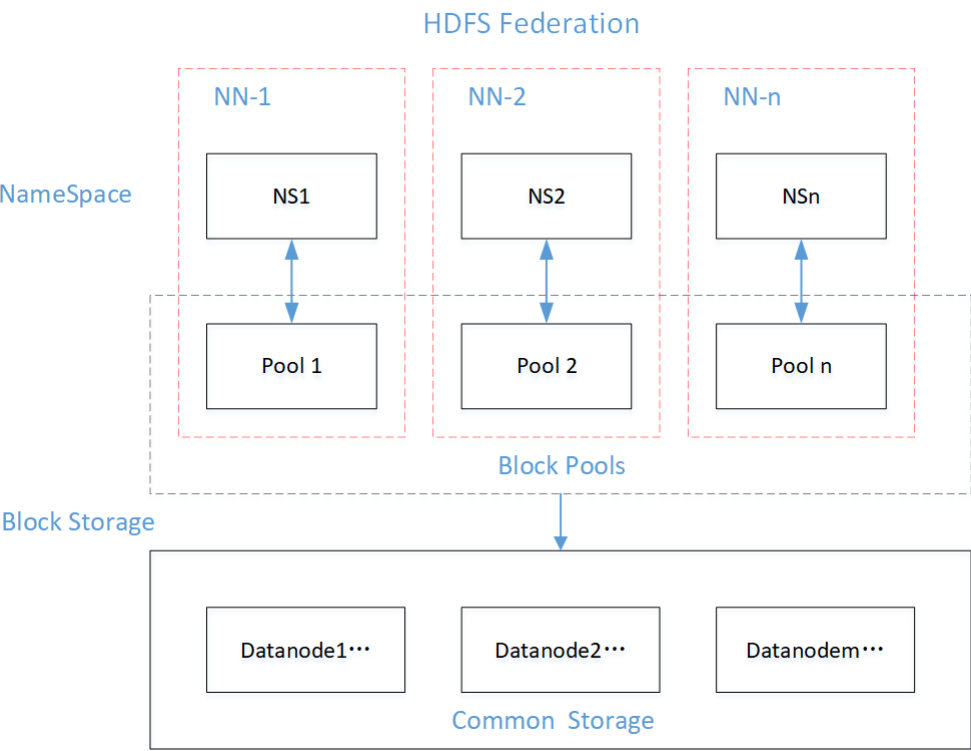
#### 3) 性能的瓶颈

由于是单个NameNode的HDFS架构，因此整个HDFS文件系统的吞吐量受限于单个NameNode的吞吐量

4.5.2 HDFS Federation架构设计

能不能有多个NameNode

NameNode	NameNode	NameNode
元数据	元数据	元数据
Log	machine	电商数据/话单数据



4.5.3 HDFS Federation应用思考

不同应用可以使用不同NameNode进行数据管理图片业务、爬虫业务、日志审计业务。Hadoop生态系统中，不同的框架使用不同的NameNode进行管理NameSpace。（隔离性）

