



1

静静



一种美德

写评论



目录



收藏



微信



微博



QQ

RSS订阅

## HDFS HA切换后missing block问题分析

2015年05月31日 21:03:31

我们Hadoop平台也从Hadoop1.2.1升级到了Hadoop2.4.0版本，当然HDFS HA 也配置到集群中。具体的配置方案（<https://issues.apache.org/jira/browse/HDFS-1623>）。感恩cloudera 这样伟大的公司，为我们提供了一个从Hadoop1.2.1迁移到Hadoop2.4.0，同时验证HDFS HA 切换是否对HBase功能产生影响。问题出现了，当我们持续写入数据（namenode），读出刚才写入的数据，悲剧了，数据读不来，HDFS web页面出现了missing block。

接下来我们分析HBase 相关日志，发现在HA 切换时写入 hfile 文件出现问题。ok，我们现在先验证HDFS HA 切换不难啊，我们用HDFS client代码持续写入大量小文件（10个字节），同时进行HA切换。悲剧再次发生了，H

There are 14 missing blocks. The following files may be corrupted:

```
blk_1075342606 /robin/our/jm.txt.48
blk_1075342607 /robin/our/jm.txt.49
blk_1075342608 /robin/our/jm.txt.50
blk_1075342609 /robin/our/jm.txt.51
blk_1075342610 /robin/our/jm.txt.52
blk_1075342611 /robin/our/jm.txt.53
blk_1075342612 /robin/our/jm.txt.54
blk_1075342613 /robin/our/jm.txt.55
blk_1075342614 /robin/our/jm.txt.56
blk_1075342615 /robin/our/jm.txt.57
blk_1075342616 /robin/our/jm.txt.58
blk_1075342617 /robin/our/jm.txt.59
blk_1075342618 /robin/our/jm.txt.60
blk_1075342619 /robin/our/jm.txt.61
```

Please check the logs or run fsck in order to identify the missing blocks. See the Hadoop FAQ for common causes and potential solutions.

现在问题复现了，先兴奋一把，我们先解决 HDFS HA 切换造成块丢失的问题。解决了这个问题再继续验证HBase

- namenode2是active Namenode，我们切换的时候直接kill掉该节点；
- namenode1是standby Namenode，我们切换后直接变为active Namenode；

### 一、问题分析

#### 1) Namenode 分析

我们选择第一个block:blk\_1075342606进行分析，切换后namenode1日志如下：

```
2015-04-19 10:53:45,789 INFO BlockStateChange: BLOCK* addBlock: block blk_1075342606 1602142 on 10.10.10.10 size 143 does not belong to any file
2015-04-19 10:53:45,789 INFO BlockStateChange: BLOCK* InvalideBlocks: add blk_1075342606 1602142 to 10.10.10.10 5001
```

「码字计划」：拿万元写作基金！

免费云主机试用一年

cloudera课程

题库下载

登录

注册

X

1 该日志说明已经给相应的block分配了相关的Datanode节点。但namenode的inode ( Fsimage ) 和 blockmap 忘记了。忘记的同学可以一下看看ggjucheng 同学对namenode源码简要分析: 2013/02/04/2889386.html。

写评论 需要记住的是：NameNode作为HDFS中文件目录和文件分配的管理者，它保存的最重要信息，就是下面所

- 文件名=>数据块：持久化为fsimage，持久化到内存中、文件系统的inode相似；
- 数据块=>DataNode列表，BLockmap保存的结构。

## 收藏 Datanode日志

量查询Datanode日志，查看blk是否写入到Datanode。方法如下

微信 `for i in $(cat ~/software/hadoop/etc/hadoop/slaves); do echo $i; ssh $i "grep blk_1075342606 ~/hadoop/logs/hadoop-bigdata-hadoop-hdfs.server.datanode.DataNode: PacketResponder: BP-2057938263-1422984781848:blk_1075342606_1602142, type=HAS_DOWNSTREAM_IN_PIPELINE terminating`

微博

QQ

```
/home/bigdata/hadoop/logs/hadoop-bigdata-datanode-1-1422984781848.log.1:2015-04-19 10:53:12,595 INFO org.apache.hadoop.hdfs.server.datanode.DataNode: PacketResponder: BP-2057938263-1422984781848:blk_1075342606_1602142, type=HAS_DOWNSTREAM_IN_PIPELINE terminating
/home/bigdata/hadoop/logs/hadoop-bigdata-datanode-2-1422984781848.log.2:2015-04-19 10:53:12,595 INFO org.apache.hadoop.hdfs.server.datanode.DataNode: PacketResponder: BP-2057938263-1422984781848:blk_1075342606_1602142, type=LAST_IN_PIPELINE, downstreams=0:[] terminating
/home/bigdata/hadoop/logs/hadoop-bigdata-datanode-3-1422984781848.log.3:2015-04-19 10:53:12,595 INFO org.apache.hadoop.hdfs.server.datanode.DataNode: PacketResponder: BP-2057938263-1422984781848:blk_1075342606_1602142, type=HAS_DOWNSTREAM_IN_PIPELINE terminating
```

这说明blk 已经成功写入到了三个Datanode节点。数据写入正常，但是namenode 1里面关于该block的元数据不

### 3) FileJournalManager同步日志

```
2015-04-19 10:53:42,415 INFO org.apache.hadoop.hdfs.server.namenode.FSImage: Start loading edits file
```

需要查看edits是否记录该block的edit log文件。在10:53:42开始大量的日志editlog同步。需要确定blk\_1075342606\_1602142

### 4) 综合分析

我们可以做如下分析：对于namenode1日志,blockmap 里面为什么没有该block信息？熟悉HDFS的同学知道，一个文件对应的blockid保存入BlocksMap，此时BlocksMap中每个block对应的Datanodes列表暂时为空。创建文件时在大神（微博名：@王志强-Austin）指导下，结合 Namenode HA方案，standby namenode同步editlog 通过Cock给 namenode1的时候，edit log还没同步过来。怎么验证这种猜想了？重启一下这个block对于的Datanode，我们重启该block report对于的datanode，oh my God！奇迹发生了，该block从页面的missing block队列里面去数据。

我们回头过来分析namenode1 block report 和edit log 的时间，发现在HDFS HA 切换期间，当Datanode 的incr时，后台standby没有完全变成activenamenode之前，会出现包含 invalidate block 的后台日志。

## 1) Datanode blockreport

Datanode会进行增量或者全量block report 给namenode。增量的block report间隔时间为100\* ( heartbeat间隔时间) 一个小时。下面是代码入口：

reportReceivedDeletedBlocks：从上次lastDeletedReport后的接收到的block汇报给namenode；所有新加入的block中；当本次汇报完了以后会清空pendingIncrementalBRperStorage队列；如果汇报失败会将合并到下一次的blockReport；从上次lastBlockReport计算是否过了dnConf.blockReportInterval（默认为1个小时）时间，进行

```
bpNamenode.blockReport(bpRegistration, bpos.getBlockPoolId(), reports);
```

## Namenode blockreport process

查看namenode如何处理IncrementalBlockReport，处理入口是：

```
for (ReceivedDeletedBlockInfo rdbi : srdb.getBlocks()) {
    switch (rdbi.getStatus()) {
        case DELETED_BLOCK:
            removeStoredBlock(rdbi.getBlock(), node);
            deleted++;
            break;
        case RECEIVED_BLOCK:
            addBlock(node, srdb.getStorage().getStorageID(),
                rdbi.getBlock(), rdbi.getDelHints());
            received++;
            break;
        case RECEIVING_BLOCK:
            receiving++;
            processAndHandleReportedBlock(node, srdb.getStorage().getStorageID(),
                rdbi.getBlock(), ReplicaState.RBW, null);
            break;
    }
}
```

根据不同的状态对block处理不同，需要关注RECEIVED\_BLOCK、RECEIVING\_BLOCK的处理：

```
1 | BlockInfo storedBlock = blocksMap.getStoredBlock(block);
2 |     if(storedBlock == null) {
3 |         // If blocksMap does not contain reported block id,
4 |         // the replica should be removed from the data-node.
5 |         toInvalidate.add(new Block(block));
6 |         return null;
7 |     }
```

如果在blocksmmap里面找不到该block，就判断该block为invalidate block；并且添加到toInvalidate队列里面。全量blockReport相似。

## 3) FileJournalManager editlog

rollEditLog生成入口如下：

```
1 | CheckpointSignature rollEditLog() throws IOException {
```

```

2      checkSuperuserPrivilege(); 3      checkOperation(OperationCategory.JOURNAL);
4      writeLock();
5      try {
6          checkOperation(OperationCategory.JOURNAL);
7          checkNameNodeSafeMode("Log not rolled");
8          if (Server.isRpcInvocation()) {
9              LOG.info("Roll Edit Log from " + Server.getRemoteAddress());
10             }
11             return getFSImage().rollEditLog();
12         } finally {
13             writeUnlock();
14         }
15     }

```

namenode正常运行过程中，每个edit log最大为1M。每两分钟生成一个新的editlog。也就是说standby namenode

ed edit log的加载入口为loadFSEdits，实现过程如下：

```

1      long startTime = now();
2      FSImage.LOG.info("Start loading edits file " + edits.getName());
3      long numEdits = loadEditRecords(edits, false, expectedStartingTxId,
4          startOpt, recovery);
5      FSImage.LOG.info("Edits file " + edits.getName()
6          + " of size " + edits.length() + " edits # " + numEdits
7          + " loaded in " + (now()-startTime)/1000 + " seconds");

```

将目前的editlog加载到内存中，尽最大可能和active namenode保持同步。下面这段代码值得关注，可以查看加载

```

1      if (FSNamesystem.LOG.isDebugEnabled()) {
2          FSNamesystem.LOG.debug(op.opCode + ": " + path +
3              " numblocks : " + addCloseOp.blocks.length +
4              " clientHolder " + addCloseOp.clientName +
5              " clientMachine " + addCloseOp.clientMachine);
6      }

```

如果需要查看standby Namenode edit log详细同步加载过程，可以打开该方法对于class文件的debug日志，具体

```
hadoop daemonlog --setlevel Namenode: 9000 org.apache.hadoop.HDFS.server.Namenode. FSEd
```

### 3. 问题解决

通过上面分析日志和源码分析我们得出如下结论：

1. 当 standby 变为 active Namenode 时候，HDFS ha发生了切换。如果Datanode 的incremental block report 先于 standby Namenode 的 block report ；当 standby 完全切换完后变为active Namenode，这些新加入的block依然是invalidate block。需要Datanode 定期向 active Namenode 汇报block。在生产集群出现这样的情况，不可能当HA发生切换的时候去重启invalidate block 对应的Datanode。可以采用如下方案：

写评论

目录

收藏

微信

微博

QQ

- Datanode发现namenode 发生切换的时候做一次全量blockreport。这种方案的优点是只需要做一次全量blockreport，缺点是每次发送Namenode HA 切换所有Datanode都要重新做一次全量blockreport，通信开销较大。
- namenode专门起一个线程；当standby namenode变为activenamenode以后发现有invalidate blockreport命令。这种方案的优点是只让有部分Datanode参与全量blockreport，减少通讯开销。维护专门的线程去查找invalidate block对应的Datanode并发送相关命令。

本文为了先解决问题，实现了第一种方案。在后续的文章中我们会继续验证第二种方案的可行性：

```
1 commit edf77140caba9761e22fb6660de4f8d11add216b
2 Author: Ricky Yang <494165115@qq.com>
3 Date: Fri Apr 24 08:13:33 EDT 2015
4 missing block when HA switch
5 diff --git a/hadoop-HDFS-project/hadoop-HDFS/src/main/java/org/apache/hadoop/HDFS/server/
6 > private DatanodeProtocolClientSideTranslatorPB activeNamenode;
7 > DatanodeProtocolClientSideTranslatorPB nowactiveNamenode = bpos.getActiveNN
8 > if((nowactiveNamenode != null) && !nowactiveNamenode.equals(activeNamenode)
9 > LOG.info("active nn changed and block report agained");
10 > lastBlockReport= now()-dnConf.blockReportInterval-1;
11 > activeNamenode = nowactiveNamenode;
12 > }
```

- 1) 在BPServiceActor.java设置全局变量保持当前activeNamenode；
- 2) 获取当前的activeNamenode 保存为nowactiveNamenode；
- 3) activeNamenode不为空并且namenode发生了切换时；
- 4) 调整lastblockreport时间，让BPServiceActor满足blockreport时间；

```
1 if (startTime - lastBlockReport <= dnConf.blockReportInterval) {
2     return null;
3 }
```

1

保存当前active Namenode;

2

由该BPServiceActor 进行block report

写评论

处巧妙修改了block report时间判断条件lastBlockReport，提前触发了全量blockreport。

目录

测试总结

收藏

编译HDFS源码，替换相应的hadoop-HDFS-2.4.0.jar包，重启HDFS和HBase。进行下面的测试：

Test-1	HDFS 可用性
测试步骤	1 ) 持续写入HDFS小文件 2 ) kill 掉active Namenode 3 ) 读出文件对比 4 ) 查看后台日志和web页面
预期	1 ) HDFS 读写正确 2 ) 后台页面也web页面正常
测试结果	符合预期

Test-1	HBase可用性
测试步骤	1 ) 持续写入HBase数据 2 ) kill 掉active Namenode 3 ) 读出刚才写成的数据对比 4 ) 查看后台日志和web页面
预期	1 ) HBase读写正确 2 ) HDFS后台页面也web页面正常
测试结果	符合预期

当kill active Namenode 时候，HDFS HA发生了切换。如果Datanode 的incremental block report 先于 standby 成了invalidate block；当 standby 完全切换完后变为active Namenode，这些新加入的block依然是invalidate block原因是edit log还没有同步过来，在namespace里面没有创建相应的inode和blockmap信息。当edit log同步完了以需要修改的代码较少，可以暂时使用该方案。在生产集群已经运行1个多月，目前没有暴露出问题。

开源软件出现问题是很正常的，学习分析日志并与代码结合不失是一种解决问题的好方法。

版权声明：本文为博主原创文章，未经博主允许不得转载。 <https://blog.csdn.net/u011291159/article/details/46289639>

文章标签：

hadoop

HDFS HA

大数据