# HBase 疑难杂症诊治

## hbase2.0处理rit状态记录

| 日期 | 版本号 | 类别 | 描述 |
|------|--------|------|------|
| 2019-07-05 | 1.0.0 | A | 排查hbase2.0的rit问题 |

**问题说明**

由于使用HDP3.0，HDP3.0使用的是hbase2.0.0版本，hbase的ui页面发现很多表出现了rit，删除表过程中，region的状态卡在了opening。先尝试使用hbck2工具进行修复，发现在hbase2.0的master的rpc方法中没有hbck2中的bypass，assign方法，通过源码发现，hbck2的rpc方法是在hbase2.0.2中才增加的，所以只能尝试手动处理：

**一.hdfs上已经没有对应目录，meta里没有对应表状态信息，存在有对应的分区信息**

1. 检查表状态

```
get 'hbase:meta','KYLIN_0054K9NLSU','table:state'
```

结果为空

2. 通过源码发现表状态

```
ENABLED,对应meta表里的值\x80\x00
DISABLED,   对应meta表里的值\x80\x01
DISABLING,  对应meta表里的值\x80\x02
ENABLING,   对应meta表里的值\x80\x03
```

3. 查找分区信息

```
scan                                                                 'hbase:meta',
{FILTER => org.apache.hadoop.hbase.filter.PrefixFilter.new(org.apache.hadoop.hbase.util.Bytes.toBytes('KYLIN_0054K9NLSU'
```

发现存在有分区记录

4. 手动修改表状态或者删除分区信息

```
put 'hbase:meta','KYLIN_0054K9NLSU','table:state','\x80\x01'
```

或者deleteall 表对应的分区信息，修改后重启hbase，发现rit状态消失

**二.hdfs上已经有对应目录，meta里有对应表状态信息和分区信息**

1. 确认一下表的信息和数据

```
hbase hbck -summary TableName
```

2. 检查表状态

```
get 'hbase:meta','KYLIN_0354K9NLSU','table:state'
```

meta表里的值\x80\x02，表的状态为DISABLING

3. 找出异常的region

```
echo                                    "scan                                     'hbase:meta',
{FILTER => org.apache.hadoop.hbase.filter.PrefixFilter.new(org.apache.hadoop.hbase.util.Bytes.toBytes('KYLIN_0354K9NLSU'
n|grep  OPENING
```

找出异常的region

4. 将region信息更新为CLOSED状态

```
put 'hbase:meta','KYLIN_0354K9NLSU,,1561953520536.30b7d24eaa3209c6e5e8de764ad04855.','info:state','CLOSED',1562117738678
```

………

5. 将表状态更新为disable

```
put 'hbase:meta','KYLIN_0354K9NLSU','table:state',"\x08\x01",1562120793251
```

重启hbase后rit消失

**存在问题**

- rit是删除表的时候出现，所以表中的数据可以忽略，上述操作也是表中没有数据时操作
- 如果是生成集群，已经存在的数据比较多，不建议直接重启，可以通过切换master的方式
- 可以使用HDP3.1.1，里面hbase版本是2.0.2，可以使用hbck2操作
- 使用hbck2的方法的话，修改meta状态后还会同步改zookeeper状态，能避免状态不一致

---

## HBase2.x之RIT问题解决

### 问题描述

Region-In-Trasition机制：从字面意思来看，Region-In-Transition说的是Region变迁机制，实际上是指在一次特定操作行为中Region状态的变迁，例如merge、split、assign、unssign等操作。RIT问题指的是在RIT过程中出现异常情况，然后导致region的状态一直保持在RIT，使得HBase出现异常。

**解决方案**

**方案一**

检查hdfs的健康度，是否有hbase的文件丢失或损坏，运行命令hadoop fsck /，结果如下：



排除hdfs丢失block的问题。如果出现hdfs的block损坏或丢失的情况，可以通过hdfs的修复命令进行修复。

**方案二**

在HBase1.x系列中RIT问题通常可以通过hbase    hbck –repair操作完成修复。但是在HBase2.x系列中，该命令还没有支持，所以暂时无法通过这种命令完成修复。结果如下：



第一次执行发现没有权限，root用户不是hdfs的超级用户，安装提示需要以hbase用户去执行该命令。修改如下：

```
su  hbase  -s  /bin/sh  -c  "hbase  hbck  -repair"
```



提示为hbase有其他线程正在执行hbck fix命令，但是其实没有这种命令，其实从这里就可以看出HBase2.x对于-repair的支持是不够的。我按照提示删除了hdfs(/hbase/.tmp/)上的hbase-hbck.lock文件，重新运行命令如下：
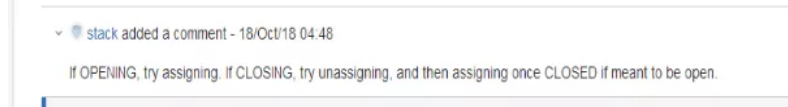


**方案三**

根据RIT状态执行assign或者unassign命令，状态信息如下：

经过查询资料，解决方案如下：



stack added a comment - 18/Oct/18 04:48

If OPENING, try assigning. If CLOSING, try unassigning, and then assigning once CLOSED if meant to be open.
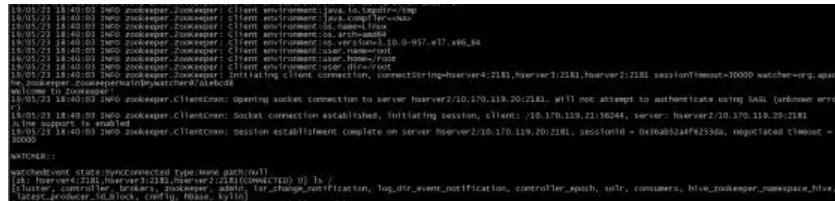
hbase shell多次执行unassign '20acfcbd68fd624a78bb34c88f9382d1'和unassign '20acfcbd68fd624a78bb34c88f9382d1'，true都超时结束，通过修改rpc和zk的超时时间都无法完成(正常超时时间为60000ms，修改为600000ms)。

## 方案四

经过多次试验，最终都无法使得HBase回复正常，最终决定删除进行测试。

Zookeeper节点删除：

通过hbase zkcli命令进入Zookeeper命令行界面：



我想删除节点 /hbase-unsecure/region-in-transition，但是发现并没有该节点，经过资料查询了解到HBase2.x的RIT状态并不像HBase1.x系列存储在Zookeeper上。经过测试删除/hbase节点重启hbase并不能解决RIT问题。

HBase表删除：

`hbase shell>disable M_TDY_PT_LCZZ`

disable失败，所以正常删除表操作无法执行。需要进行暴力删除，暴力删除指的是从元数据进行删除。

先停掉HBase



删除hdfs表目录(记得先备份，等下恢复用)

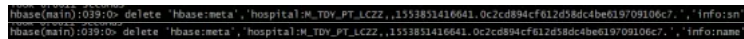`hdfs dfs -cp /hbase/data/hospital/P_TDY_DASC_DE02_01_039_63 /`
`hdfs dfs -cp /hbase/data/hospital/M_TDY_PT_LCZZ /`
`hdfs dfs -rm -r -skipTrash /hbase/data/hospital/P_TDY_DASC_DE02_01_039_63`
`hdfs dfs -rm -r -skipTrash /hbase/data/hospital/ M_TDY_PT_LCZZ`
`delete 'hbase:meta','rowkey', 'column'`

Rowkey信息可以通过hbase的UI看到：



然后重启HBase，但是发现问题没有解决。

hbase shell查询数据看到hbase的meta删除失败了，原本的meta信息还在：



然后重新删除(delete命令换成deleteall命令)：



再删除Zookeeper中的/hbase节点，重启HBase发现RIT问题已经解决了。

后续就是重建表，然后恢复数据。


## Phoenix故障处理笔记

1. Timeline

- 06-26 16:00 Phoenix使用方反馈慢；
- 06-26 16:02 同事通过监控看到Phoenix HBase集群一个对应的RegionServer，QueueSize过高，此bug基本是Butch Put Lock在高并发写入的问题，我们已在下个版本中增加信息日志定位此问题；
- 06-26 16:05 同事重启该队列过高的RegionServer；
- 06-26 16:10 同事跟我说，好多Phoenix的Region处于RIT状态；
- 06-26 17:00 暂停该Phoenix集群所有的写入；
- 06-26 20:00 跟业务沟通，可能会正常影响一段时间，经同意。至此各种hbck，各种重启RegionServer&Master不怎么管用，RIT数量升至550个；
- 06-27 12:00 尝试修复；
- 06-27 15:00 问题修复。

## 2. 处理流程

### 2.1 异常现象

1. 大量Region无法上线(NotServingRegionException)



2. Phoenix的SYSTEM.CATALOG系统表也无法上线。



### 2.2 处理流程

1. 手动assign SYSTEM.CATALOG 表的Region上线，并且跟踪Master&对应RegionServer的日志。整个offline&open流程都正常。但是中间会由于各种其他的表不在线failover后close掉；

2. 打印jstack, 感觉这几个线程有问题，都在waiting；

```
RS_OPEN_REGION█▐▟ ▜▟▗▟▄▟ J
priority:5 - threadId:0x00007f178ecc3800 - nativeId:0xfff1e - nativeId (decimal):1048350 - state:WAITING
stackTrace:
java.lang.Thread.State: WAITING (parking)
at sun.misc.Unsafe.park(Native Method)
- parking to wait for <0x0000000282c5e738> (a com.google.common.util.concurrent.AbstractFuture$Sync)
at java.util.concurrent.locks.LockSupport.park(LockSupport.java:175)
at java.util.concurrent.locks.AbstractQueuedSynchronizer.parkAndCheckInterrupt(AbstractQueuedSynchronizer.java:836)
at java.util.concurrent.locks.AbstractQueuedSynchronizer.doAcquireSharedInterruptibly(AbstractQueuedSynchronizer.java:997)
at java.util.concurrent.locks.AbstractQueuedSynchronizer.acquireSharedInterruptibly(AbstractQueuedSynchronizer.java:1304)
at com.google.common.util.concurrent.AbstractFuture$Sync.get(AbstractFuture.java:275)
at com.google.common.util.concurrent.AbstractFuture.get(AbstractFuture.java:111)
at org.apache.phoenix.hbase.index.parallel.BaseTaskRunner.submit(BaseTaskRunner.java:66)
at org.apache.phoenix.hbase.index.parallel.BaseTaskRunner.submitUninterruptible(BaseTaskRunner.java:99)
at org.apache.phoenix.hbase.index.write.recovery.TrackingParallelWriterIndexCommitter.write(TrackingParallelWriterIndexCommitter.java:197)
at org.apache.phoenix.hbase.index.write.IndexWriter.write(IndexWriter.java:185)
at org.apache.phoenix.hbase.index.write.RecoveryIndexWriter.write(RecoveryIndexWriter.java:75)
at org.apache.phoenix.hbase.index.Indexer.preWALRestore(Indexer.java:504)
at org.apache.hadoop.hbase.regionserver.RegionCoprocessorHost$58.call(RegionCoprocessorHost.java:1422)
at org.apache.hadoop.hbase.regionserver.RegionCoprocessorHost$RegionOperation.call(RegionCoprocessorHost.java:1663)
at org.apache.hadoop.hbase.regionserver.RegionCoprocessorHost.execOperation(RegionCoprocessorHost.java:1738)
at org.apache.hadoop.hbase.regionserver.RegionCoprocessorHost.execOperation(RegionCoprocessorHost.java:1695)
at org.apache.hadoop.hbase.regionserver.RegionCoprocessorHost.preWALRestore(RegionCoprocessorHost.java:1413)
at org.apache.hadoop.hbase.regionserver.HRegion.replayRecoveredEdits(HRegion.java:4062)
at org.apache.hadoop.hbase.regionserver.HRegion.replayRecoveredEditsIfAny(HRegion.java:3919)
at org.apache.hadoop.hbase.regionserver.HRegion.initializeRegionStores(HRegion.java:986)
at org.apache.hadoop.hbase.regionserver.HRegion.initializeRegionInternals(HRegion.java:858)
at org.apache.hadoop.hbase.regionserver.HRegion.initialize(HRegion.java:832)
at org.apache.hadoop.hbase.regionserver.HRegion.openHRegion(HRegion.java:5973)
at org.apache.hadoop.hbase.regionserver.HRegion.openHRegion(HRegion.java:5937)
at org.apache.hadoop.hbase.regionserver.HRegion.openHRegion(HRegion.java:5907)
at org.apache.hadoop.hbase.regionserver.HRegion.openHRegion(HRegion.java:5861)
at org.apache.hadoop.hbase.regionserver.HRegion.openHRegion(HRegion.java:5810)
at org.apache.hadoop.hbase.regionserver.handler.OpenRegionHandler.openRegion(OpenRegionHandler.java:356)
at org.apache.hadoop.hbase.regionserver.handler.OpenRegionHandler.process(OpenRegionHandler.java:126)
at org.apache.hadoop.hbase.executor.EventHandler.run(EventHandler.java:128)
at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1142)
at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:617)
at java.lang.Thread.run(Thread.java:745)
```

- 通过上面的信息看，open region确实有问题。查看Phoenix Indexer Observer源码就会知道是在根据Recover WAL Entry构建索引；
- 修改hbase.regionserver.executor.openregion.threads数,此配置是负责open region的handler数：

```xml
<property>
    <name>hbase.regionserver.executor.openregion.threads</name>
    <value>50</value>
</property>
```

默认 3， 我们这边的hbase版本(1.0.0-cdh5.4.4)

- **重启RegionServer;**
- assign SYSTEM.CATAOG 表的Region，上线成功；
- 修修补补，fixMeta fixAssignments就ok了。

## 3．原理分析

1．重启RegionServer， 会造成该RegionServer上面的Region下线，并且被重新Balance到新的RegionServer中。

2．Region在新的RegionServer中open过程会找到该Region下的recover.edits 文件，进行replay；

3.Phoenix表使用HBase的协处理类型之Observer，具体使用查看示例

org.apache.phoenix.hbase.index.Indexer，此用作根据WAL构建索引的，具体参考Phoenix的相关材料。

4．在SYSTEM.CATALOG 的打开过程中，会查询其他的里面表，其他的表也处于RIT未恢复。然而其他的表Region在open的过程也需要构建Index，尚且有一部分在openregion的队列里面。最终SYSTEM.CATALOG无法上线(此处不准确，纯属囫囵吞枣似的查看源码推测)。

5．增加open region handler数之后，重启RegionServer后，需要进行一些hbck -fixMeta -fixAssginment 将一些未上线的Region上线，就ok了。

6．如果出现个别的Region还是上线失败，那就手动解决吧！个人认为比hbck -repair暴力修复靠谱。