

# Cassandra 数据一致性修复 repair 来龙去脉

来自：阿里云数据库 ⌚ 2020-06-12

**简介：** 文章分3块：1.为什么需要repair？；2.repair大概流程？；3.repair可能的问题。

## Cassandra repair 流程

文章分3块：1.为什么需要repair？；2.repair大概流程？；3.repair可能的问题。

### 为什么需要repair

#### cassandra 如何修复副本数据

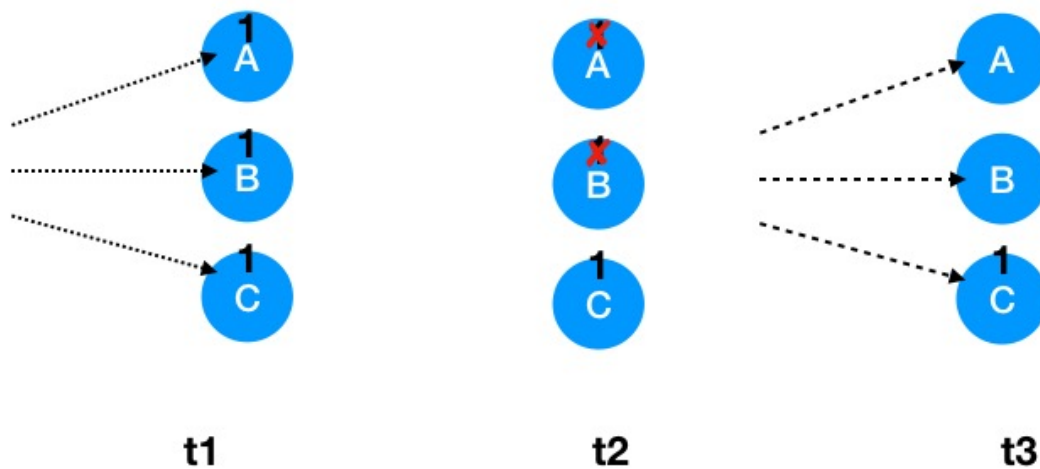
我们知道Cassandra是一个强调最终一致性的系统，副本间的数据并不能保证强一致(但是从客户端角度，通过QUORUM等级别读写还是可以保证客户端视角的强一致[1])。因为副本间数据是最终一致，所以Cassandra通过hinted-handoff、read-repair、repair 进行副本间数据补齐；这三个方式各有优缺点：

- hinted-handoff 用来补齐节点挂掉期间的数据，但是挂掉时间太久这个特性会失效且hint机制存在数据可靠性风险；
- read-repair：只有被读到的不一致数据才会被修复，那么如果没有被读到的数据很多怎么办？
- repair：全量的修复方案，保证做多个副本在某个时间点（触发repair时刻）前的数据全量修复一致。

从上述描述看出，repair是一个兜底修复副本数据的方案，那么既然上面说了，通过quorum可以保证客户端视角的强一致，我们还需要通过repair来修复全量副本数据么？答案是：必须要，而且官方也建议在一定周期（gc\_grace\_seconds）内必须要做一次[2]，才能保证系统的正确性。

### 如果不做repair会怎么样？

delete .....  
Read - - - - -

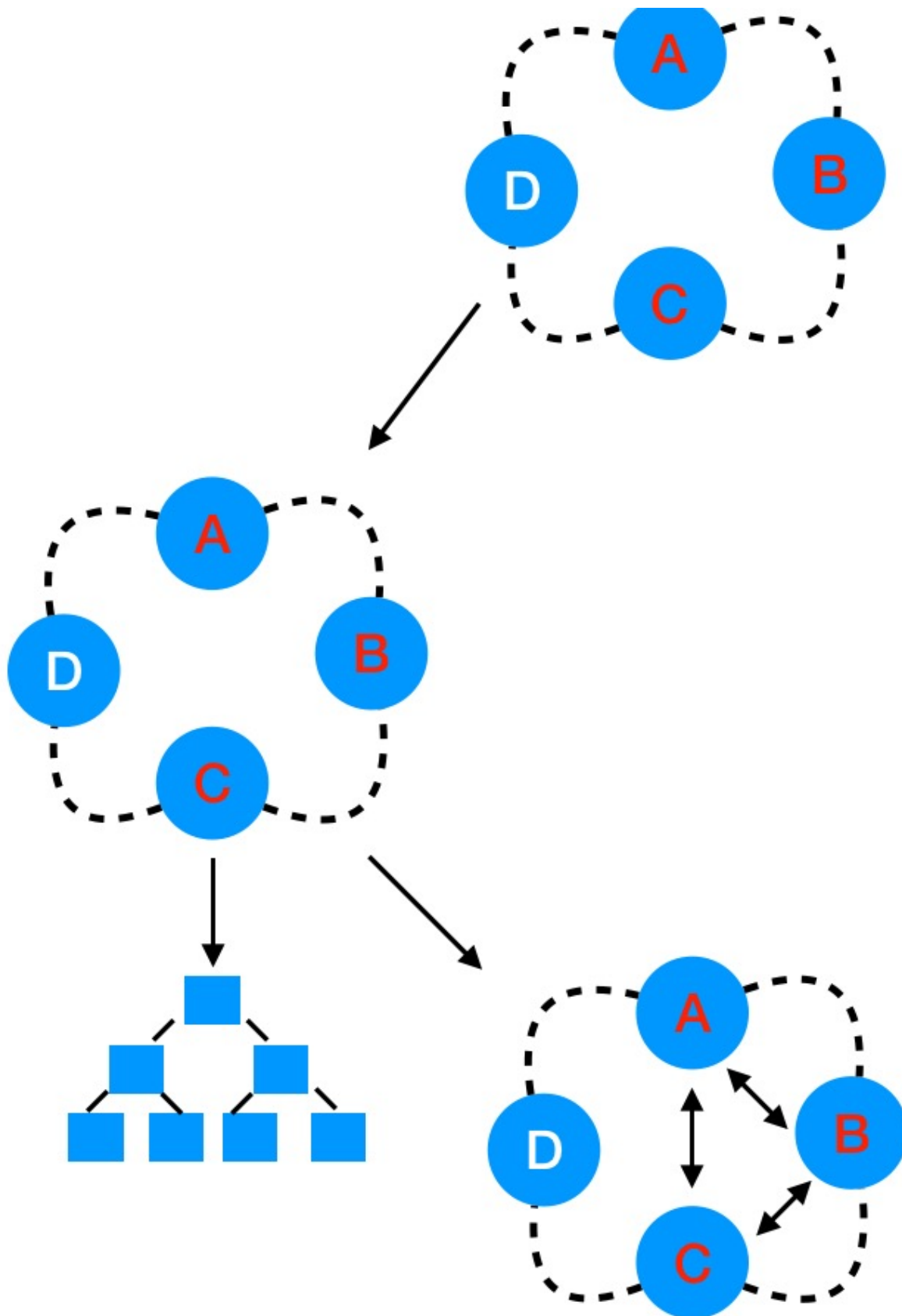


假设3个副本(A/B/C)，quorum级别的读写删数据，时刻t1用户最初以quorum级别成功写入数据1，假设A/B/C都写入成功；时刻t2用户quorum级别删除1数据，假设这里C副本删除失败，但是客户依旧显示删除成功；时刻t3（这里假设 $t3 - t2 > gc\_grace\_seconds$ ，且用户没有做repair且期间）。那么A/B副本会compaction把1数据合并删除掉，但是c副本没有删除mark。

最终结果一条被用户认为删除成功的key，“死灰复燃”的读到了。

所以：repair必须要做。

## repair大概流程



全量数据repair需要人为手工通过节点nodetool提交外围任务，具体的nodetool 命令行参数 可以下次介绍，这里大概分享下一轮repair在副本节点之间进行的流程。假设开始做A 节点负责的数据，对应副本涉及B/C节点，那么一次执行的修复链路是：

- 计算A/B/C三个相关副本数据的全量merkle-tree[3]，这是一个二分hash树；从底往上计算hash、叶子节点是小range范围的数据的sha2 hash值，内部节点是其左右子节点的hash值（xor）；

- 通过两两对比merkle-tree, 可以知道具体节点间不一致的数据范围;
- 两两走内部stream 拖数据流程补齐相关数据。

## repair可能的问题

- 运维复杂:
  - 节点数据量如果较大, 整个执行过程可能会耗时很久, 时间越久出现问题的可能性就越大;
  - 为了避免repair对集群影响较大, repair需要针对节点差异化执行, 那么对运维复杂性会带来挑战;
  - 需要表级别gc\_grace\_seconds 内做一次, 如果表过多, 会造成运维差异化难度较大;
- 资源消耗较大:
  - 一般cassandra被用于在线服务场景, 但是做repair会带来瞬时资源较大开销: cpu、io、网络, 影响服务稳定性;

现在社区解决方案有: incremental repair、schedule repair等方案[4],此外datastax公司也有nodesync[5],scylladb 公司有row-level repair。对应我们也有相关的定制功能。目的是降低运维复杂度, 降低repair时候对在线服务的影响。

引用:

[1] [https://www.allthingsdistributed.com/2008/12/eventually\\_consistent.html](https://www.allthingsdistributed.com/2008/12/eventually_consistent.html)

[2] <https://cassandra.apache.org/doc/latest/operating/repair.html>

[3] [https://en.wikipedia.org/wiki/Merkle\\_tree](https://en.wikipedia.org/wiki/Merkle_tree)

[4] <https://issues.apache.org/jira/browse/CASSANDRA-14346>

[5] [https://docs.datastax.com/en/dse/6.0/dse-dev/datastax\\_enterprise/config/aboutNodesync.html](https://docs.datastax.com/en/dse/6.0/dse-dev/datastax_enterprise/config/aboutNodesync.html)

[6] <https://docs.scylladb.com/operating-scylla/procedures/maintenance/repair/>