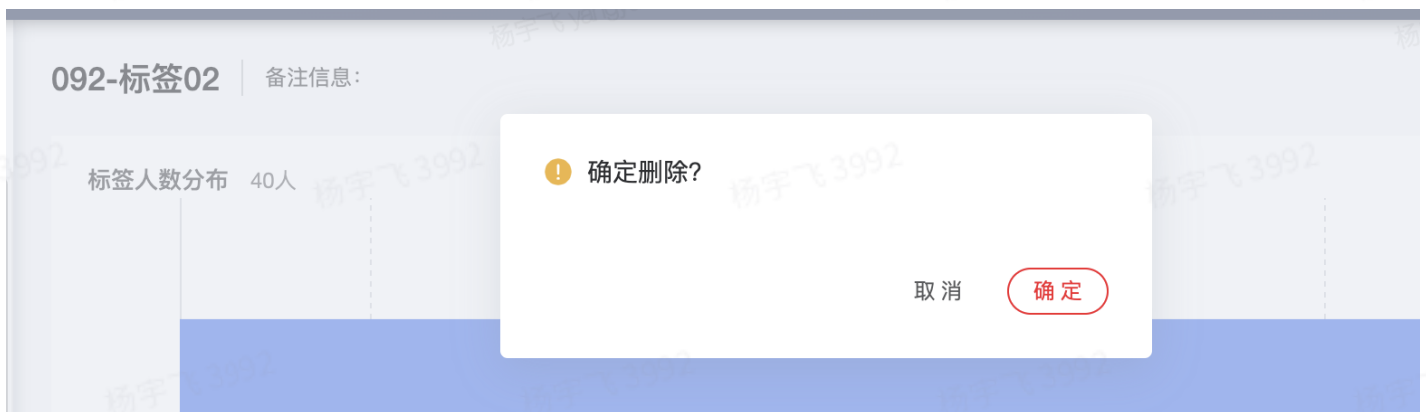
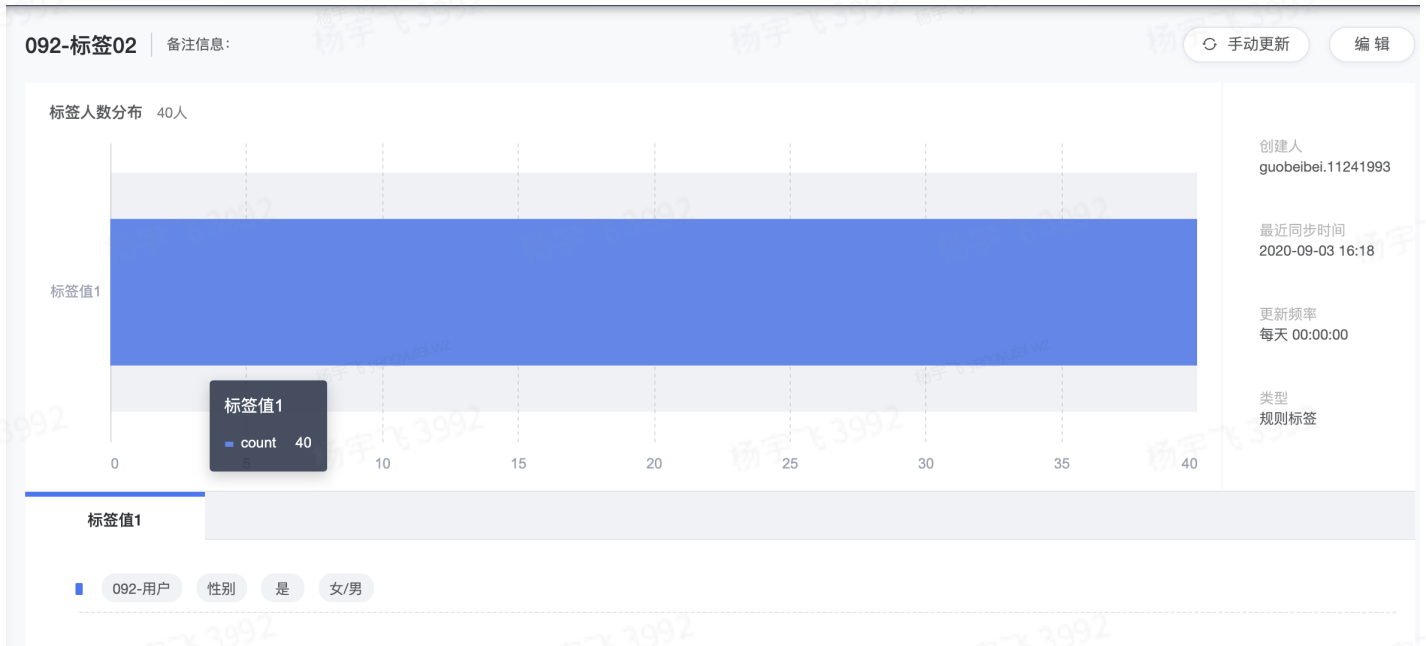


# 一场由删除标签引发的悲剧

## 背景介绍



## v1.0 仅删除mysql记录

### 方法实现:

- 删除列: 仅删除mysql记录, 将标签状态由0改为1表示删除
- 创建标签, 分配列:
  - 首先, 获取clickhouse表结构, 然后获取mysql中没有被删除的标签
  - (1) 如果clickhouse中的列都被分配, 则重新创建列分配标签
  - (2) 如果clickhouse中有列未分配, 则从未分配列中选择一个分配给新创建的标签

## 导致问题

- 历史数据遗留
  - 由于删除标签时没有删除clickhouse中的数据，当下次创建标签还可以分配到这一列
  - 此时，上一个标签的数据仍然存在，如果用该标签创建群体就可以看到历史标签值

|      | Mysql标签                                   | 新增标签前Clickhouse                              | 新增标签           | 新增标签后Clickhouse   |
|------|---|--|----------------|---|
| 正常场景 | 1 label_1 0<br>2 label_2 0<br>3 label_3 0 | label_1 label_2<br>label_3<br>标签值1 标签值2 标签值3 | 4 label_4<br>0 | label_1 label_2 label_3<br>label_4<br>标签值1 标签值2 标签值3 标签值4 |
| 问题场景 | 1 label_1 0<br>2 label_2 1<br>3 label_3 0 | label_1 label_2<br>label_3<br>标签值1 标签值2 标签值3 | 4 label_2<br>0 | label_1 label_2 label_3<br>标签值1 标签值2 标签值3<br>标签值4         |

## 解决方案

- 删除标签时，将数据列一并删除

## v2.0 删除标签同时删除clickhouse列

### 方法实现：

- 调用clickhouse的drop\_column接口进行删除

## 导致问题

- 列删除仍可能被分配，但是最终标签计算失败
  - 由于Clickhouse删除列有延迟，删除标签后马上创建新的标签时，clickhouse中的列还没有被删除
  - 此时，创建标签获取clickhouse表结构中仍然包含刚被删除的列
  - 但是数据库中的记录已经被删除，所以利用上面的分配列策略，仍然可能分配到这一列
  - 从而，导致虽然分配列成功了，但是过一段时间这一列就会从clickhouse中删除，最终导致导入标签数据失败

|      | Mysql标签                    | 新增标签前Clickhouse                           | 新增标签        | 新增标签后Clickhouse              |
|------|----------------------------|---|-------------|------------------------------|
| 问题场景 | 1 label_1 0<br>2 label_2 1 | label_1 label_2 label_3<br>标签值1 标签值2 标签值3 | 4 label_2 0 | label_1 label_3<br>标签值1 标签值3 |

|  |             |  |  |
|--|-------------|--|--|
|  | 3 label_3 0 |  |  |
|--|-------------|--|--|

## 解决方案

- 设置正在删除标志位
- 异步检查clickhouse中是否删除成功，删除成功后才将标签状态修改为删除

## v3.0 设置正在删除标志位

### 方法实现：

- 删除标签时，首先将状态改为2，表示正在删除
- 发送clickhouse删除列请求后，开启一个线程异步检测列是否删除成功，如果删除成功再将标签的状态改为1
- 同时，创建标签时，获取mysql中没有被删除的标签会将正常标签和删除中标签，从而与clickhouse的表结构保持一致

|                | Mysql标签  | Clickhouse                                | 新增标签        |                                 |
|----------------|--|---|-------------|---------------------------------|
|                |  |   |             | 异步线程检测<br>clickhouse中标<br>签未删除  |
| 删除标签后，<br>分配列时 | 1 label_1 0<br>2 label_2 2<br>3 label_3 0                | label_1 label_2 label_3<br>标签值1 标签值2 标签值3 | 4 label_4 0 |                                 |
|                |  |   |             | 异步线程检测<br>clickhouse中标<br>签列已删除 |
| 分配列后，<br>导入数据时 | 1 label_1 0<br>2 label_2 1<br>3 label_3 0<br>4 label_4 0 | label_1 label_3 label_4<br>标签值1 标签值3 标签值4 |             |                                 |

### 导致问题：

- 列删除仍可能被分配
- 多线程同步问题：

|  |  |  |  |
|--|--|--|--|
|  |  |  |  |
|--|--|--|--|

|                                 | Mysql标签  | Clickhouse                                   | 新增标签        |                                |
|---------------------------------|--|--|-------------|--------------------------------|
|                                 |  |  |             | 异步线程检测<br>clickhouse中标<br>签未删除 |
| 删除标签后，分<br>配列时获取<br>clickhouse列 | 1 label_1 0<br>2 label_2 2<br>3 label_3 0                | label_1 label_2<br>label_3<br>标签值1 标签值2 标签值3 |             |                                |
|                                 |  |  |             | 异步线程检测<br>clickhouse中标<br>签已删除 |
| 删除标签后，分<br>配列时获取mysql<br>标签并分配列 | 1 label_1 0<br>2 label_2 1<br>3 label_3 0                | label_1 label_3<br>标签值1 标签值3                 | 4 label_2 0 |                                |
| 分配列后，<br>导入数据时                  | 1 label_1 0<br>2 label_2 1<br>3 label_3 0<br>4 label_2 0 | label_1 label_3<br>标签值1 标签值3                 |             |                                |

- 新问题：两个标签同时创建，请求分配到不同的服务器，可能分配到同一列

#### 解决方案：

- 修改分配列逻辑
- 设置分布式锁

## v4.0 调整删除策略并添加分布式锁

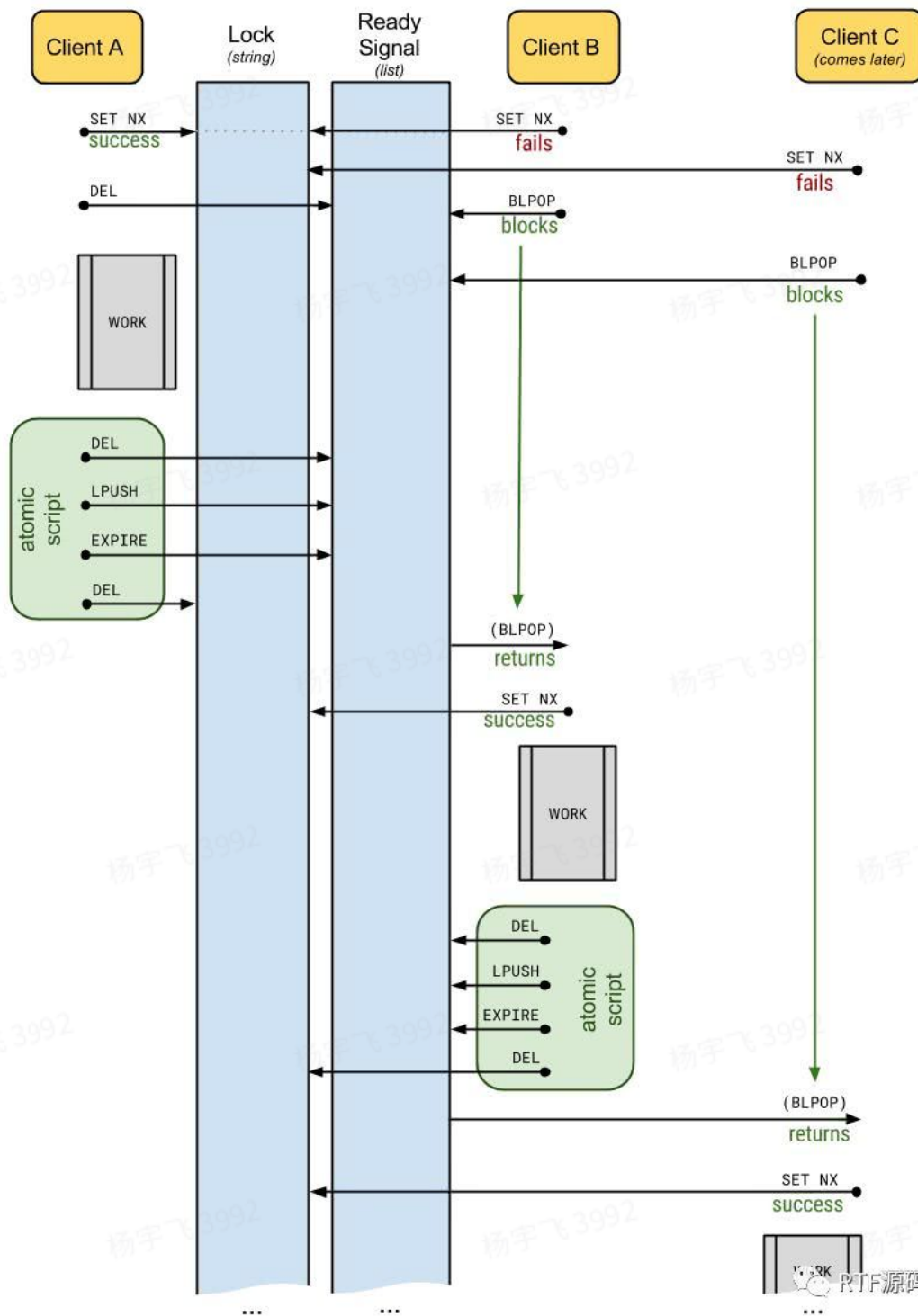
#### 方法实现：

- 修改分配列逻辑：由于删除标签不会再次分配，所以获取mysql中App下所有的标签（包括正常标签和已删除标签），只分配clickhouse中空闲，未在mysql中有记录的列

|      | Mysql标签                                   | 新增标签前Clickhouse                              | 新增标签           | 新增标签后Clickhouse   |
|------|---|--|----------------|---|
| 正常场景 | 1 label_1 0<br>2 label_2 0<br>3 label_3 0 | label_1 label_2<br>label_3<br>标签值1 标签值2 标签值3 | 4 label_4<br>0 | label_1 label_2 label_3<br>label_4<br>标签值1 标签值2 标签值3 标签<br>值4 |

|      |   |  |                |  |
|------|---|--|----------------|--|
| 删除场景 | 1 label_1 0<br>2 label_2 1<br>3 label_3 0 | label_1 label_2<br>label_3<br>标签值1 标签值2 标签值3 | 4 label_4<br>0 | label_1 label_3 label4<br>标签值1 标签值3 标签值4 |
| 删除场景 | 1 label_1 0<br>2 label_2 1<br>3 label_3 0 | label_1 label_3<br>标签值1 标签值3                 | 4 label_4<br>0 | label_1 label_3 label4<br>标签值1 标签值3 标签值4 |

- 设置分布式锁：为了避免分布式场景下，分配到同一个列，使用redis实现分布式锁进行解决



- acquire\_lock

```

1 def acquire_lock(cls, cluster_name, db_name=None, table_name=None,
2   data_type_name=None, lock_timeout=60, acquire_timeout=60):
3     k = cls.gen_key(cluster_name=cluster_name, db_name=db_name,
4       table_name=table_name,
5         data_type_name=data_type_name)
6     identifier = str(uuid.uuid4())
7     start = time.time()
8     end = time.time() + acquire_timeout
  
```

```

8     while time.time() < end:
9         if cls.client.setnx(k, identifier, _raise=True):
10            cls.client.expire(k, lock_timeout)
11            logger.info('locking time: {}'.format(time.time() - start))
12            logger.debug('acquire lock_lock: {}'.format(time.time()))
13            return identifier
14        elif not cls.client.ttl(k):
15            cls.client.expire(k, lock_timeout)
16
17        time.sleep(0.1)

```

- release\_lock

```

1 def release_lock(cls, identifier, cluster_name, db_name=None, table_name=None,
2   data_type_name=None):
3     while True:
4         try:
5             iden = cls.get(cluster_name=cluster_name, db_name=db_name,
6   table_name=table_name, data_type_name=data_type_name, _raise=True)
7             if iden and iden.decode('utf-8') == identifier:
8                 cls.del_one(cluster_name=cluster_name, db_name=db_name,
9   table_name=table_name, data_type_name=data_type_name)
10                return True
11            break
12        except Exception as e:
13            logger.info('release_lock exception: {}'.format(e))
14            raise e
15    return False

```

## 存在问题和下一步改进:

- 释放锁并非原子操作:
  - 事务方法: 公司pyredis不支持事务, 因为redis的事务处理并不好
  - 执行LUA脚本方法: 目前风神的redis集群不支持, 需要升级
- Redis的Master节点宕机: [分布式锁proposal](#)
  - RedLock算法
  - 使用etcd分布式kv存储: 采用 `raft` 协议作为一致性算法
  - 使用zookeeper框架

## 参考:

[Redis 新人学习相关资料汇总](#)

<https://blog.halukshan.com/blog/2020/03/31/redis-distributed-lock-python/>

<https://zhuanlan.zhihu.com/p/112016634>

<https://www.cnblogs.com/angelyan/p/11523846.html>

杨宇飞 3992

杨宇飞 3992

杨宇飞 3992

杨宇飞?

3992

杨宇飞 3992

杨宇飞 3992

杨宇飞 3992

杨宇飞 3992

杨宇飞 3992

杨宇飞 3992

杨宇飞?

3992

杨宇飞 3992

杨宇飞 3992

杨宇飞 3992

杨宇飞 3992

杨宇飞 3992

杨宇飞 3992

杨宇飞?

3992

杨宇飞 3992

杨宇飞 3992

杨宇飞 3992

杨宇飞 3992

杨宇飞 3992

杨宇飞 3992

杨宇飞?

3992

杨宇飞 3992

杨宇飞 3992

杨宇飞 3992

杨宇飞 3992

杨宇飞 3992

杨宇飞 3992

杨宇飞?