

携程 ClickHouse 集群 数据管理工具 实战

林东煜

2025.02.29

目录

CONTENTS

- 1 背景与问题
- 2 设计与调研
- 3 目标与成果

个人简介

林东煜-携程-系统研发部

- ClickHouse 日志
- OLAP Paas平台

背景与问题

- 背景:

- 日志Clickhouse集群 20.3版本
- Alma系统服务器升级
- 上云/下云/迁云
- 测试环境机房搬迁

- 痛点:

- 原地**升级**无法灰度, 用户有感
- 集群无法**缩容**
- 集群数据不**均衡**
- 机器空间常年处于80%以上, **压缩**程度不够
- 旧版本人力**运维**成本高Metadata on replica is not up to date with common metadata in Zookeeper

缺乏弹性

服务器: 1000+

数据量: 120w亿行

数据表: 1w+

存储空间: 70+PB

查询: 10M+/天

集群数量: 40+

P90: 500ms

P99: 3s

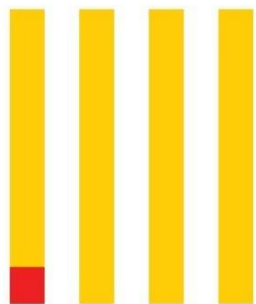
为什么要升级

运维效率
与新特性

• 20.3 VS 24.3

- Atomic Engine
- Projection
- JSON type
- lightweight delete
- Variant type
- New analyzer
- EXPLAIN
- NEW FUNCTIONS
- NEW COMANDS FOR OPS
- S3 ZERO COPY
- async_insert

• 元数据不一致 <https://github.com/ClickHouse/ClickHouse/issues/54902>



ClickHouse

为什么要数据均衡

空间与负载

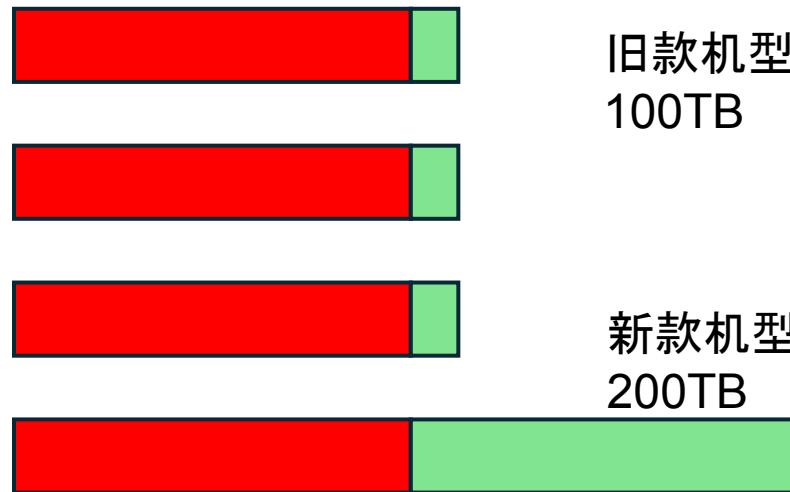
Case1:

新扩机器数据不均衡
(空间,cpu,磁盘io)



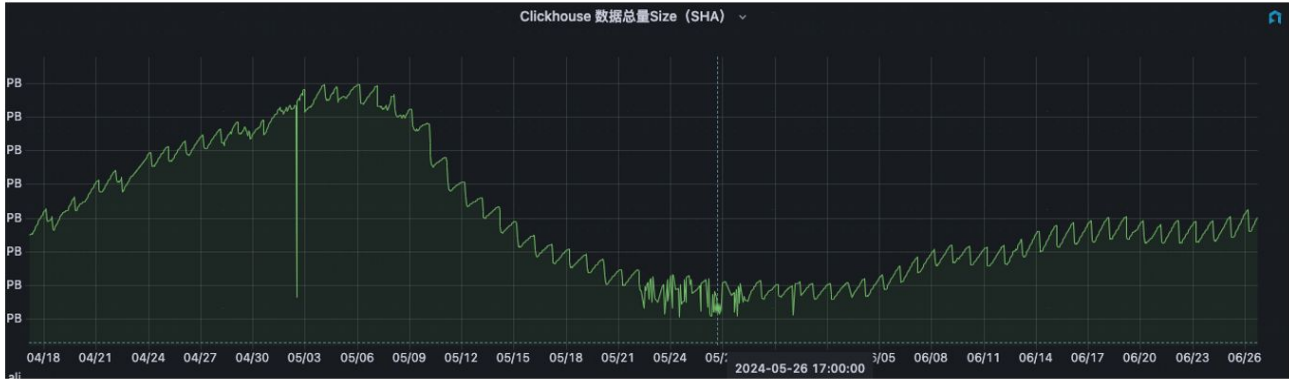
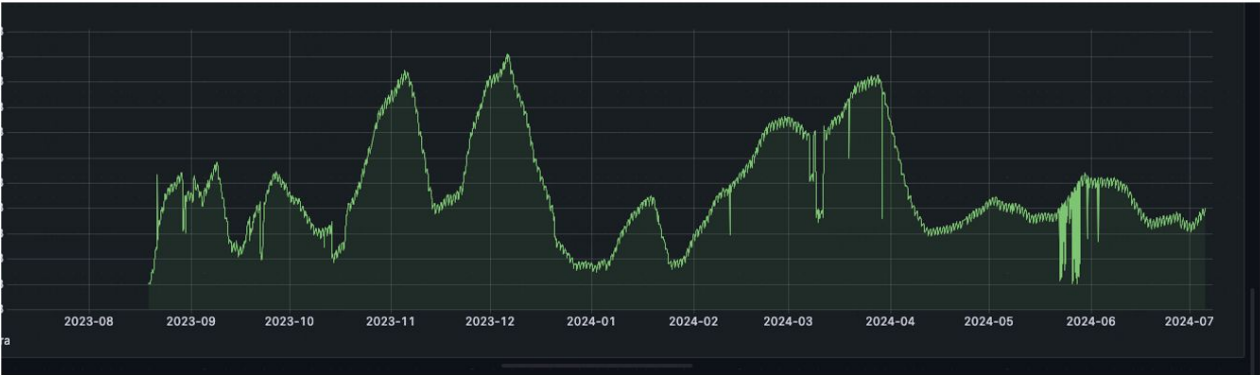
Case2:

机型不统一数据不均衡
(空间,cpu,磁盘io)



为什么要缩容

30%以上的
波动



设计与调研

方案调研：

- Copier
- insert into ...
Select
- File copy
- Fetch

设计与调研

File

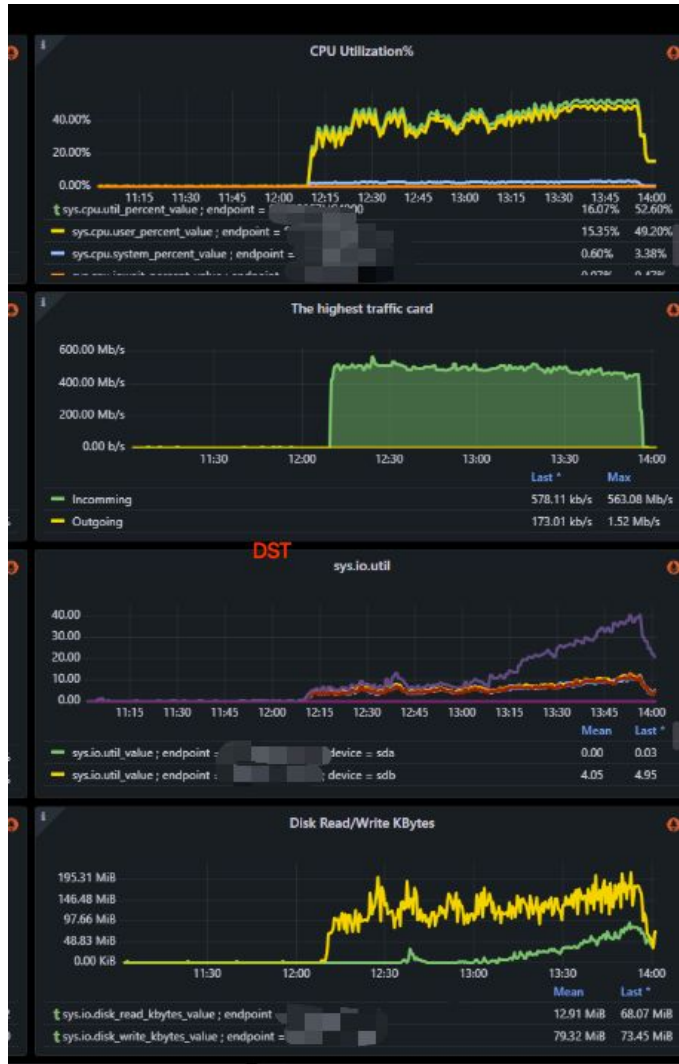
- JuiceFS
- 命令下发
- 文件离散、打包
- 并发传输
- 任务管控

设计与调研

insert into ... Select

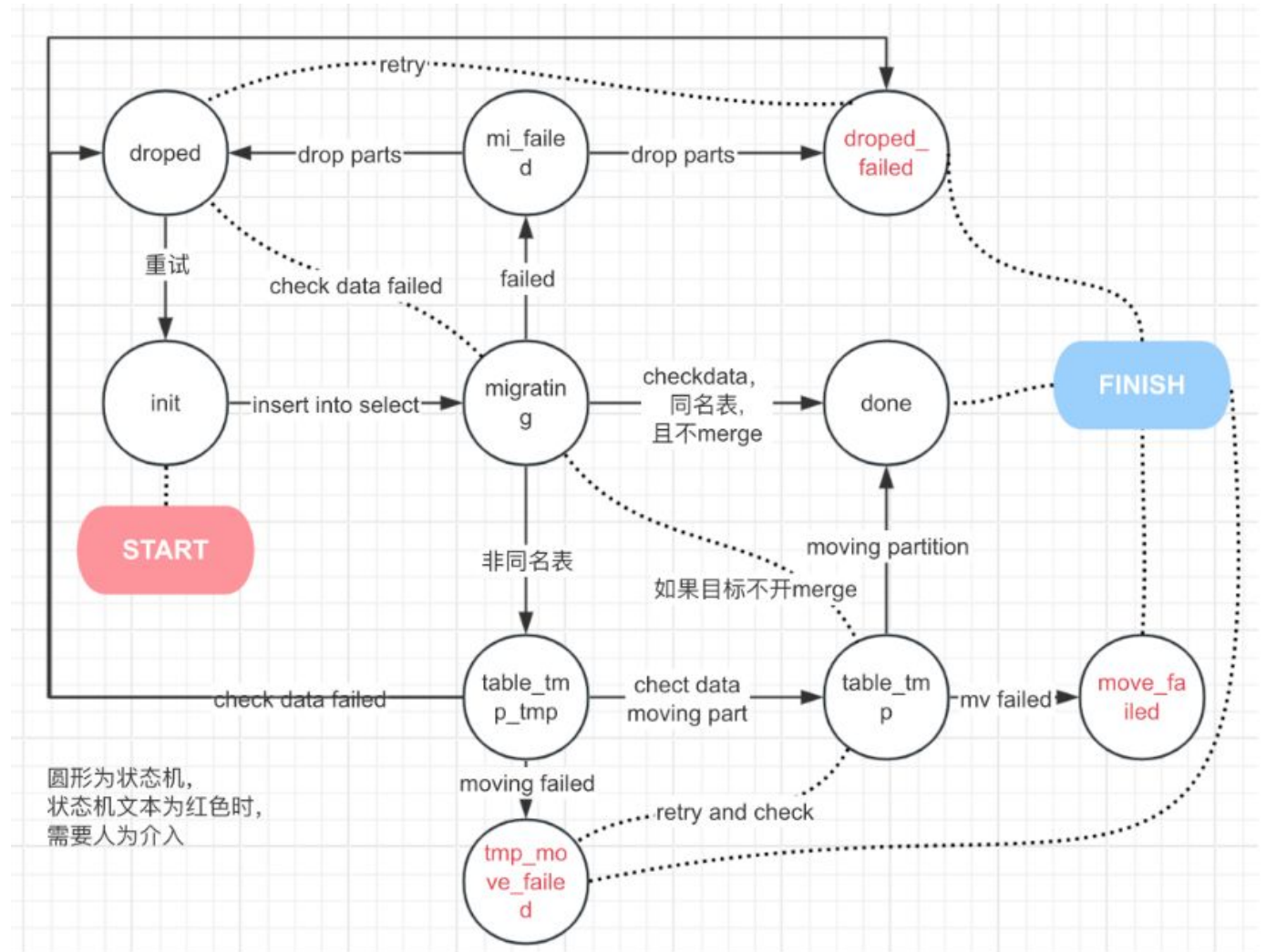
- virtual-columns(_part/_partition_id)
- insert_quorum/optimize_on_insert
- min_insert_block_size_rows/min_insert_block_size_bytes
- max_execution_time
- max_bytes_to_read
- max_threads/max_threads_insert
- min_compress_block_size/max_compress_block_size
- Queryid、part_log
- system FLUSH LOGS 刷system.query_log system.part_log
- insert 操作和 DDL 是否冲突
- 性能测试

设计与调研



设计与调研

状态机



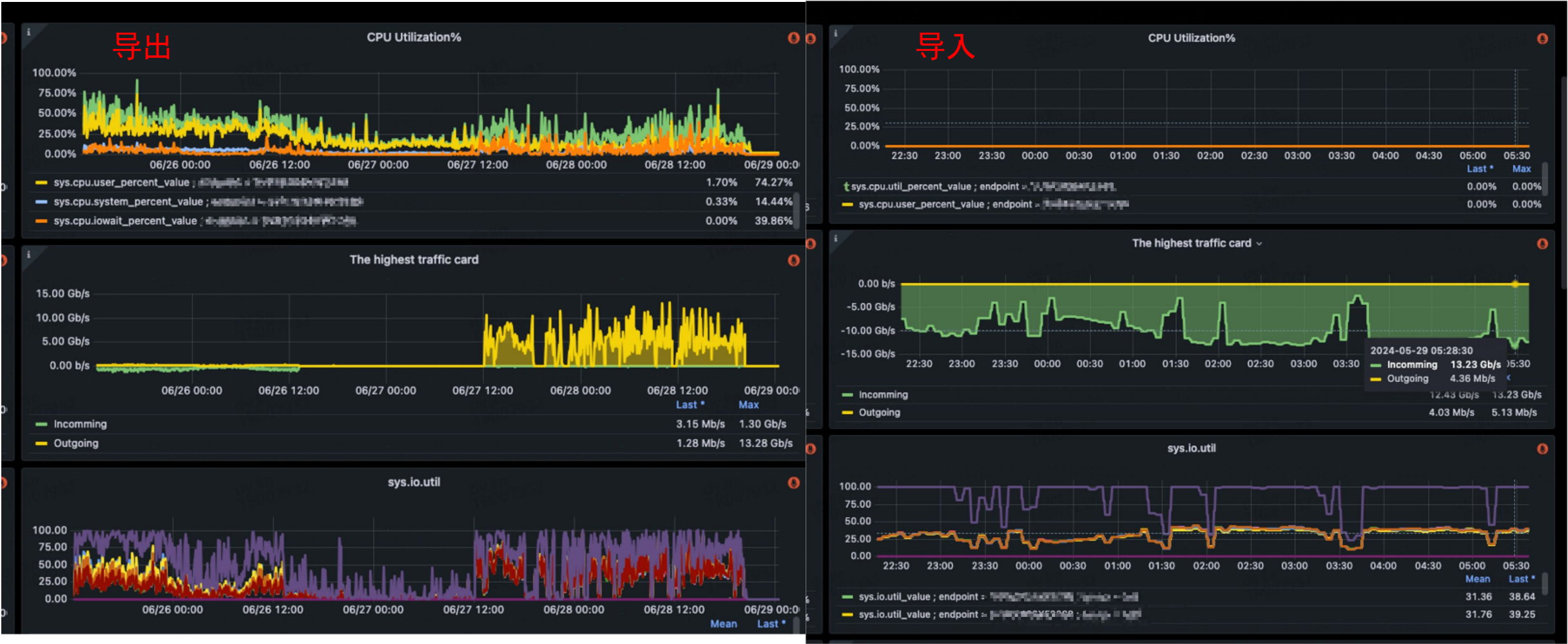
设计与调研

Fetch

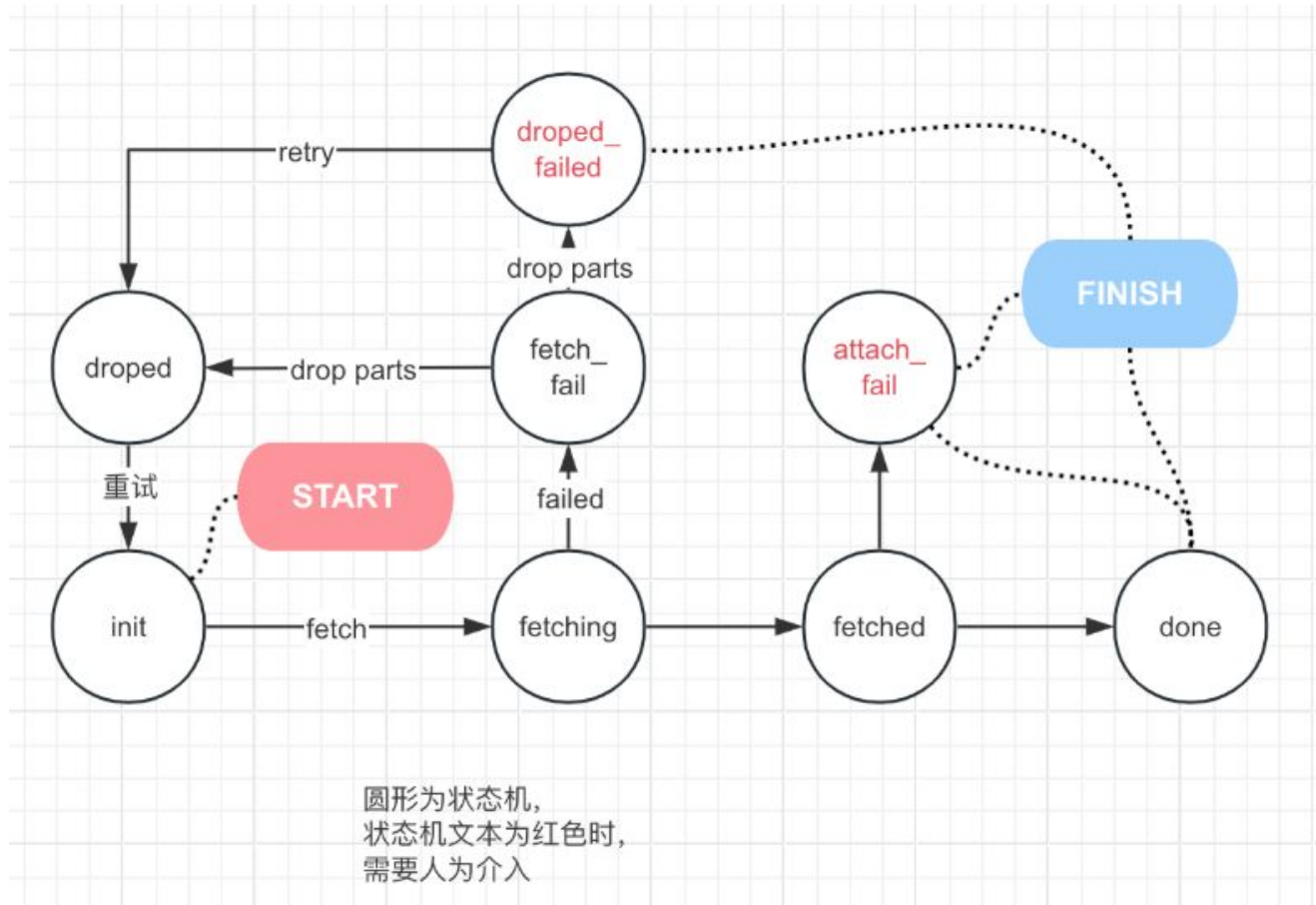
- FETCH FROM zkpath
 - 不同字段表是否能够fetch
- 压缩格式不同是否能够attach(批量改压缩)
- 不同版本数据兼容
- 字段不同能否attach成功(批量加字段)
- fetch 过程中对dst 表加字段后, 是否能正常attach(DDL)
- Attaching_part bug
- Inmemory part无法fetch
- Stop merge / ttl
- 动态auxiliary_zookeepers
- ALTER TABLE log.xx FETCH partition '20240212' FROM 'test_zookeeper:/clickhouse/tables/cluster1-shard1/tabletest'

性能与效率

成果：集群升级



设计与调研



设计与调研

查询兼容测试:

- allow_experimental_analyzer
- 特殊语法
- SQL扫描

1、**multif**、**case when**语句包含**in**的, 如下:

```
multif(  
  pageld in ('220236','401138'), 'h5',  
  pageld in ('103816','103817'), 'pc',  
  '未知')
```

修改后, 如下:

```
case when (pageld = '220236' or pageld = '401138') then 'h5'  
when (pageld = '103816' or pageld = '103817') then 'pc'  
else '未知'  
end
```

2、**if**语句包含**in**的, 如下:

```
count(DISTINCT if(page IN ('10650039006', '10650045146',  
'10650034951'), vid, NULL)) / count(DISTINCT if(page IN  
('10650040697', '10650039004', '109908'), vid, NULL)) AS countpage
```

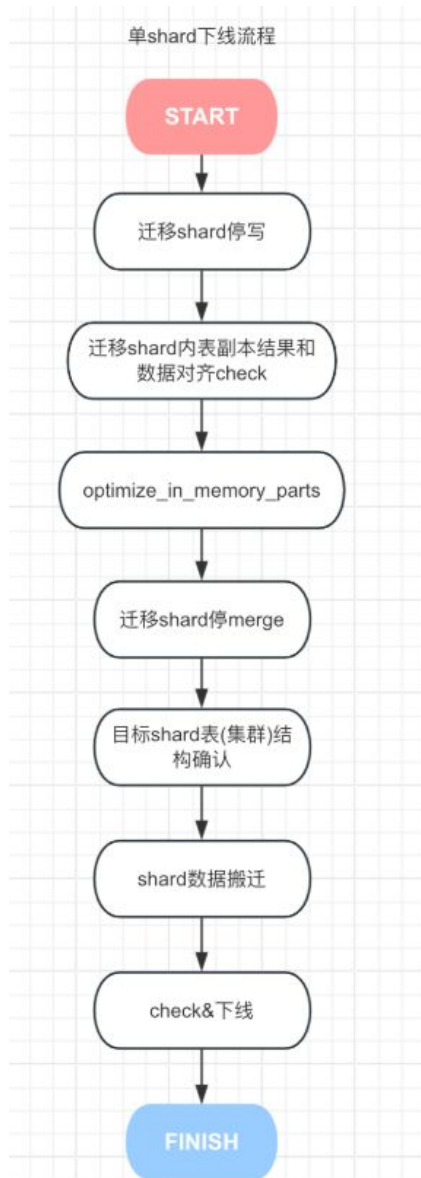
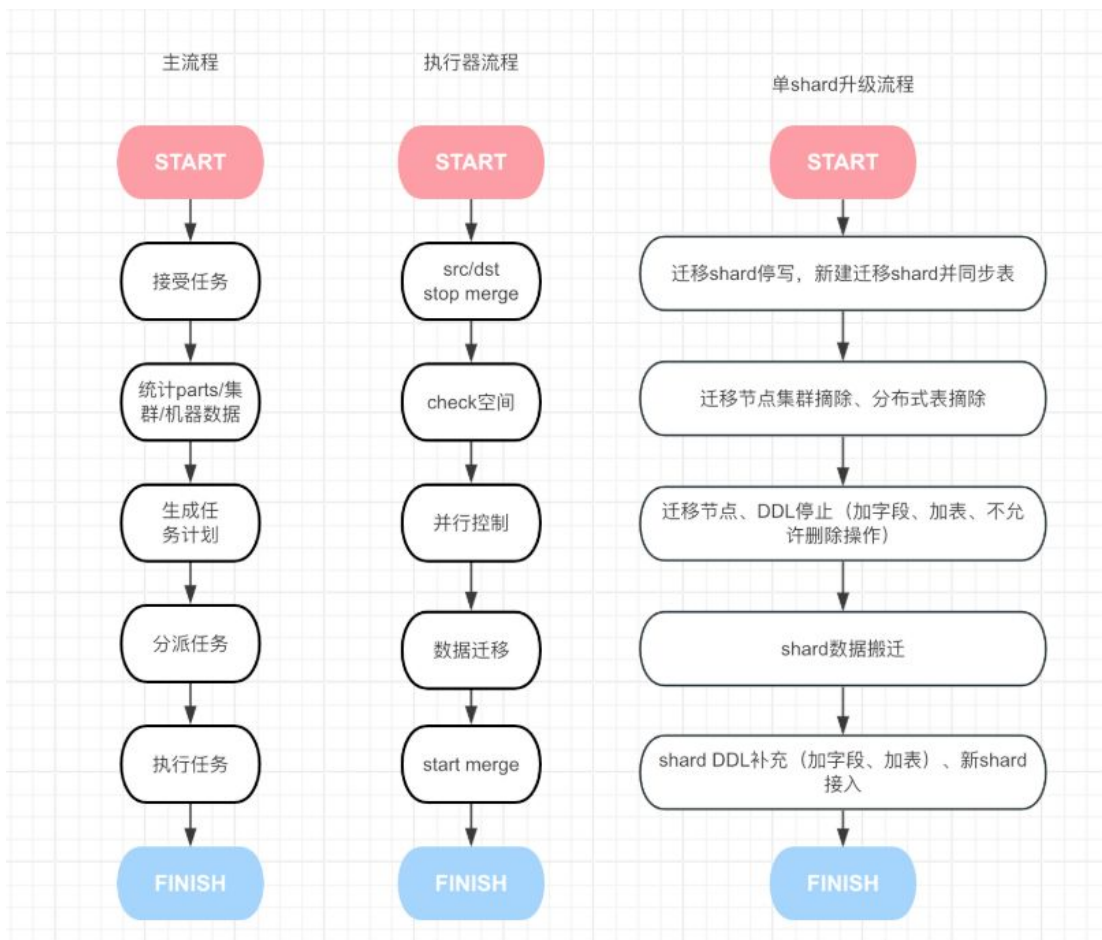
修改后, 如下:

```
count(DISTINCT case when (page ='10650039006' or page =  
'10650045146') then vid  
  when page ='10650034951' then vid  
  else NULL end) / count(DISTINCT case when (page  
='10650040697' or page = '10650039004' ) then vid  
  when page ='109908' then vid  
  else NULL end) AS countpage
```

注意: case when 中 or 、and不能超过2个, 超过的用when再分开

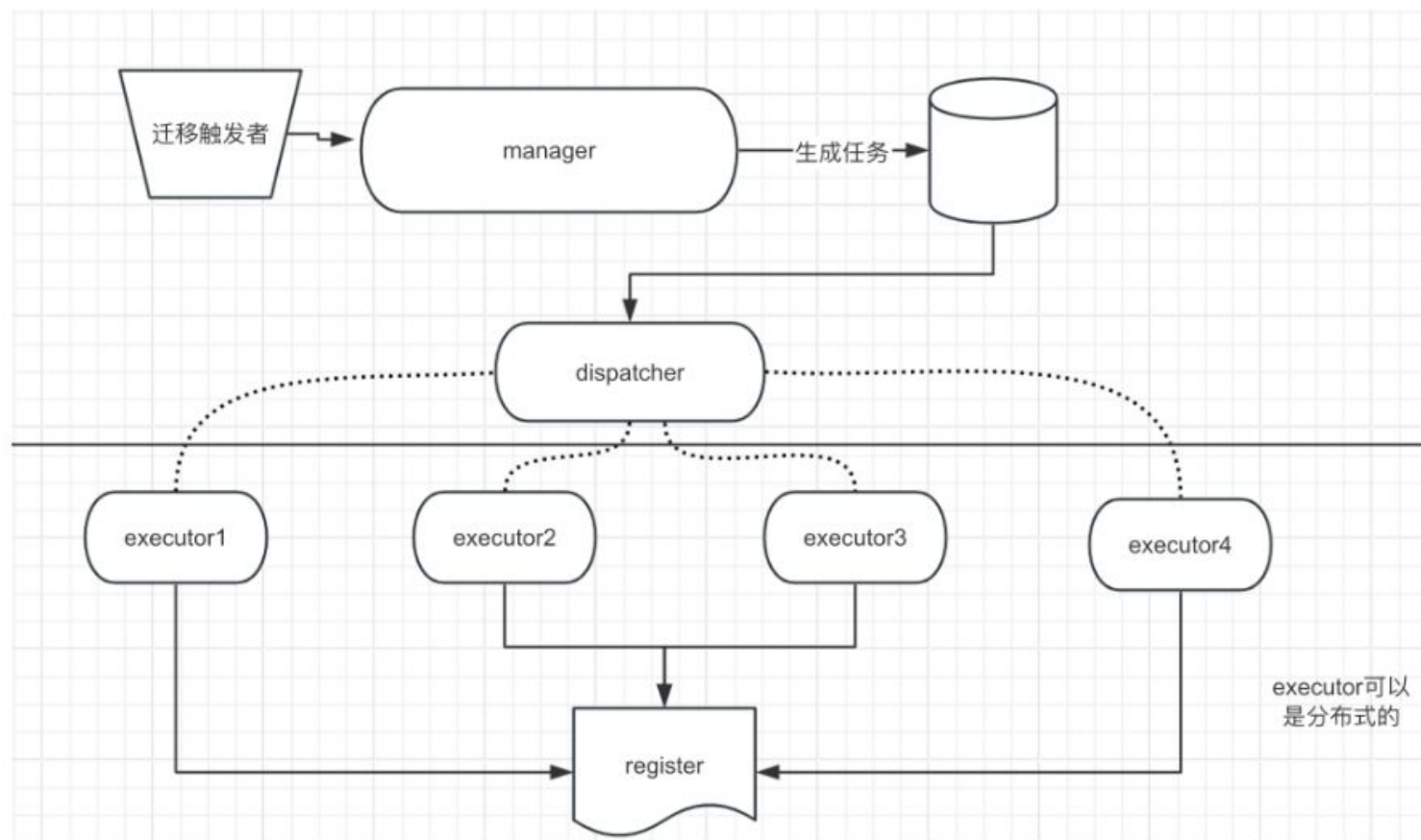
设计与调研

流程设计



设计与调研

架构设计



- 执行器与shard级别的任务绑定

成果:集群均衡与缩容



成果：集群升级

迁移：

涉及70+PB数据搬迁，每日搬迁上PB数据，过程用户无感

空间节省：

压缩比从4.3提升20%到5.56

稳定性：运维报障数(7-9月 TS技术支持统计)

91->67-> 35 (次)

