



ClickHouse在腾讯音乐敏捷数 据分析中的实践和思考

zelus

2021.02

生态全景



一体化的音乐娱乐，高参与度、强社交性和有趣的用户体验



目录

- 01 快速的业务迭代与激增的数据需求
- 02 数据平台架构实践
- 03 平台思考

ClickHouse 小调查

call 1

简单搭建使用，性能测试对比

call 2

落地生产环境大规模应用

call 3

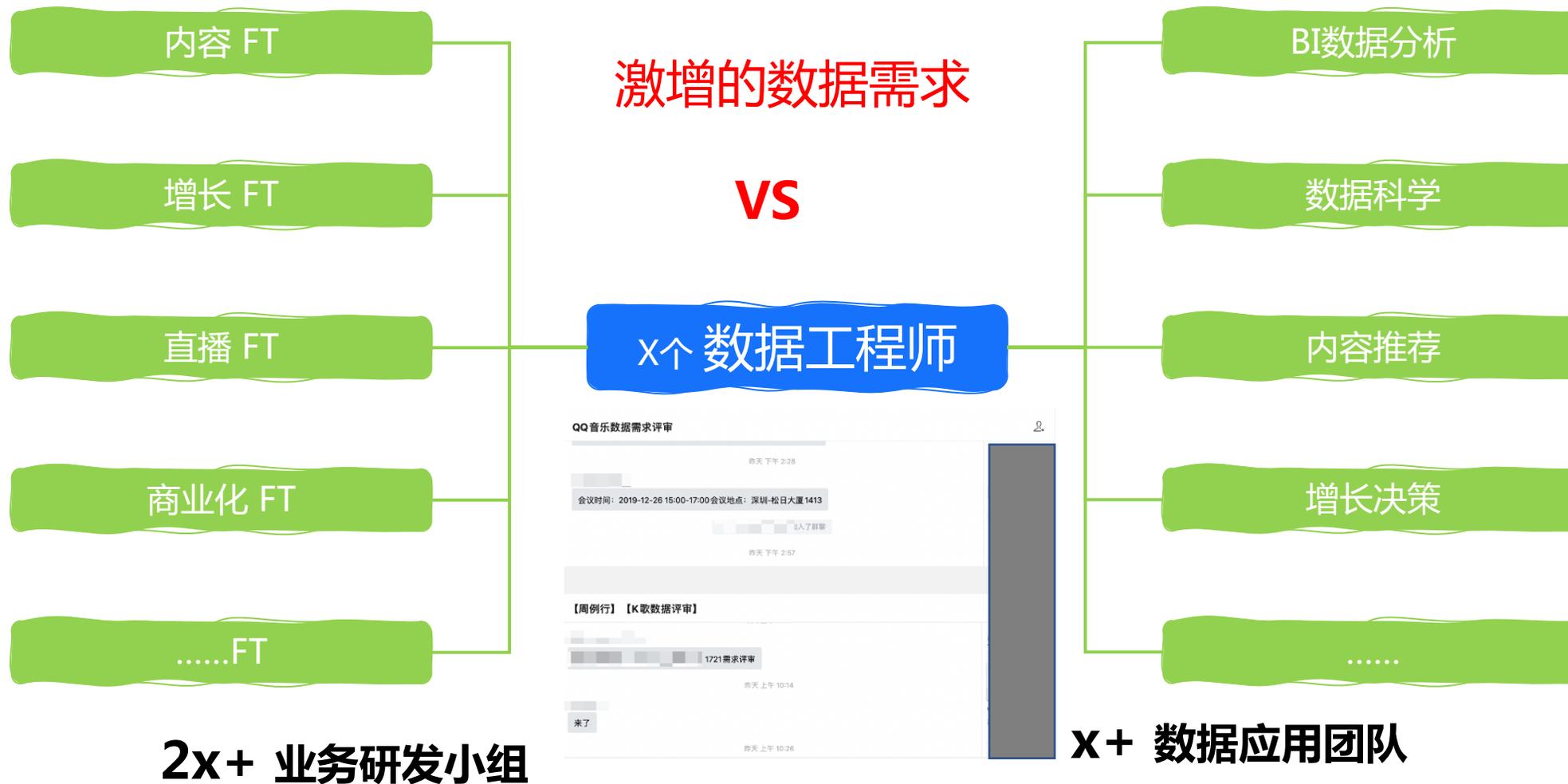
社区源码贡献者



01 快速的业务迭代与激增的数据需求



数据需求中心化处理瓶颈



02 数据平台演进实践



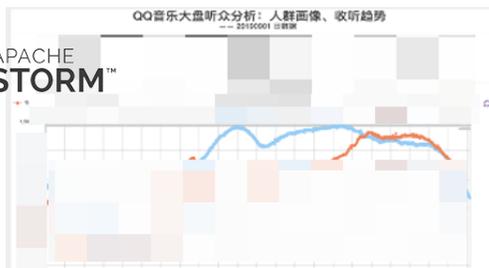
数据平台架构演进

1.0 BI数据分析

- 数据需求驱动，罗盘BI数据报表建设

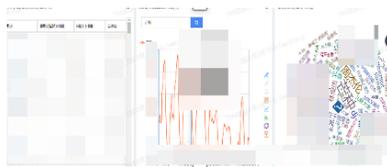


- 探索TRC (Storm) 实时计算



2.0 业务数据分析平台

- 大规模实时流式内容分析计算应用

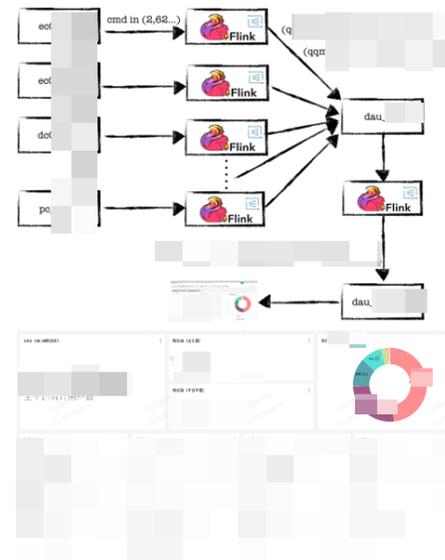


- Hermes画像在线分析千万级用户圈层
- Kylin预计算解决内容类多维度分析需求

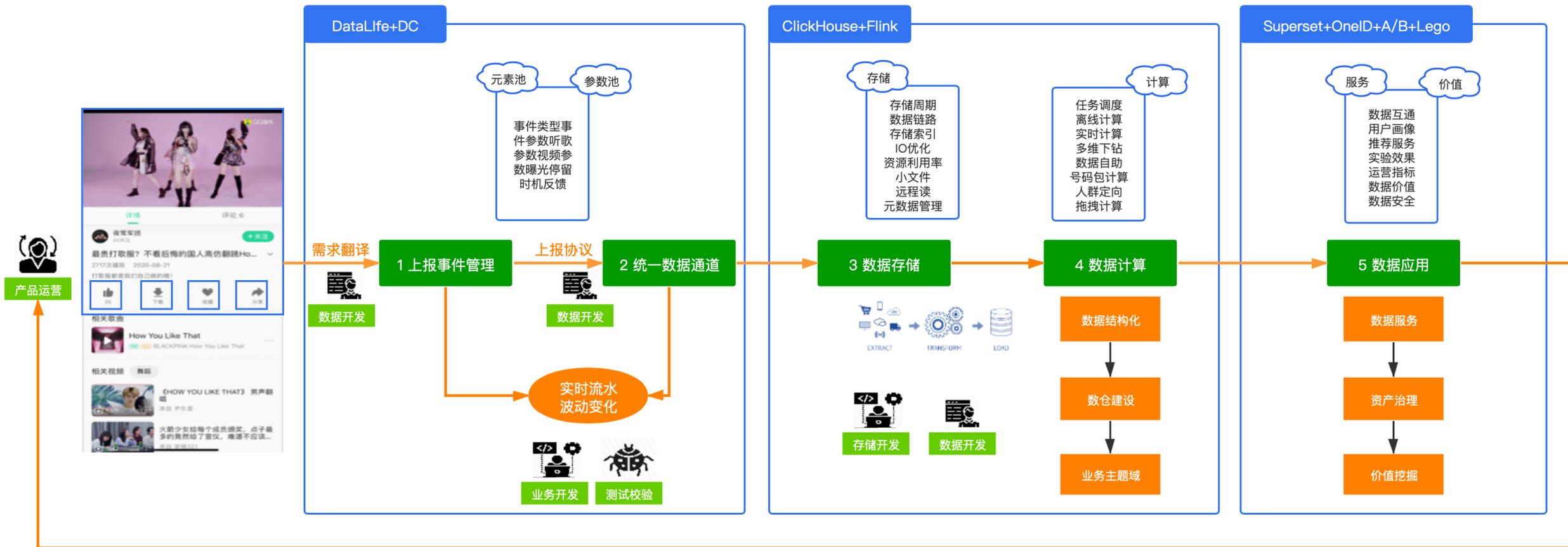


3.0 自助数据分析平台

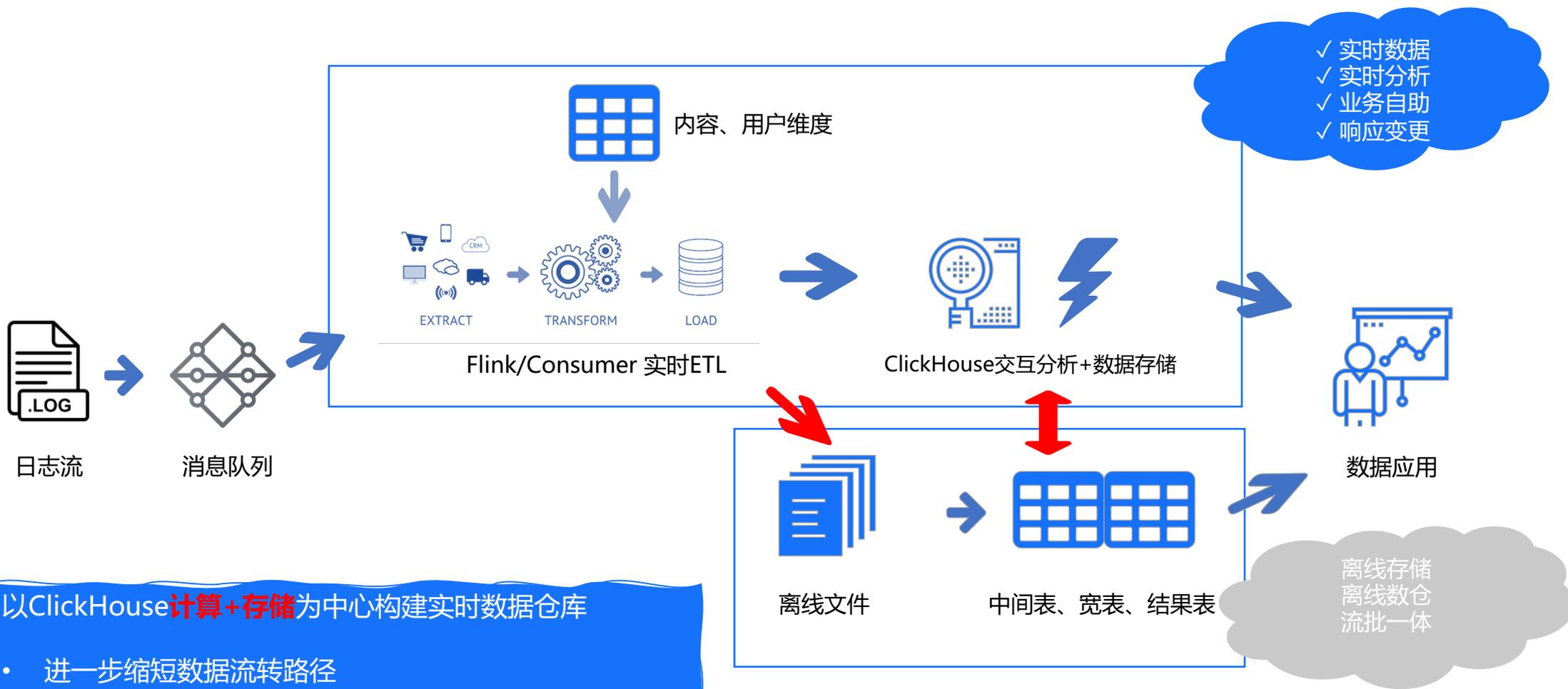
- Clickhouse同时满足大数据量即席查询与实时计算
- Superset一站式的数据探索+数据可视化平台
- Flink实时基础实时数仓建设



数据全链路生产加工



3.0交互式数据流批一体计算

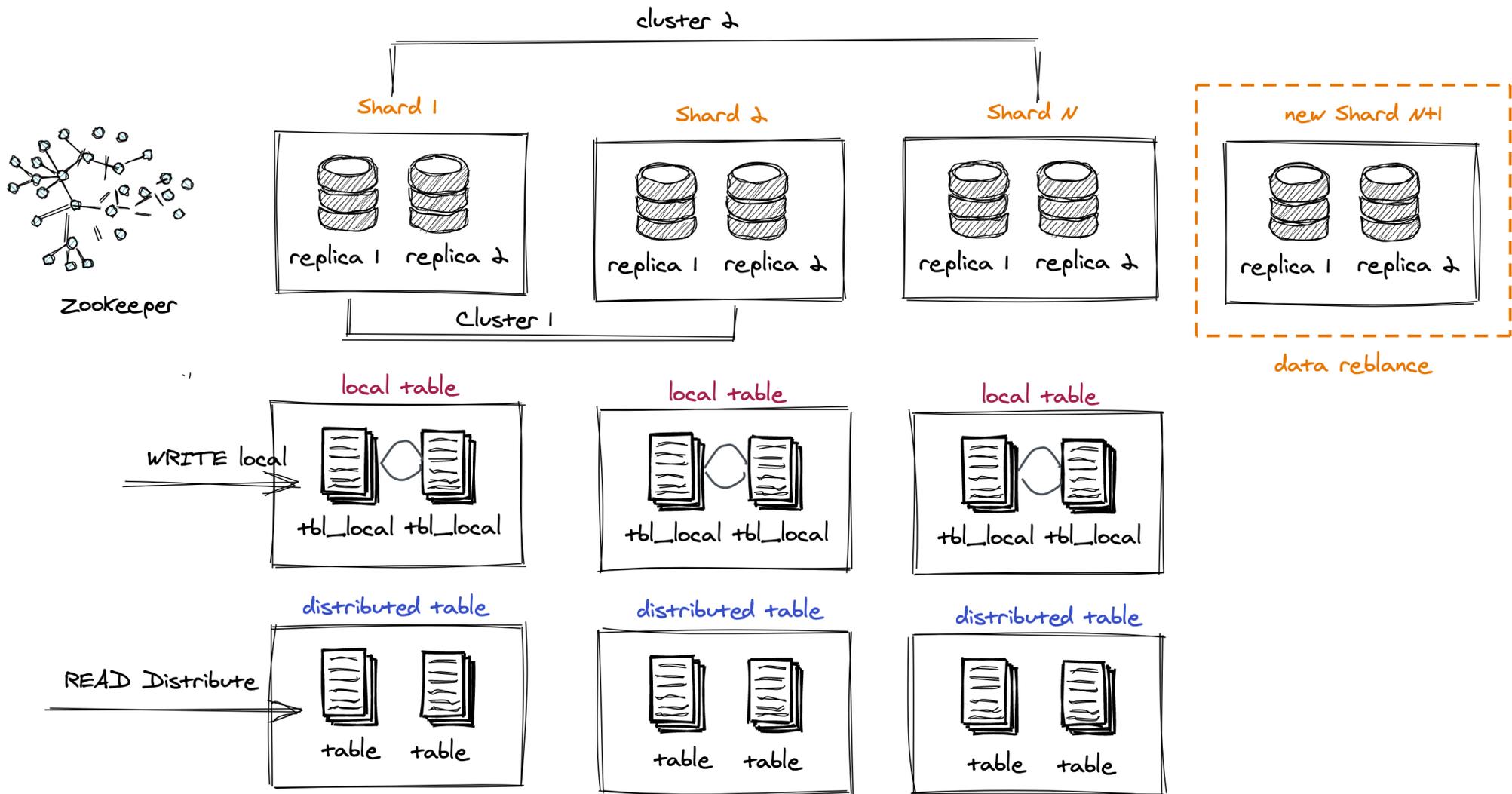


以ClickHouse**计算+存储**为中心构建实时数据仓库

- 进一步缩短数据流转路径
- 计算代码复用，流批计算一体化
- 离线大规模存储

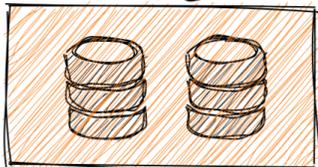


K8容器化多集群架构部署

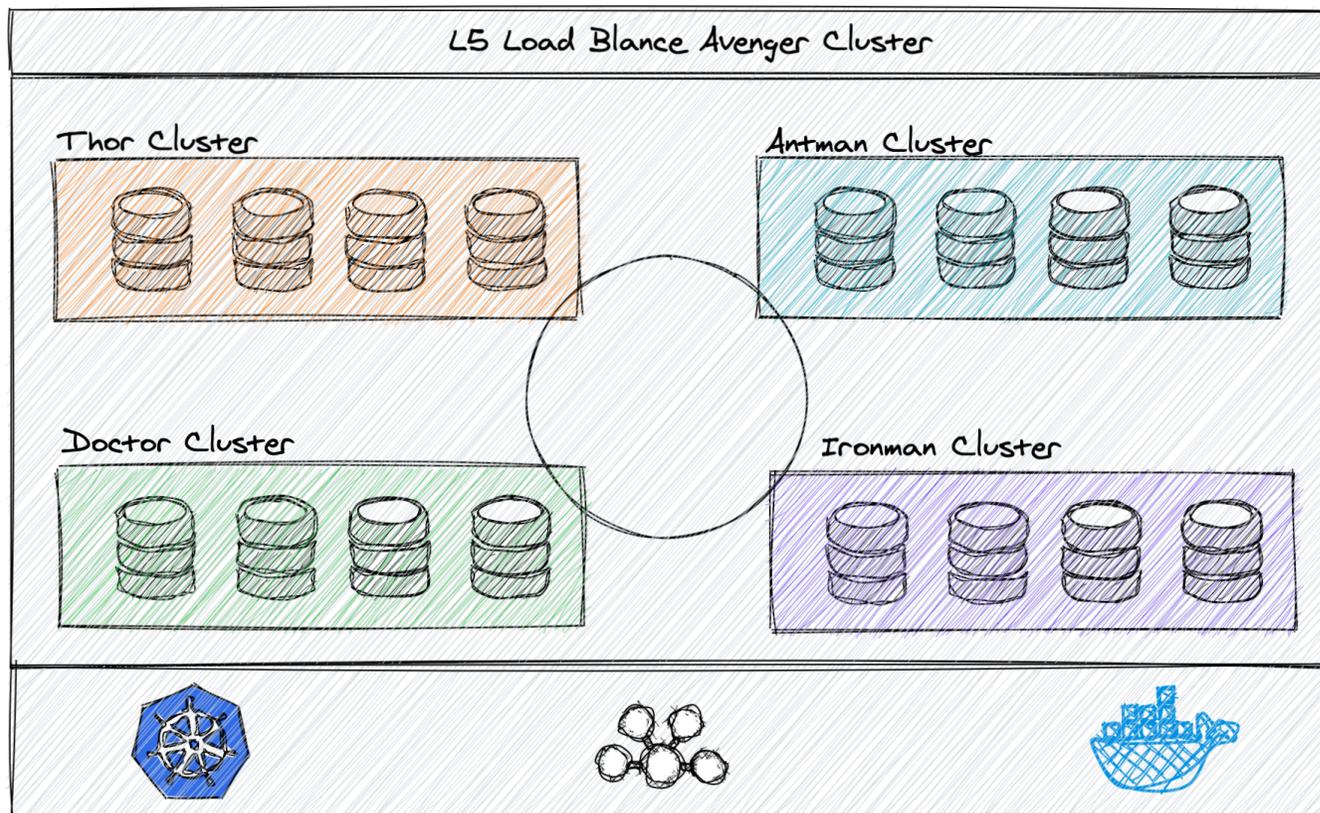
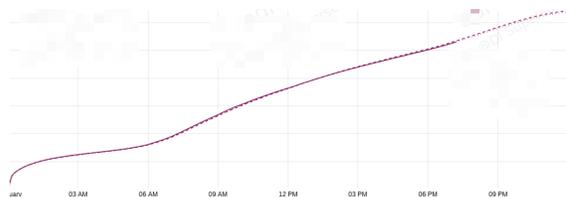


K8容器化多集群架构部署

Thor Cluster



- 秒级累计全量DAU计算
- 实时消费流水xxxxW/s
- 实时OLAP需求激增



1. 数据多分片
2. 多数据中心
3. 多集群相互可见
4. K8容器化管理
5. 状态存活检测
6. Server熔断
7. 资源限制
8. 负责均衡
9. 多维度监控



为什么用K8s

v1.0



clickhouse(Docker手动管理)

juno(bash手动管理)



原来如此!

熔断

资源使用
断/就绪检测

Log	Node-X	Ch-X	Event-X
日志	节点监控	ch监控	事件监控

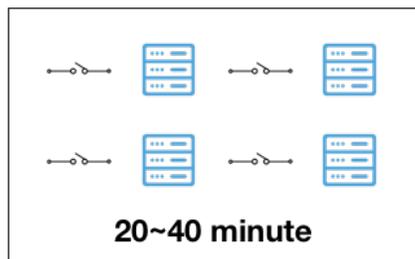
监控报警

工作重心转移

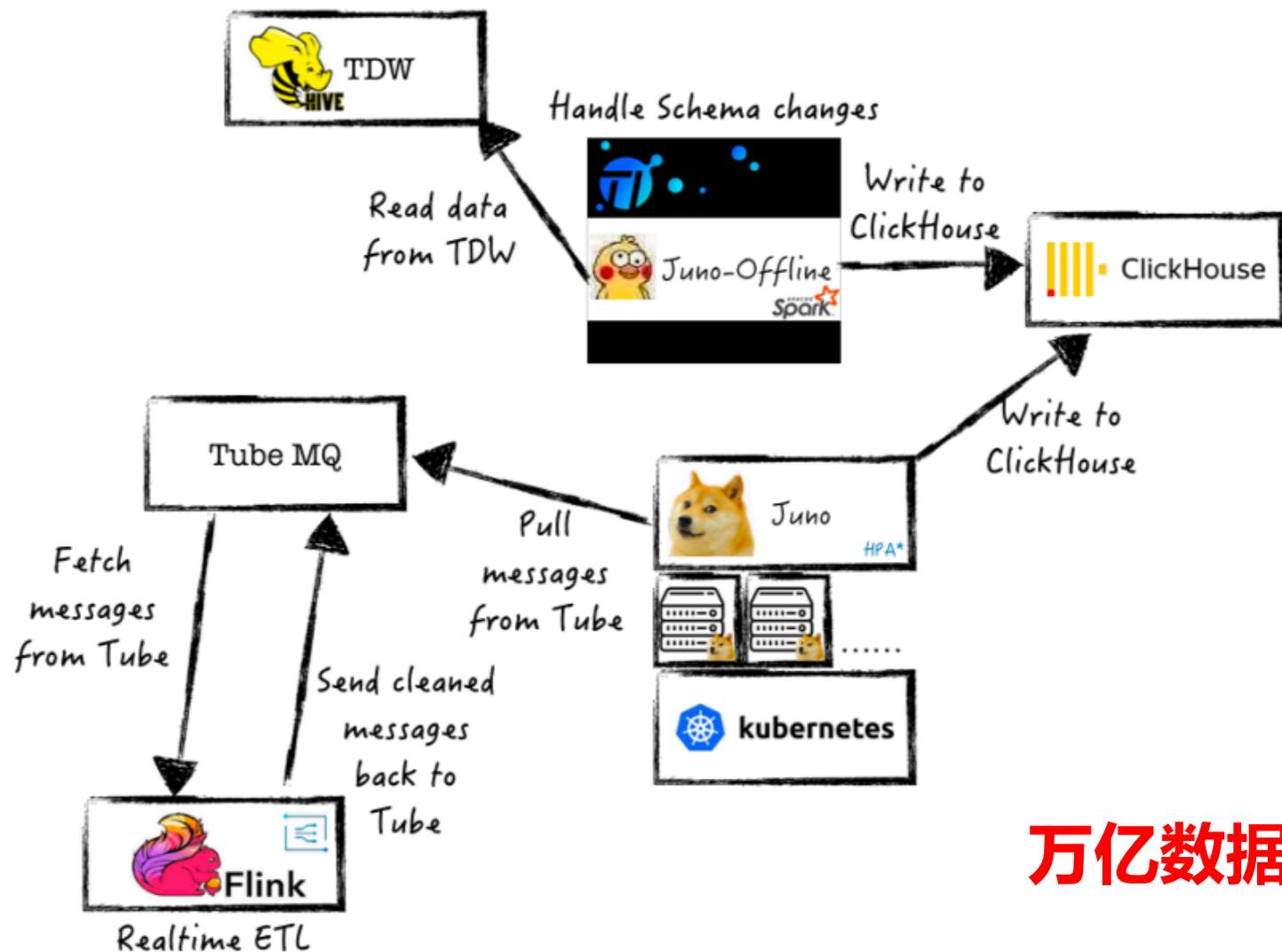


Clickhouse

自动化熔断修复



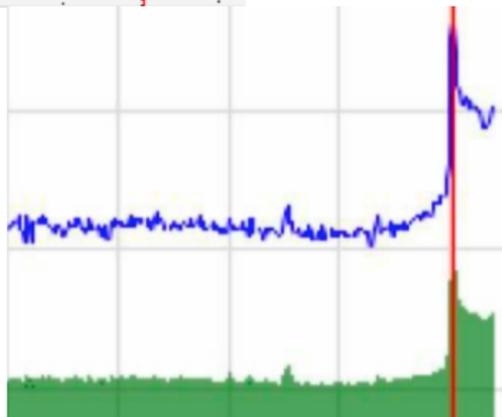
离线+实时 写入方案



万亿数据同步方案



HPA 弹性伸缩



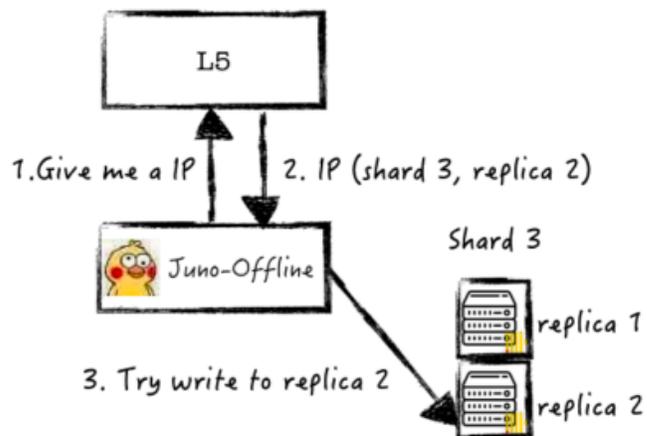
QQ音乐直播：陈奕迅《Live is so much better with music》慈善音乐会
07-11 17:00



TME live:五月天 突然好想见到你 mayday
2020 live in the sky 线上演唱会
05-31 20:00



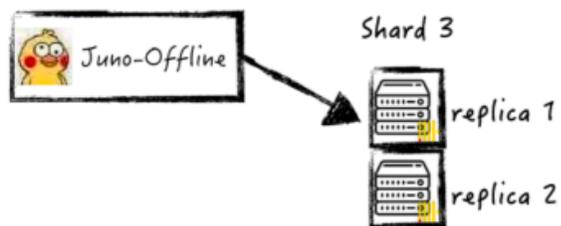
spark 离线写



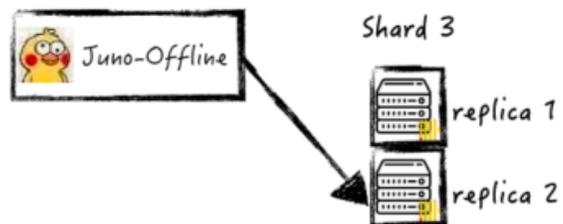
Each block of data is written atomically. The INSERT query is divided into blocks up to `max_insert_block_size = 1048576` rows. In other words, if the `INSERT` query has less than 1048576 rows, it is made atomically.

Data blocks are deduplicated. For multiple writes of the same data block (data blocks of the same size containing the same rows in the same order), the block is only written once. The reason for this is in case of network failures when the client application doesn't know if the data was written to the DB, so the `INSERT` query can simply be repeated. It doesn't matter which replica `INSERT`s were sent to with identical data. `INSERT`s are idempotent. Deduplication parameters are controlled by `merge_tree` server settings.

4. If 3 failed, try write to replica 1



5. If 4 failed, wait for 10 mins and back to 3



写性能瓶颈

1、分区数过多

Info

A merge only works for data parts that have the same value for the partitioning expression. This means you shouldn't make overly granular partitions (more than about a thousand partitions). Otherwise, the **SELECT** query performs poorly because of an unreasonably large number of files in the file system and open file descriptors.

小时分区 => 天分区

2、zk同步瓶颈

SSD 提升IO读写性能，JVM资源调优

拆分集群

多zk observer 节点

不同表指向不同zk集群，auxiliary_zookeepers

自定义同步环

3、大表批量，网卡打满，Merge Parts太慢

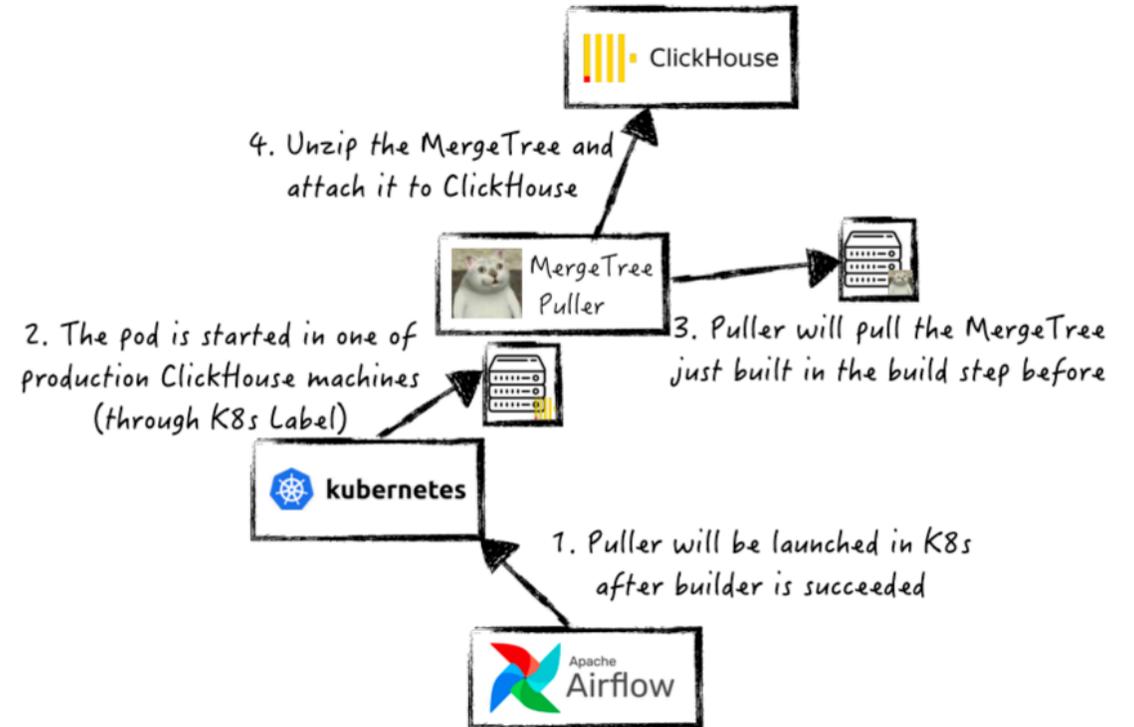
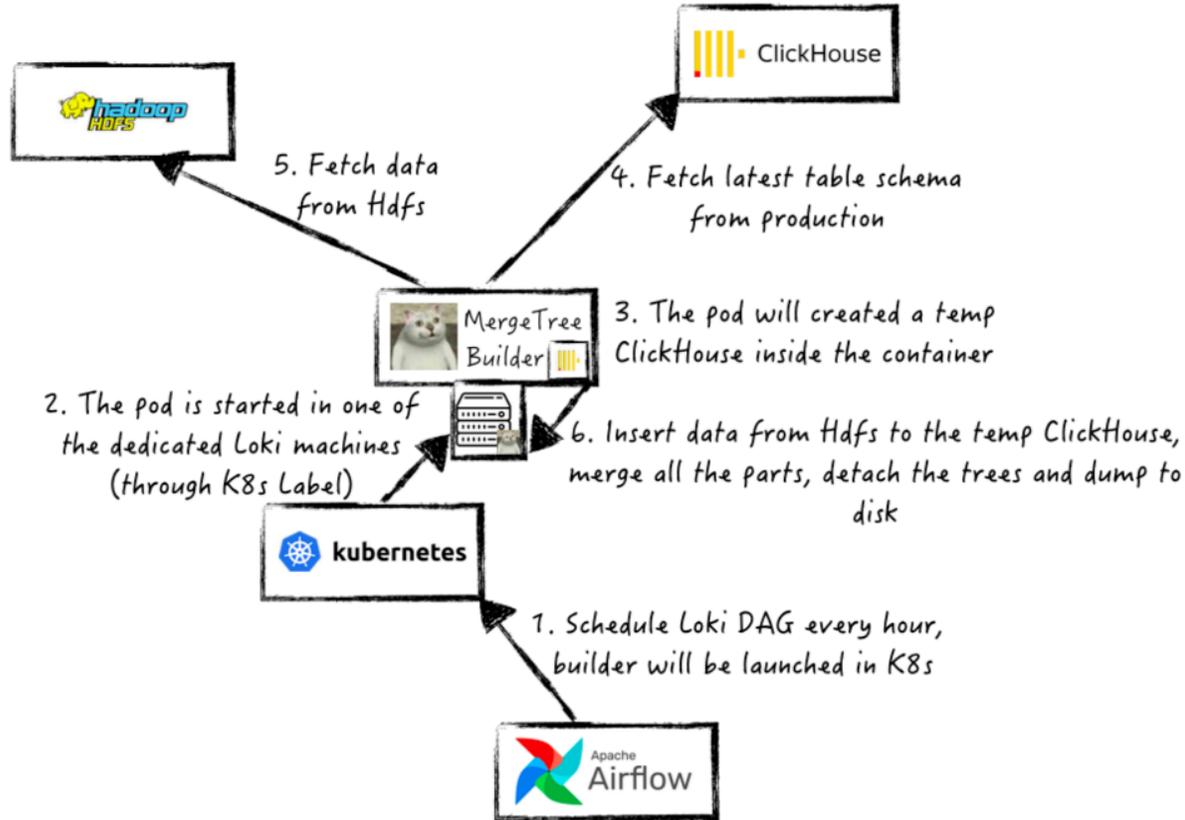
读写分离方案



Timestamp	Count
2020-03-18 11:27:11	count:166
2020-03-18 11:29:38	count:178
2020-03-18 11:32:04	count:193
2020-03-18 11:34:31	count:208
2020-03-18 11:36:57	count:221
2020-03-18 11:39:24	count:236
2020-03-18 11:41:50	count:249
2020-03-18 11:44:17	count:266
2020-03-18 11:46:43	count:278



读写分离

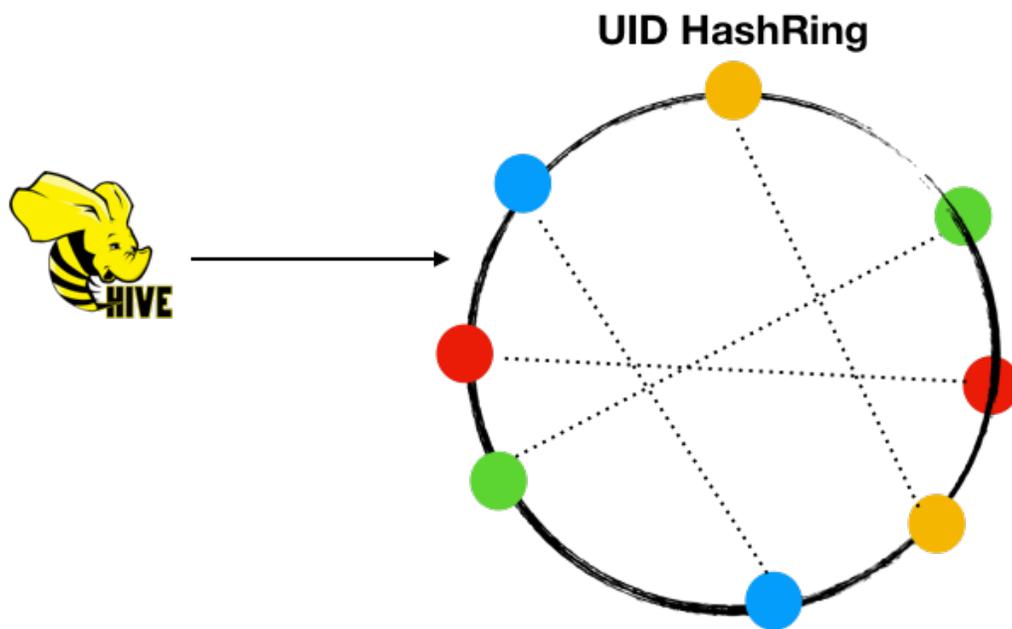


进一步缩短数据导入链路

老链路



新链路

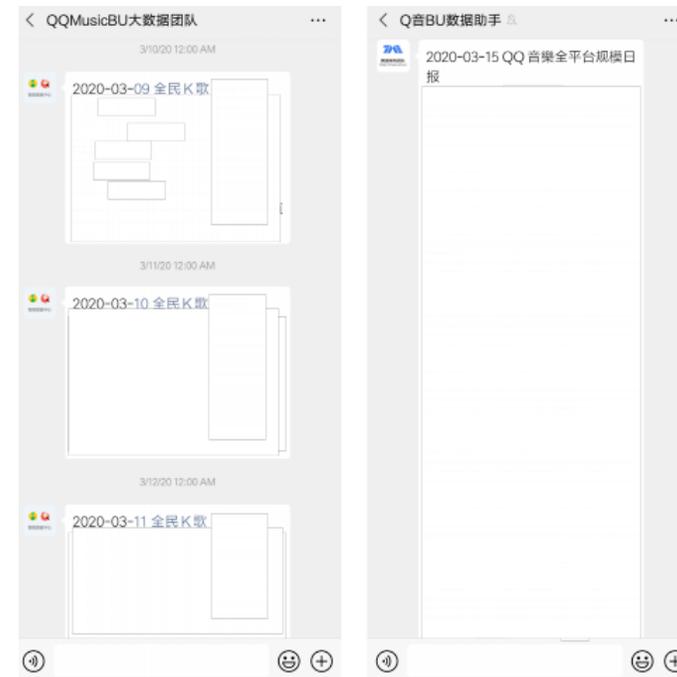
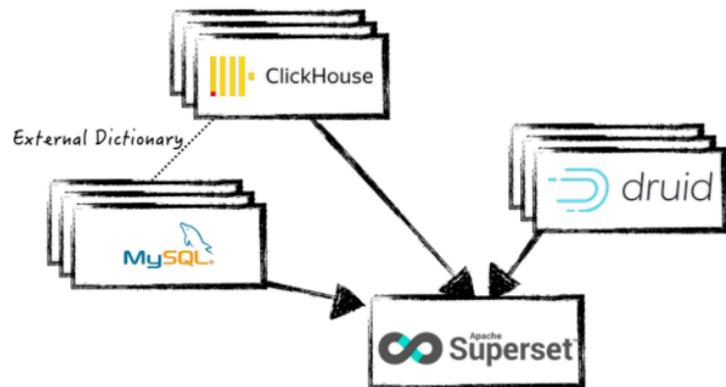


1. 一致性Hash算法，按用户uid分成1024桶
2. 同一用户行为数据落在同一个节点，提升用户历史行为聚合性能
3. 一致性hash架构，加速集群扩容和数据迁移
4. 写入端完成数据双写，不依赖zookeeper

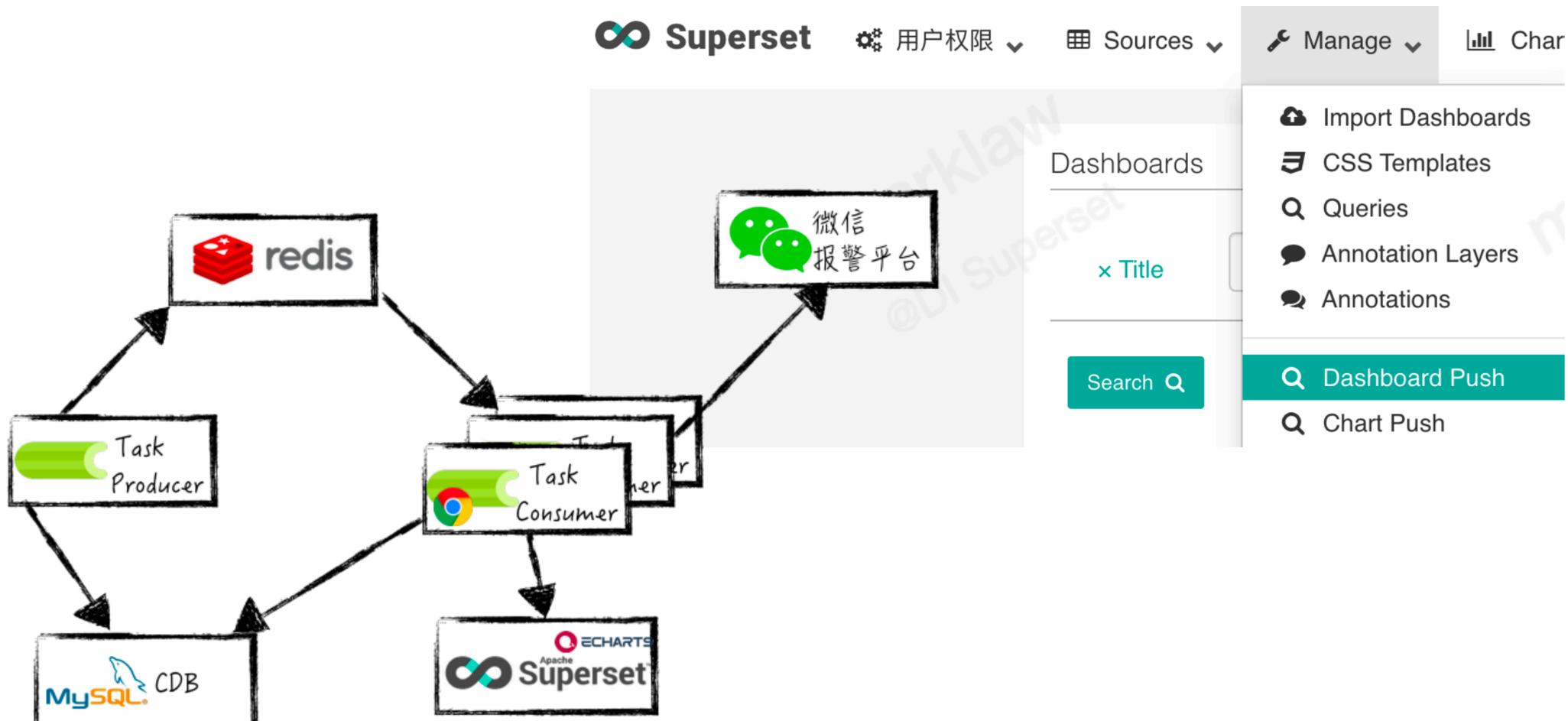
<https://github.com/ClickHouse/ClickHouse/issues/16798>



敏捷自助数据分析



数据订阅推送



03 平台思考



云原生存储计算分流

data
Service



query
processing



data
Storage



Local Machine 96c256g 30T

```
CREATE TABLE lineorderflat_multithread_insert_bx2  
ENGINE = MergeTree  
PARTITION BY toYear(LOORDERDATE)  
ORDER BY (LOORDERDATE, LOORDERKEY) AS  
SELECT *  
FROM lineorderflat  
SETTINGS max_insert_threads = 60
```

Ok.

0 rows in set. Elapsed: 1605.129 sec. Processed 600.04 million rows, 140.52 GB (373.83 thousand rows/s., 87.55 MB/s.)

Cloud CFS 3 shard 1 replication

```
0 rows in set. Elapsed: 70.678 sec. Processed 600.04 million rows, 140.41 GB (8.49 million rows/s., 1.99 GB/s.)
```

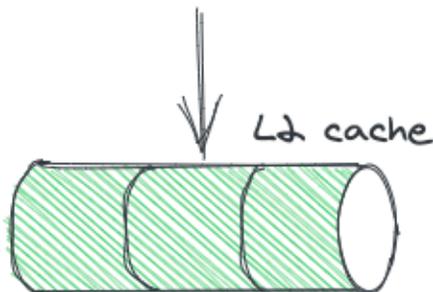
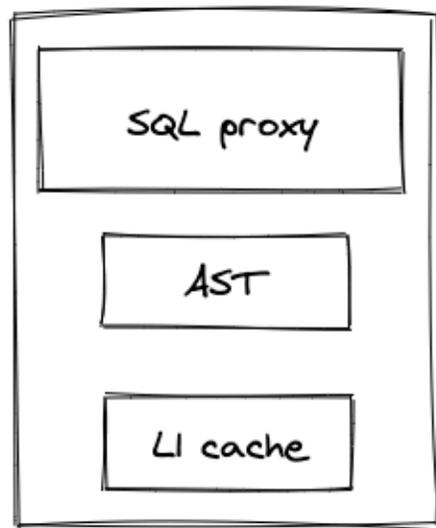
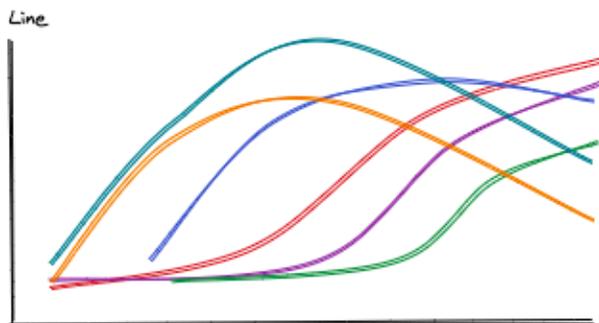
```
0 rows in set. Elapsed: 71.301 sec. Processed 600.04 million rows, 140.41 GB (8.42 million rows/s., 1.97 GB/s.)
```



Cache加速



```
SELECT toStartOfDay(toDateTime(Ptime)) AS __timestamp,  
max(expo_first_uv) AS first_expo_uv,  
max(expo_uv) AS Peak_expo,  
max(expo_s_uv) AS Fifth_expo_uv,  
max(play_first_uv) AS first_play_uv,  
max(complete_play_first_uv) AS complete_play_first_uv_sun  
FROM dtable1  
WHERE Ptime >= toDate('2021-01-03')  
AND Ptime <= toDate('2021-02-03')  
GROUP BY toStartOfDay(toDateTime(Ptime))  
ORDER BY first_expo_uv DESC
```



ClickHouse
< 30 day



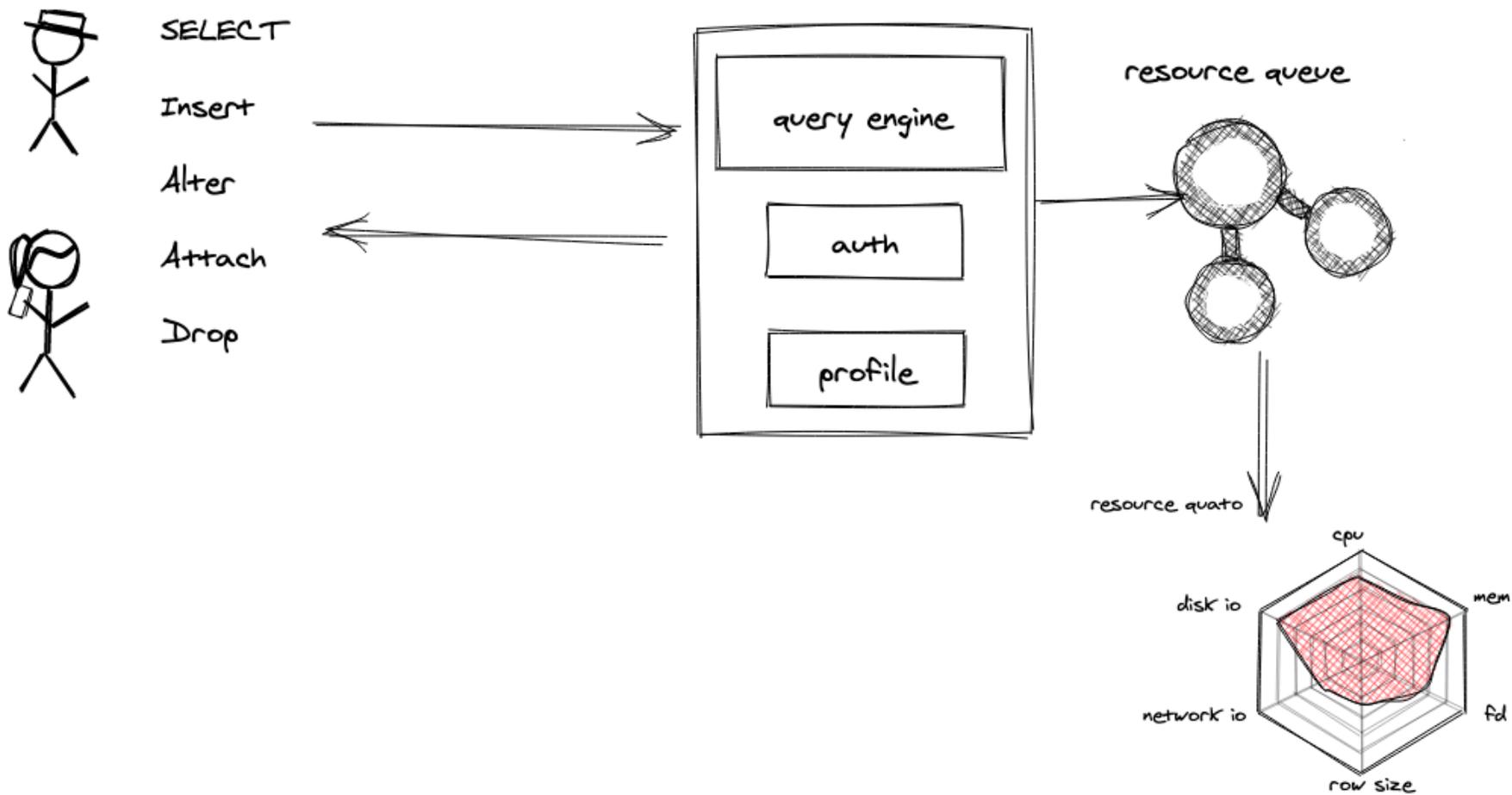
druid
30~90 day



spark sql
>90 day



多用户资源隔离



平台演进方法论

平台建设

数据
时效

+

数据
质量

=

价值
倍增

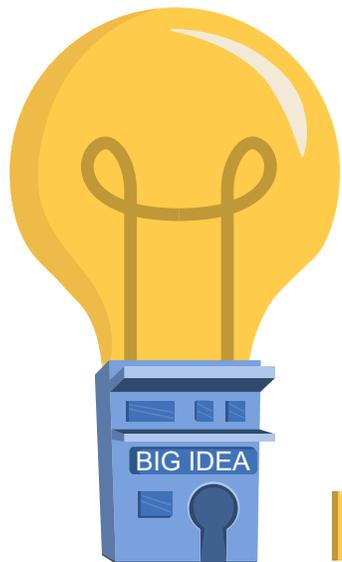
组织保障

高效
组织

+

有效
方法

- 平台的演进贴合业务并赋能业务
- 平台的演进兼顾多方协同共享共建
- 平台的演进促进组织架构优化



More Heros



欢迎自荐/推荐各类大数据人才

创造音乐无限可能

CREATING ENDLESS
OPPORTUNITIES WITH MUSIC