

ClickHouse在携程PB级日志系统实践

携程 林东煜

2023.03.25

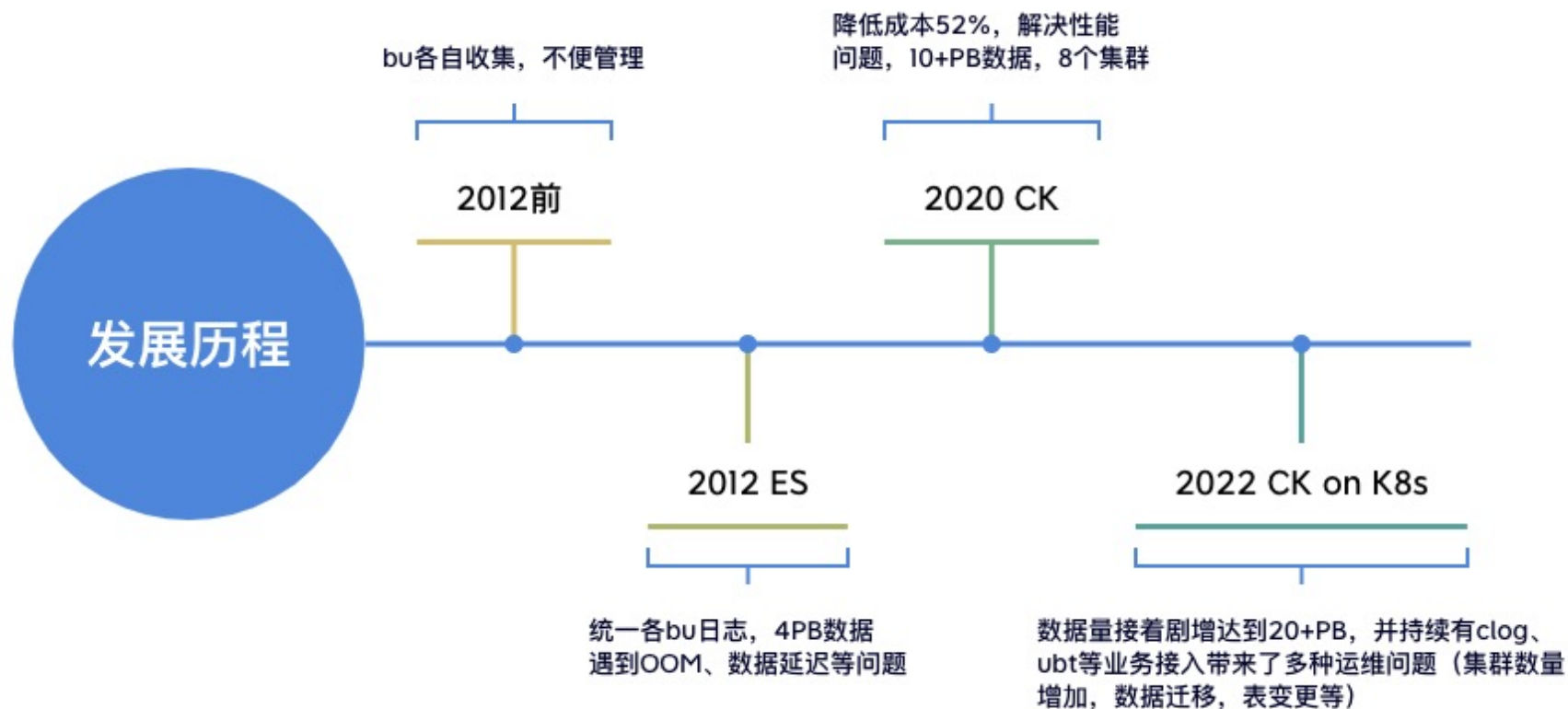
CONTENTS

目录

- 1 历史与现状
- 2 日志系统实践
- 3 痛点与解决方案
- 4 未来展望

历史与现状

日志平台发展历程



集群现状

服务器：600+

数据量：60w亿行

数据表：1w+

存储空间：30+PB

查询：10M+/天

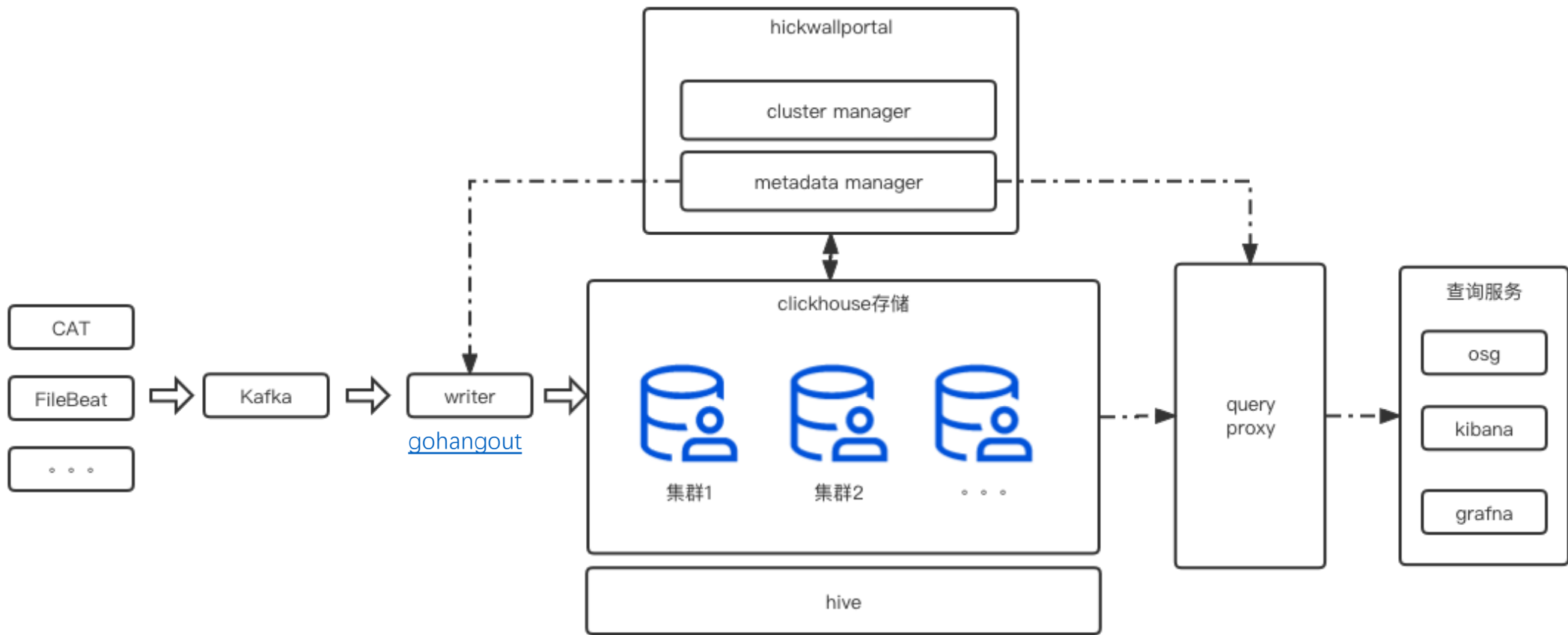
集群数量：20+

P90：500ms

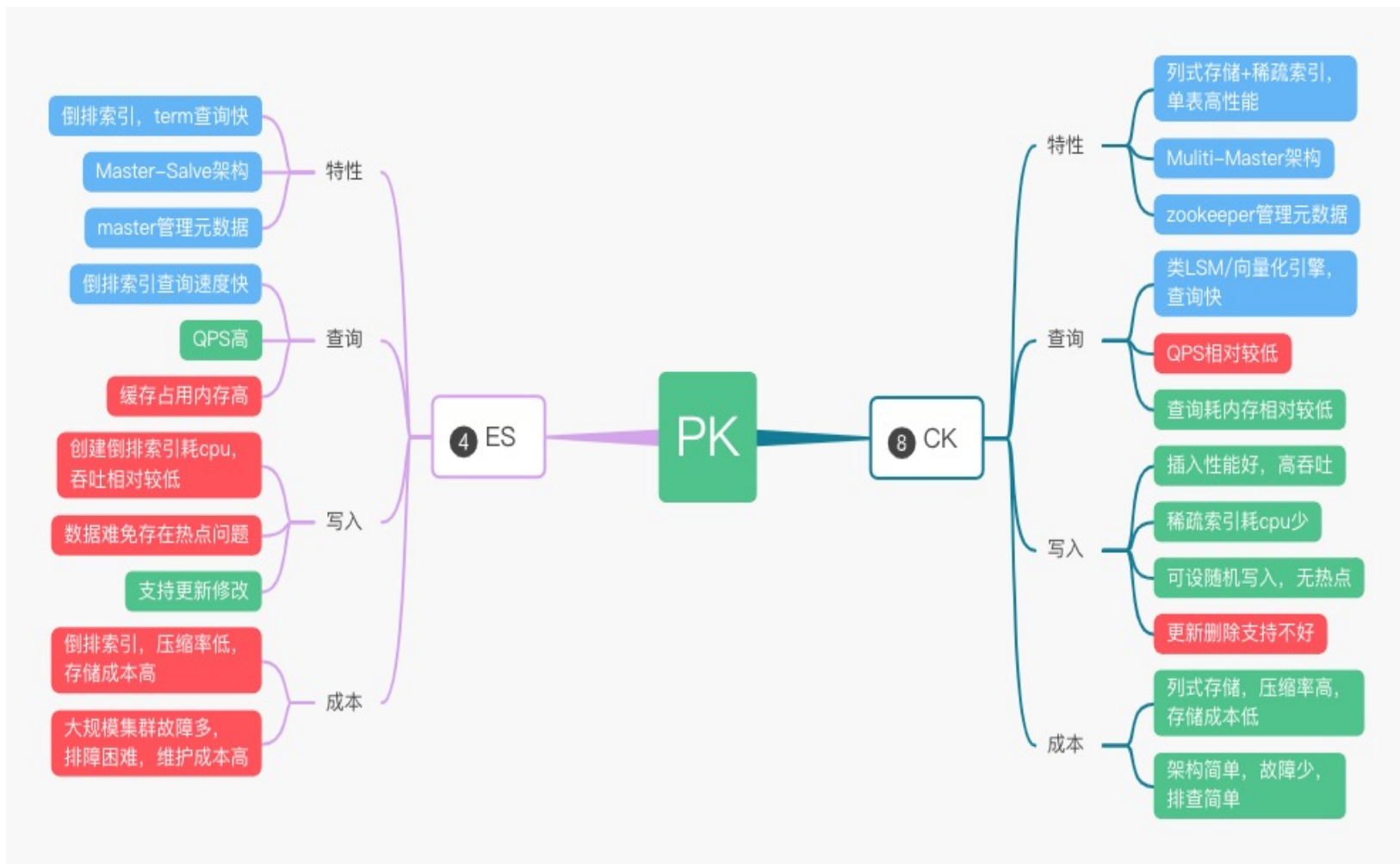
P99：3s

日志系统实践

整体部署架构



ClickHouse与ElasticSearch对比



库表设计

```

CREATE TABLE log.xxxx (
  `timestamp` DateTime,
  `_log_increment_id` Int64,
  `host_name` LowCardinality(String),
  `log_level` LowCardinality(String),
  `message` String,
  `message_prefix` String MATERIALIZED substring(message, 1, 256),
  `_tag_keys_str` Array(String),
  `_tag_vals_str` Array(String),
  `_tag_keys_float` Array(Float64),
  .....
  INDEX idx_message_prefix message_prefix TYPE tokenbf_v1(8192,2,0) GRANULARITY 16,
  .....
)
ENGINE = ReplicatedMergeTree('/clickhouse/tables/{shard}/xxxx', '{replica}')
PARTITION BY toYYYYMMDD(timestamp)
ORDER BY(timestamp, _log_increment_id, host_name, log_level)
TTL timestamp + toIntervalHour(168)
SETTINGS index_granularity = 8192

```

1. 双list存tags
2. 按天分区，时间排序
3. Tokenbf_v1优化term查询
4. 全局唯一递增id(翻页，明细)
5. ZSTD

集群管理平台

集群部署

集群管理

节点管理

DDL管理

监控告警

用户管理

集群管理

Clickhouse管理

IDC管理

Node管理

日志存档

日志迁移

集群管理

数据库

表结构

账户

数据管道

DDL

[cluster]

cluster

> [shard] 5e713d98-92d7-4a6e-8e36-83ea5cb4838e (RBLOG) (2 nodes;1)

> [shard] 2d9404d4-37bb-473f-b0f7-6f3b4cdb50b1 (RBLOG) (2 nodes;1)

▼ [shard] b1a84fb7-85e9-43fd-b109-d843a36e9a1d (RBLOG) (2 nodes;1)

修改

绑定

删除

- [node] (SV) /93TB/2021-12-31)

解绑

- [node] /93TB/2022-03-03)

> [shard] df427648-ed21-4707-a420-f3644b427f79 (RBLOG) (2 nodes;1)

▼ [shard] 63b5062f-1570-4a8b-9e3f-42ed4076c574 (RBLOG) (2 nodes;1)

- [node] /2022-01-30)

- [node] /2022-12-08)

> [shard] a16bfe1d-7ac0-42a3-ad1d-9e85423d043d (RBLOG) (2 nodes;1)

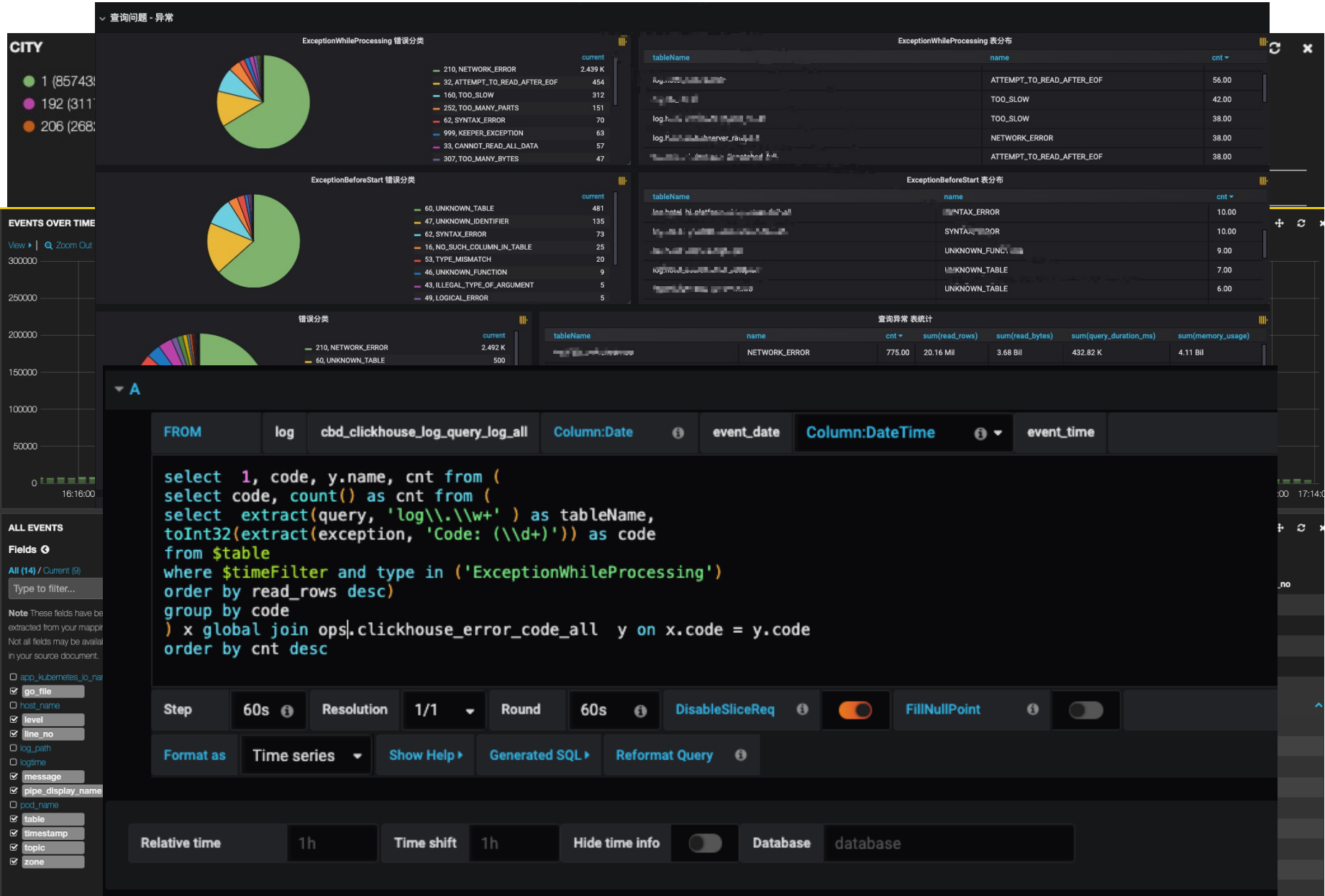
> [shard] a1355f15-ed4b-4fd5-adcd-9f7df4afe619 (RBLOG) (2 nodes;1)

> [shard] 63179837-9320-4b64-a240-fe306f5ea514 (RBLOG) (2 nodes;1)

> [shard] 8319dd7d-ac70-433f-8de1-6ad9977175b0 (RBLOG) (2 nodes;1)

> [shard] 439ab92c-6407-4775-8d01-c8fca5126512 (RBLOG) (2 nodes;1)

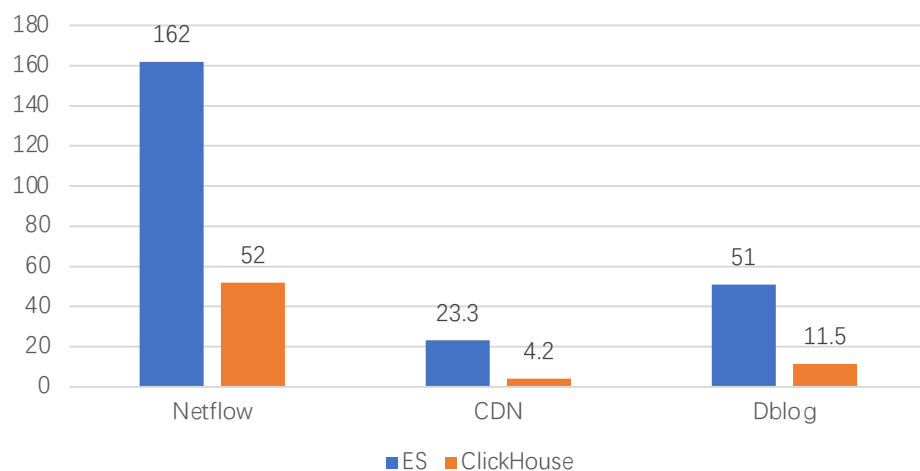
> [shard] ec90dbb6-ec43-4710-9473-523443a4c170 (RBLOG) (2 nodes;1)



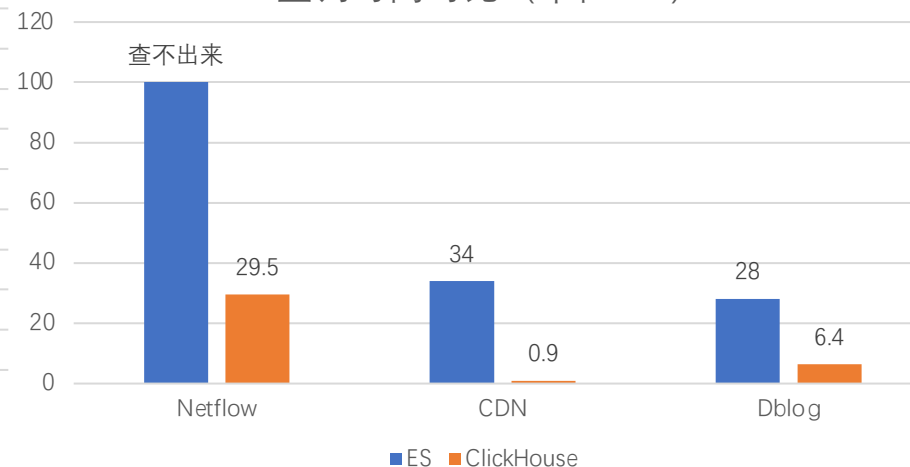
成果

- 迁移过程自动化程度超过95%
- 基本实现对用户透明
- 存储空间节省50+%，用原有ES的服务器支撑了4倍业务量的增长
- 查询速度比ES快4~30倍，查询P90小于300ms，P99小于1.5s

ES vs ClickHouse 磁盘空间对比（单位：T）



ES vs ClickHouse 查询时间对比（单位：S）



痛点与解决方案

两大痛点

性能与功能痛点：

单集群规模太大，
zookeeper性能瓶颈DDL超
时异常

删字段，字段数据量大时，超
时，容易导致元数据不一致

用户索引设置不佳，查询慢
重建会导致删除历史数据

查询层缺少限流、防呆、自
动优化，导致查询不稳定

运维痛点：

表与集群严格绑定，集群磁
盘满后，只能通过双写迁移

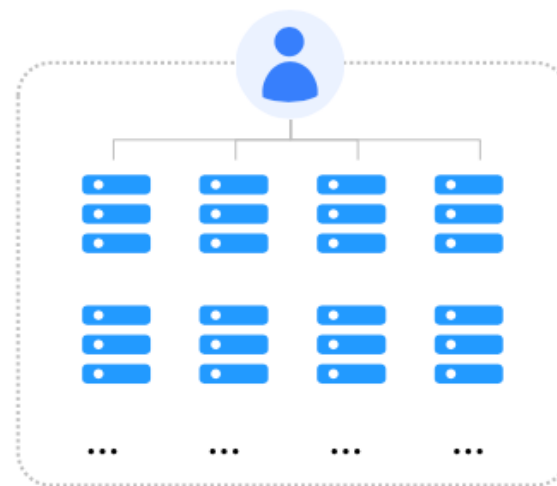
集群搭建依赖ansible，部署周期
长（数小时）

与社区版本脱节，目前部署模
式不便版本更新

ClickHouse on K8s

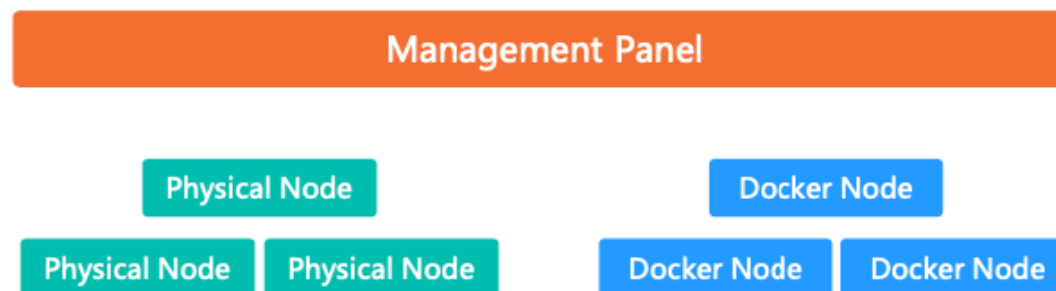
单集群规模太大，
zookeeper性能瓶颈DDL超
时异常

- 期望收益：
 - 提高资源利用率
 - 降低运维成本，释放人力



(图片来源于网络)

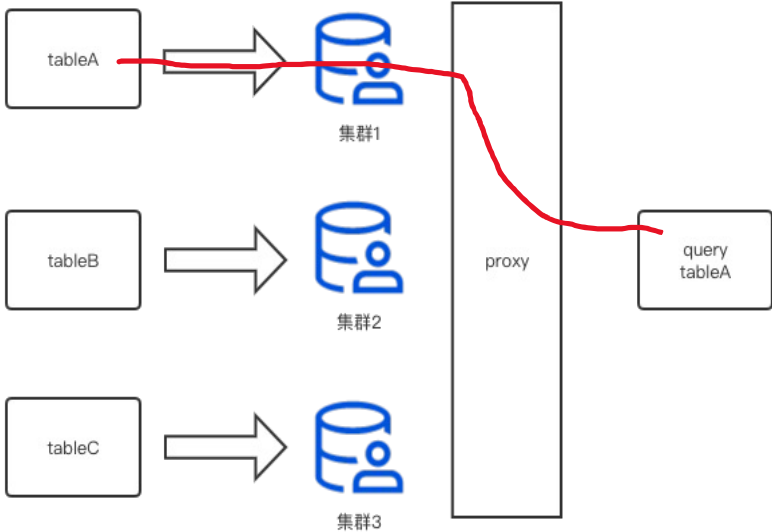
- 思路与方法：
 - sts、cm
 - 亲和性
 - topolvm



如何解决数据跨集群

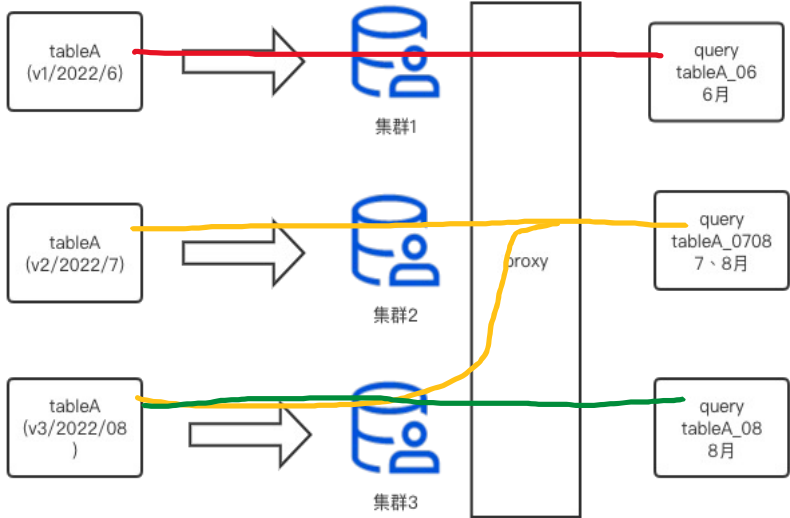
与社区版本脱节，目前部署模式不便版本更新

思考前后，各个组件需要做哪些改造？



改造前

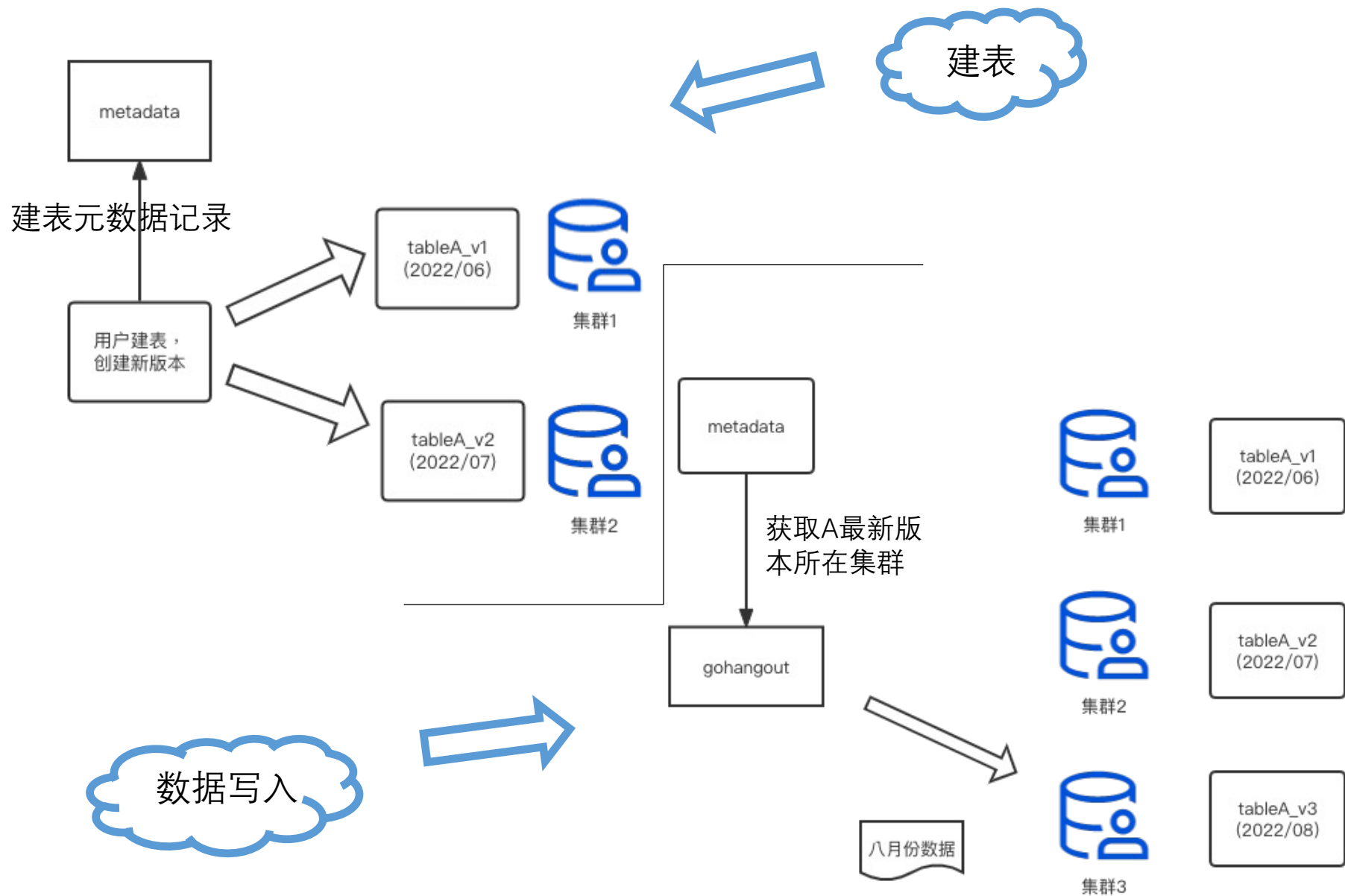
- 1. 统一的元数据管理
- 2. 写入集群选择
- 3. 查询代理



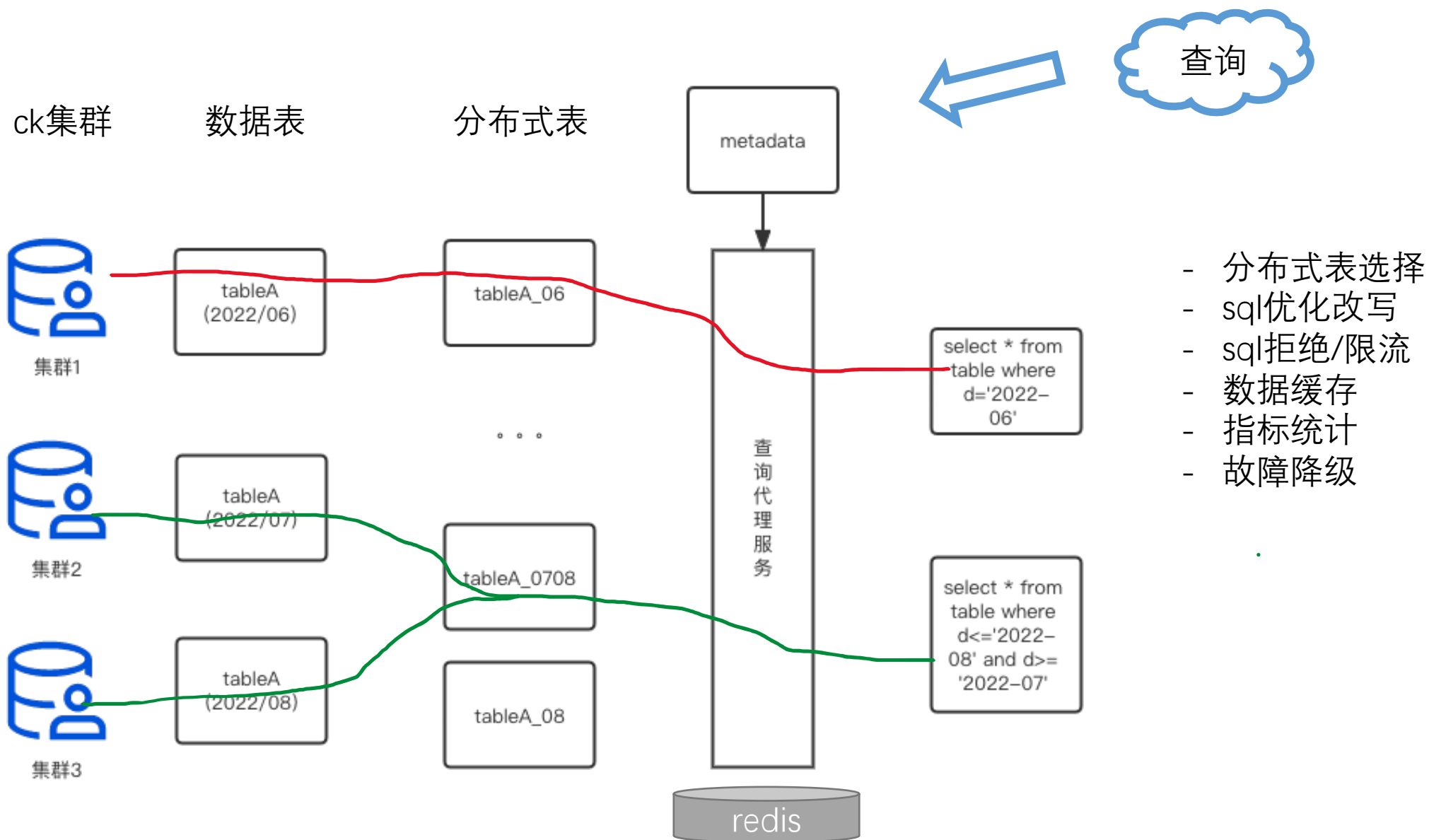
改造后

注意：这里以月份切分，实际上可以精准到分钟

如何解决数据跨集群



如何解决数据跨集群



基于antlr4的sql解析

```

CREATE TABLE test_local_1 ON
cluster ck100032722 ( `db_id` Int64 COMMENT 'db的id',
`desc` String COMMENT '描述',
`db_location_uri` String COMMENT '路径',
`name` String COMMENT '名称',
`dfs_quotas` String COMMENT 'quotas',
`dfs_usage` String COMMENT '使用量',
`owner_name` String COMMENT 'owner名称',
`owner_type` String COMMENT '类型',
`d` String COMMENT '日期分区键',
`test_string_1` String,
`test_string_2` String,
`test_string_3` String,
`test_1` String,
`test_2` String,
`test_3` String ) ENGINE = ReplicatedMergeTree('/clickhouse/tables/{replica}')
ORDER BY
d PARTITION BY d SETTINGS index_granularity = 8192
    
```



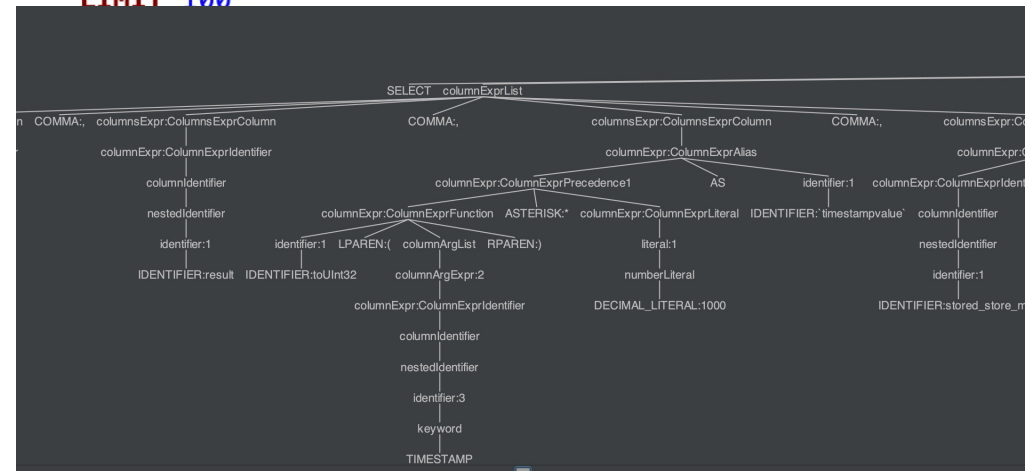
Select * from
tableA where d
= '2022-06'



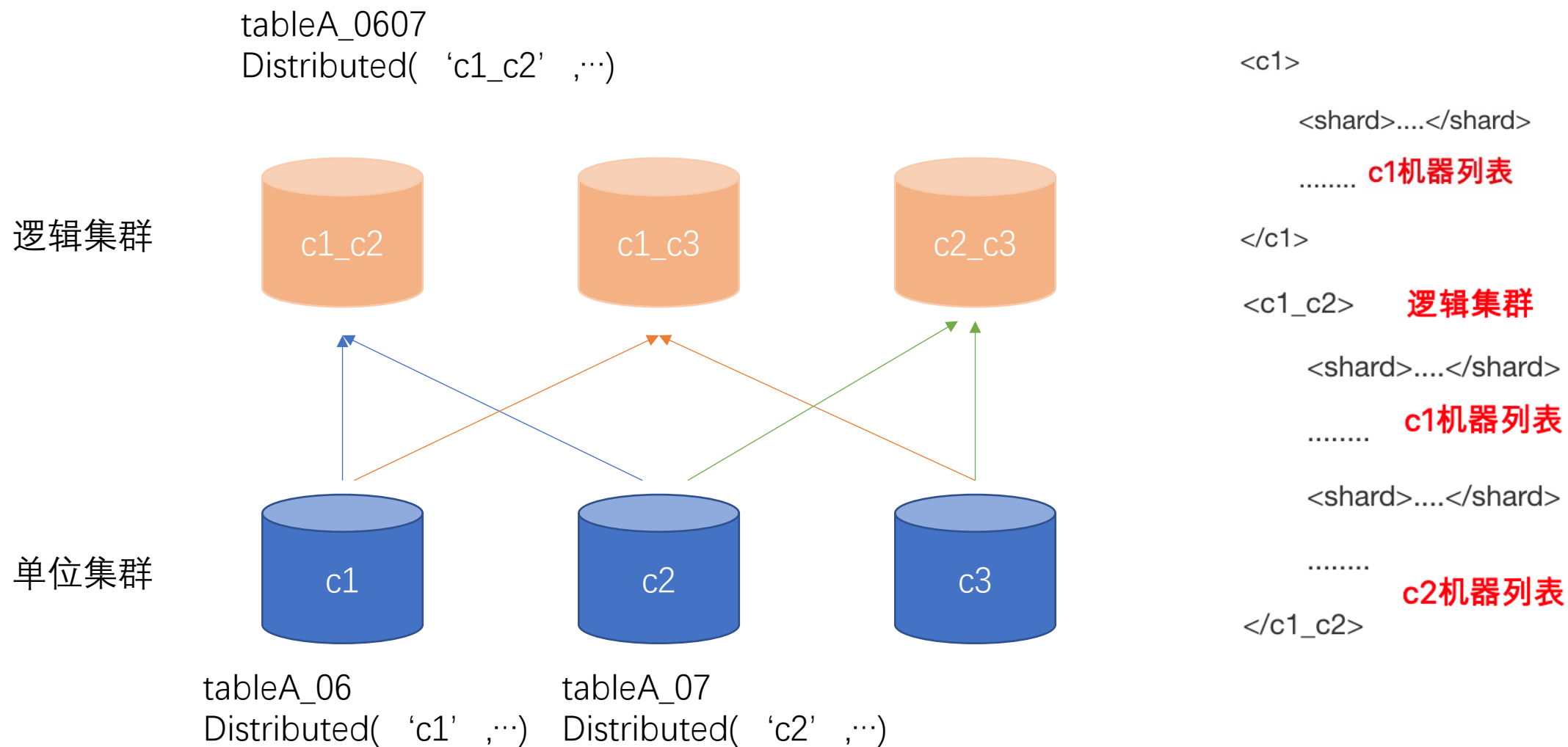
Select * from
tableA_06
where d =
'2022-06'

```

...
SELECT
orderid,
result,
toUInt32(timestamp) * 1000 as `timestampvalue`,
stored_store_msg as `message`
from
log.booking_schedule_all
where
timestamp >= toDateTime('?')
and timestamp <= toDateTime('?')
and project in ('transfersvrAsync')
and orderid = '?'
ORDER BY
timestamp desc
LIMIT 100
    
```



集群组合



未来展望

未来展望

- 降低成本
 - 冷数据存储
 - Juicefs、对象存储
 - 超冷机型
- 弹性混合云
- 日志数据预聚合

— Thanks —