

数据来源：数据库产品上市商用时间



第十三届中国数据库技术大会

DATABASE TECHNOLOGY CONFERENCE CHINA 2022

数据智能 价值创新



线上直播 | 2022/12/14-16



Apache Doris 在日志存储与分析场景的实践

肖康 SelectDB 联合创始人

01

Apache Doris 基本介绍

Introduction

| 关于 Apache Doris

Apache Doris 是一个基于 MPP 架构的高性能、实时的分析型数据库，以极速易用的特点被人们所熟知，仅需亚秒级响应时间即可返回海量数据下的查询结果，不仅可以支持高并发的点查询场景，也能支持高吞吐的复杂分析场景。基于此，Apache Doris 在多维报表、即席查询、用户画像、实时大屏、日志分析、数据湖查询加速等诸多业务领域都能得到很好应用。

Apache Doris 于 2022 年 6 月成功从 Apache 孵化器毕业，正式**成为 Apache 顶级项目**，截止目前 Apache Doris 社区已经聚集了来自不同行业百余家企业的超 400 位贡献者，每月活跃贡献者人数也接近 100 位。

Apache Doris 如今在中国乃至全球范围内都拥有着广泛的用户群体，截止目前，Apache Doris 已经在全球范围内**1000 家企业的生产环境中得到应用**，在中国市值或估值排行前 50 的互联网公司中，有超过 80% 长期使用 Apache Doris，包括百度、美团、小米、京东、字节跳动、腾讯、快手、网易、微博、新浪、360 等，同时在一些传统行业如金融、能源、制造、电信等领域也有着丰富的应用。

定位：极速易用实时统一的湖仓分析引擎

Data Sources

Database

Web Log

Mobile Log

Time Series Data

Data API

Data Ingestion and Processing

CDC (Change Data Capture)

ETL

Streaming ETL (Flink)

Batch ETL (Spark)

ETL Tools (DBT)

Data Integration Tools

Unified Data Warehouse



Real-Time Data Warehouse

Data Lake (Hive / Iceberg / Hudi)

Usage Story

Report Analysis

Ad-hoc Analysis

Federated Query

Machine Learning

02

典型日志存储与分析场景

Log Storage and Analysis

日志存储与分析场景



日志对于保障系统、业务稳定性至关重要，常用于故障排查、监控告警等。

特点：

1. 数据写入吞吐量大，还要实时可见
2. 数据存储量大，还要成本低
3. 交互式查询速度快，且支持文本检索、时间排序

典型方案与不足

	ES为代表的 倒排索引 检索架构	Loki为代表的 元数据索引/无索引架构
实时写入吞吐	低 (建倒排索引慢)	高
存储规模	中 (本地存储)	大 (存算分离)
存储成本	高 (多份数据存储)	低
semi & free text 交互式查询性能	快	慢 (数据无索引, 字符串匹配)
总结	优化查询性能 牺牲写入性能和存储空间	优化写入性能和存储空间 牺牲查询性能

典型方案与不足

优化这个牺牲那个，是不是头痛医头脚痛医脚了？

倒排索引是0-1选择吗，它是问题的全部吗？

向量化计算成熟前，是不是都认为OLAP加速要靠预计算？

03

日志场景解决方案

Solution

SelectDB日志场景解决方案

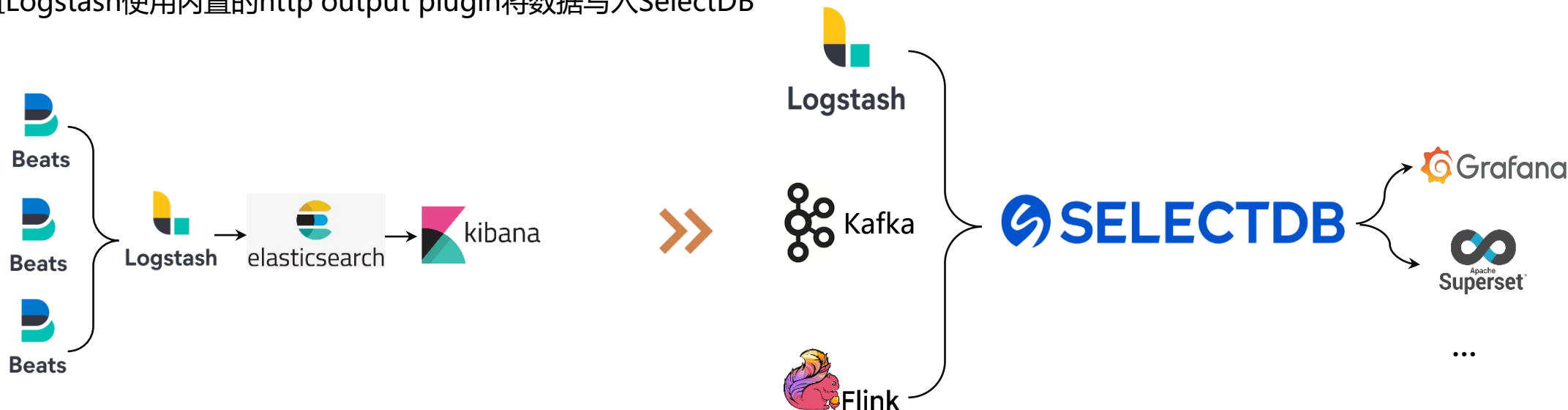
Doris 高性能向量化引擎底座 + SelectDB 存算分离架构 轻量级倒排索引 时序数据管理



SelectDB日志场景解决方案

上游写入

配置Logstash使用内置的http output plugin将数据写入SelectDB



下游查询

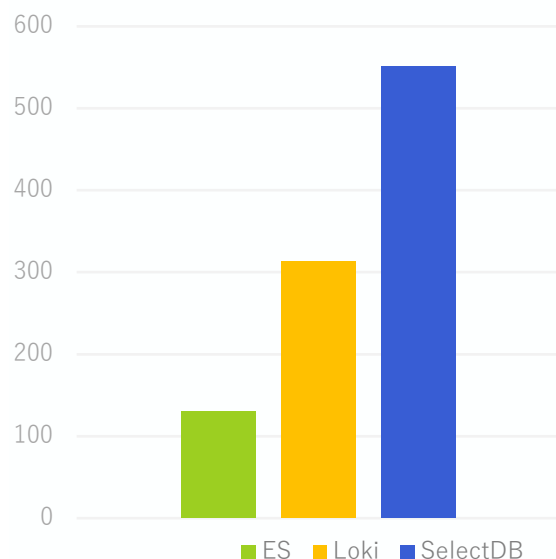
可观测性： Grafana中使用内置MySQL数据源，导入已有模板配置可视化日志看板、检索界面

商业智能： Superset等BI工具通过MySQL协议，即可开箱即用访问SelectDB进行可视化BI分析

性能测试

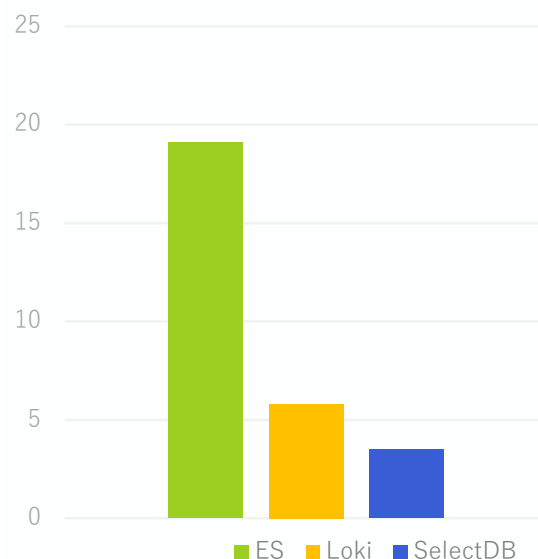
4.2倍

写入速度(MB/s)



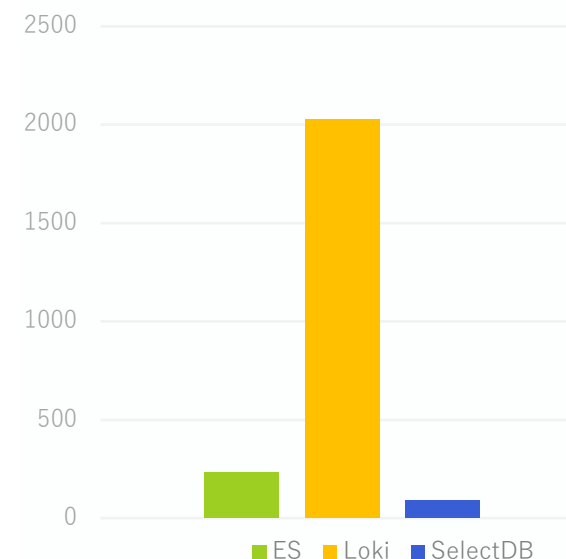
5.4倍

存储空间(GB)



2.3倍

查询时间(秒)



测试说明:

1. 测试环境是3台16c 64g云主机组成的集群
2. 测试数据和测试case来源于ES官方性能benchmark中http_logs, 数据总量32GB, 2.47亿行
3. 查询时间是ES官方性能benchmark中的11个query, 每个串行执行100次的总时间
4. 写入速度越高越好, 磁盘空间越低越好, 查询时间越低越好

04

关键技术解析

Technology

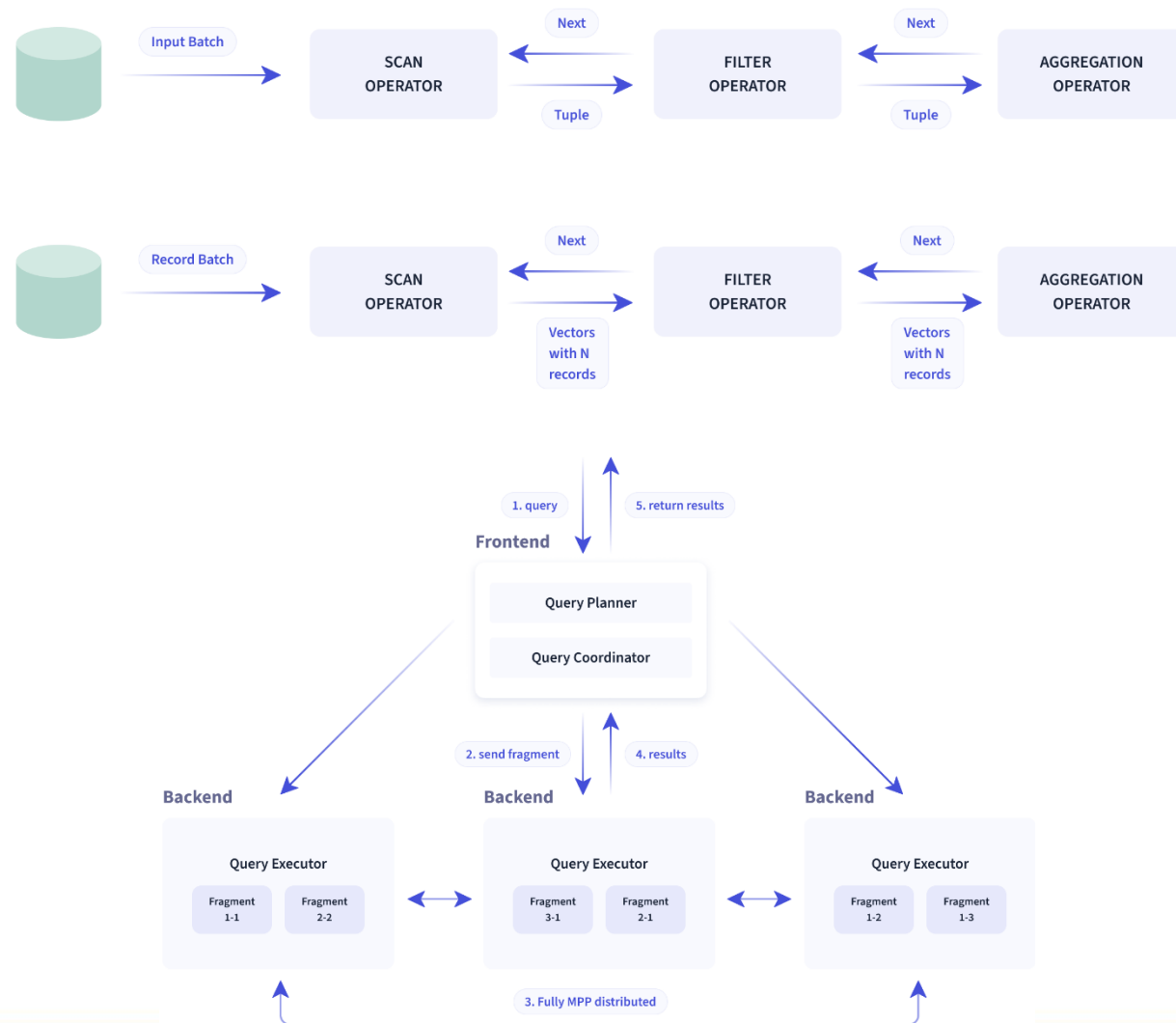
关键技术：MPP查询与向量化引擎

向量化

- 列式内存布局，向量化计算框架
 - ✓ 大幅减少虚函数调用
 - ✓ 大幅提升cache命中率
 - ✓ 高效利用SIMD指令
- 在宽表聚合场景下性能提升5-10倍

MPP查询

- 分布式MPP的查询框架，节点间和节点内都并行执行，大幅提升效率
- 支持大表的shuffle 分布式join



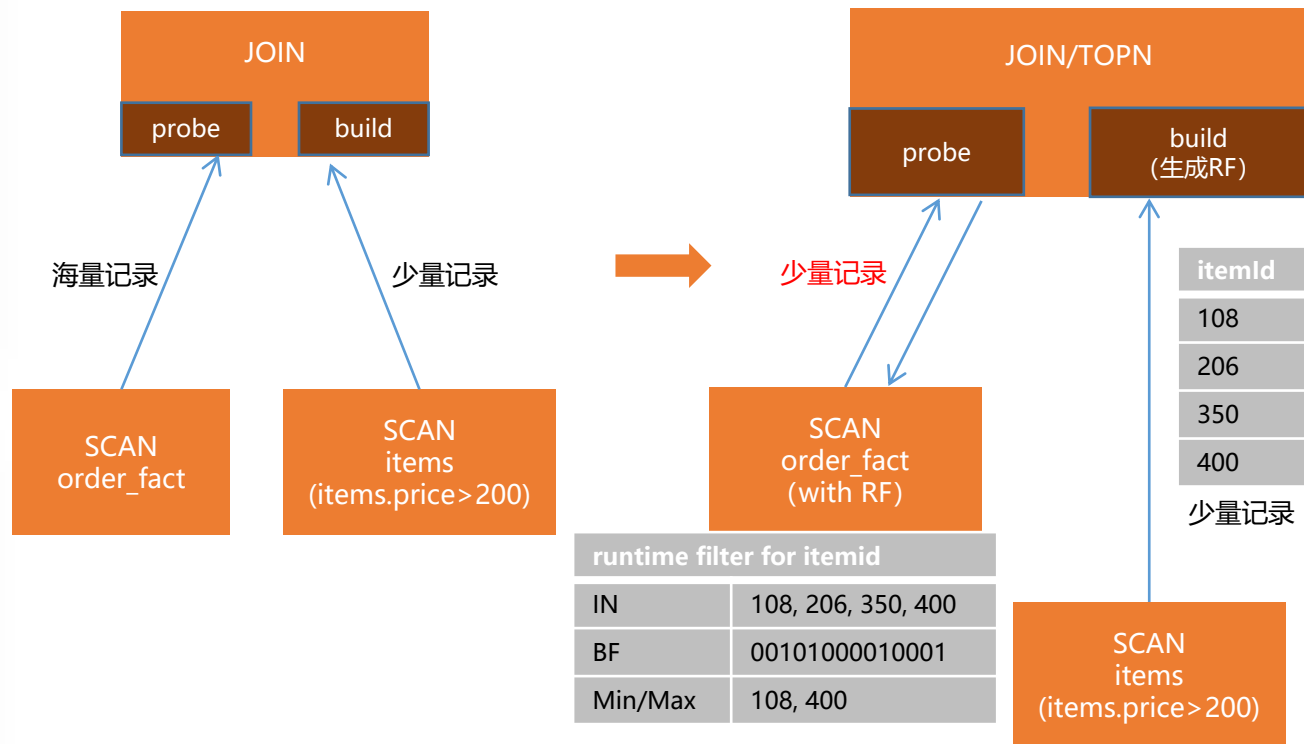
关键技术：多重算子优化与查询优化器

算子优化

- 自适应两阶段聚合算子优化
- JOIN/TOPN runtime filter优化
 - ✓ 为连接列生成filter推到左表
 - ✓ 支持in/min/max/bf等filter
 - ✓ filter自动穿透到最底层
- SSB部分查询依赖RF有2-10倍提升

优化器

- CBO和RBO结合的优化器
- RBO常见规则常量折叠、子查询改写、谓词下推等
- CBO支持Join Reorder
- 新一代智能优化器 (Nereids)



DTCC 2022
第十三届中国数据库技术大会
DATABASE TECHNOLOGY CONFERENCE CHINA 2022

ClickBench — a Benchmark For Analytical DBMS

[Architecture](#) | [Reproduce and Validate the Results](#) | [Add a System](#) | [Report](#) | [Hardware Benchmark](#)

System: [All] [Alpha partitioned] [Alpha single] [Aurora for MySQL] [Aurora for PostgreSQL] [Synapse] [Clou] [clickhouse local partitioned] [clickhouse local single] [ClickHouse] [ClickHouse IntD] [ClickHouse DvD] [CrateDB] [Databend] [Dorisaurus]

[Druid] [DuckDB] [Elasticsearch] [Elasticsearch tunnel] [Greenplum] [Hadoop] [InfluxDB] [MariaDB ColumnStore] [MareDB] [MonetDB] [MongoDB] [MySQL/MariaDB] [MySQL] [Presto] [PostgreSQL] [QuestDB (partitioned)] [QuestDB] [Redshift] [SAP HANA] [Singapore]

[Snowflake] [SQLite] [StarRocks (tuned)] [StarRocks] [TimescaleDB (compression)] [TimescaleDB]

Type: All [streaming managed] [range column-oriented C++] [MySQL compatible now-oriented] [C PostgreSQL compatible] [Clickhouse derivative] [embedded] [flink search document time-series]

Machine: All [serverless] [6acu L M S XS vfa-Axlarge-500gb-g2] [fa-metal, 500gb g2] 16 threads 20 threads 24 threads 28 threads 30 threads r5e x2 cfa-Axlarge-1000gb-g2 rx3.7xlarge rx3.xlarge rx3.xplus S24 S2 2PL 3XL 4XL

Cluster size: All [1] [2] [4] [8] [12] [16] [32] [64] [128] [serverless undefined]

Metric: Cold Run Hot Run Last Time Storage Size

System & Machine

SelectDB (efa-Axlarge, 500gb g2i7) +1%

Clickhouse (tunnel) (efa-Axlarge, 500gb g2i5) -1%

StarRocks (efa-Axlarge, 500gb g2i5) -2%

Clickhouse (efa-Axlarge, 500gb g2i5) -2%

MonetDB (efa-Axlarge, 500gb g2i5) -2%

Singapore (efa-Axlarge, 500gb g2i7) -5%

clickhouse-local (partitioned) (efa-Axlarge, 500gb g2i5) -5%

QuestDB (efa-Axlarge, 500gb g2i7) -11%

Presto (efa-Axlarge, 500gb g2i7) -12%

Dremppion (efa-Axlarge, 500gb g2i5) -12%

QuestDB (efa-Axlarge, 500gb g2i5) -12%

CrateDB (efa-Axlarge, 500gb g2i5) -12%

Mariadb Columnstore (efa-Axlarge, 500gb g2i7) -12%

clickhouse-local (single) (efa-Axlarge, 500gb g2i7) -12%

TimescaleDB (compression) (efa-Axlarge, 500gb g2i5) -12%

Druid (efa-Axlarge, 500gb g2i7) -17%

HeavyH (efa-Axlarge, 500gb g2i7) -46%

Giga (efa-Axlarge, 500gb g2i5) -113%

InfluxDB (efa-Axlarge, 500gb g2i7) -113%

TimescaleDB (efa-Axlarge, 500gb g2i5) -113%

SQLite (efa-Axlarge, 500gb g2i5) -113%

PostgreSQL (efa-Axlarge, 500gb g2i5) -113%

MySQL/MariaDB (efa-Axlarge, 500gb g2i5) -113%

MySQL (efa-Axlarge, 500gb g2i5) -113%

MySQL/MariaDB (efa-Axlarge, 500gb g2i5) -113%

Relative time (lower is better)

C6A.4XLARGE, 500GB GP2



查询总耗时远低于行业竞品

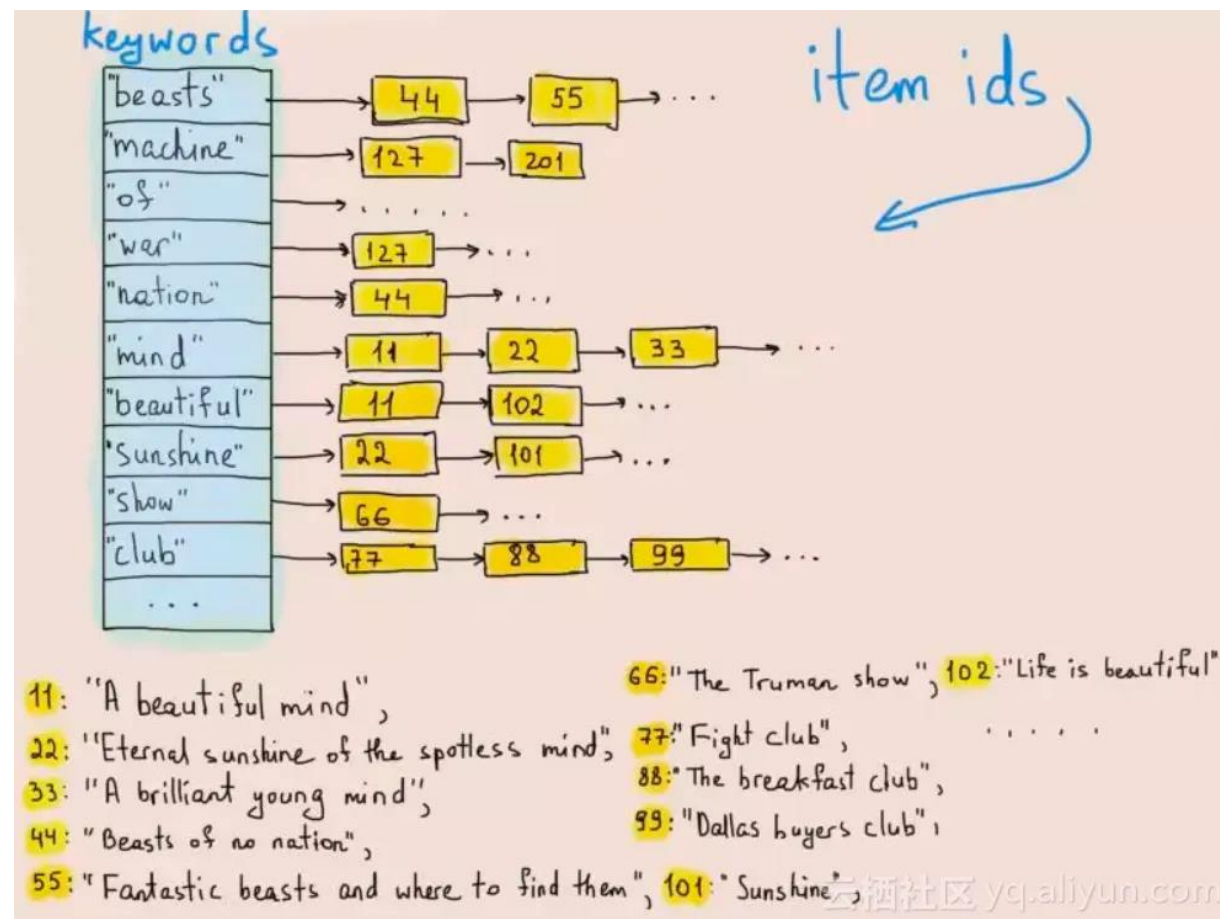
关键技术：轻量级倒排索引

支持快速检索

- 支持文本检索、普通数值/日期查找
- 支持多条件AND OR组合

扩展数据库引擎，内置倒排索引

- 避免了外挂式的跨系统通信、冗余存储



关键技术：轻量级倒排索引

为日志场景精简优化索引结构

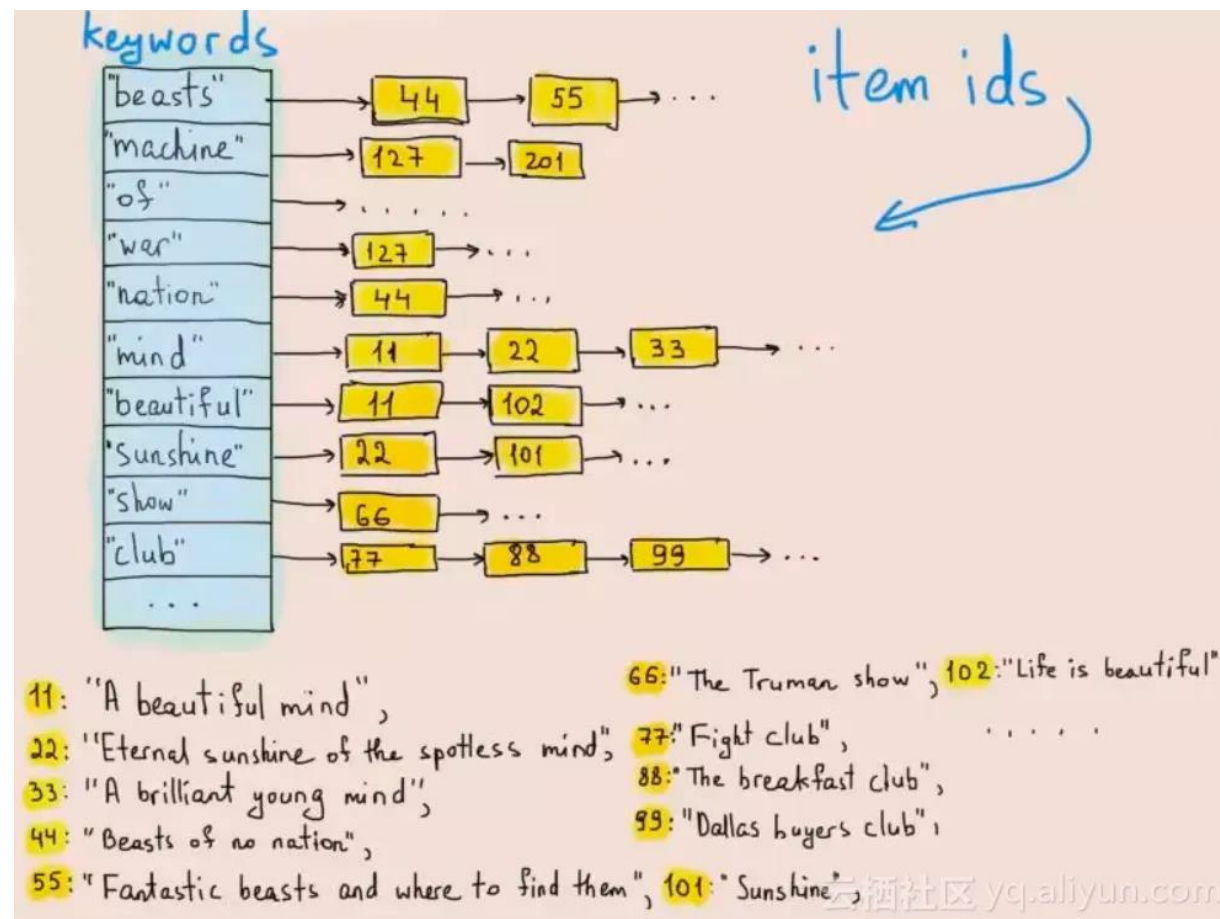
- 简化norm、score等
- bitmap等优化倒排表

列式存储 + ZSTD高效压缩算法

- 列式存储压缩率高
- zstd比gzip快5倍，压缩率更高
- 数据和索引都采用

采用C++和向量化的高性能实现

- 单核吞吐>20MB/s vs ES 5MB/s



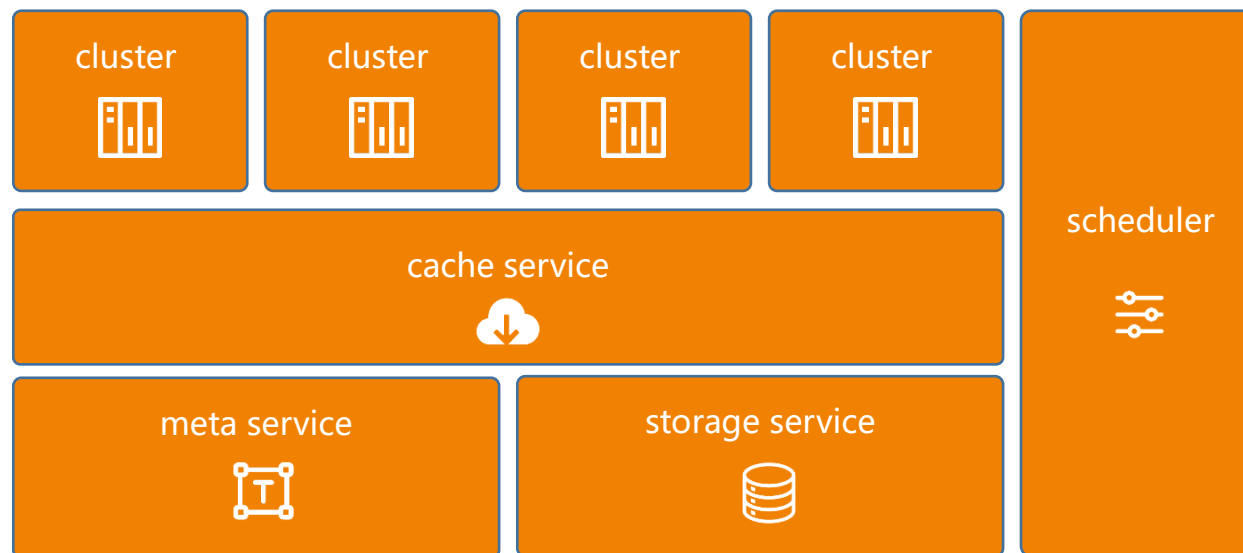
关键技术：存算分离云原生架构

存算分离，以对象存储为主存储

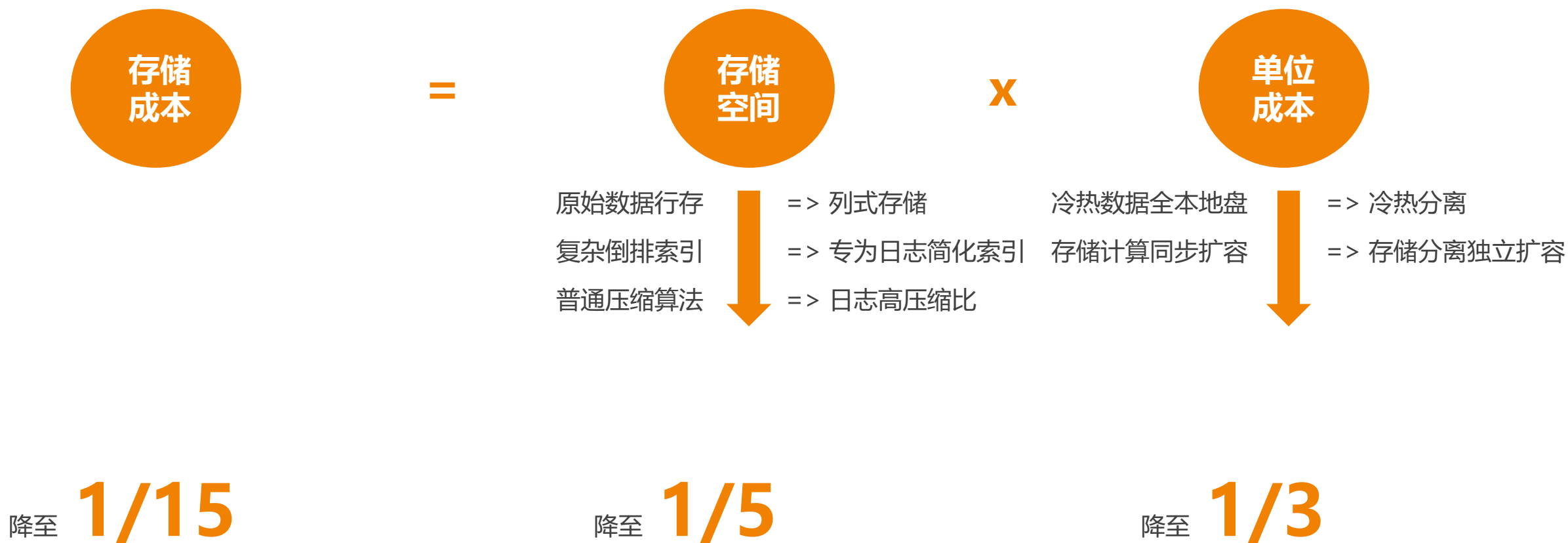
共享缓存，写入即缓存提高性能

弹性扩展，利用云的弹性加速查询

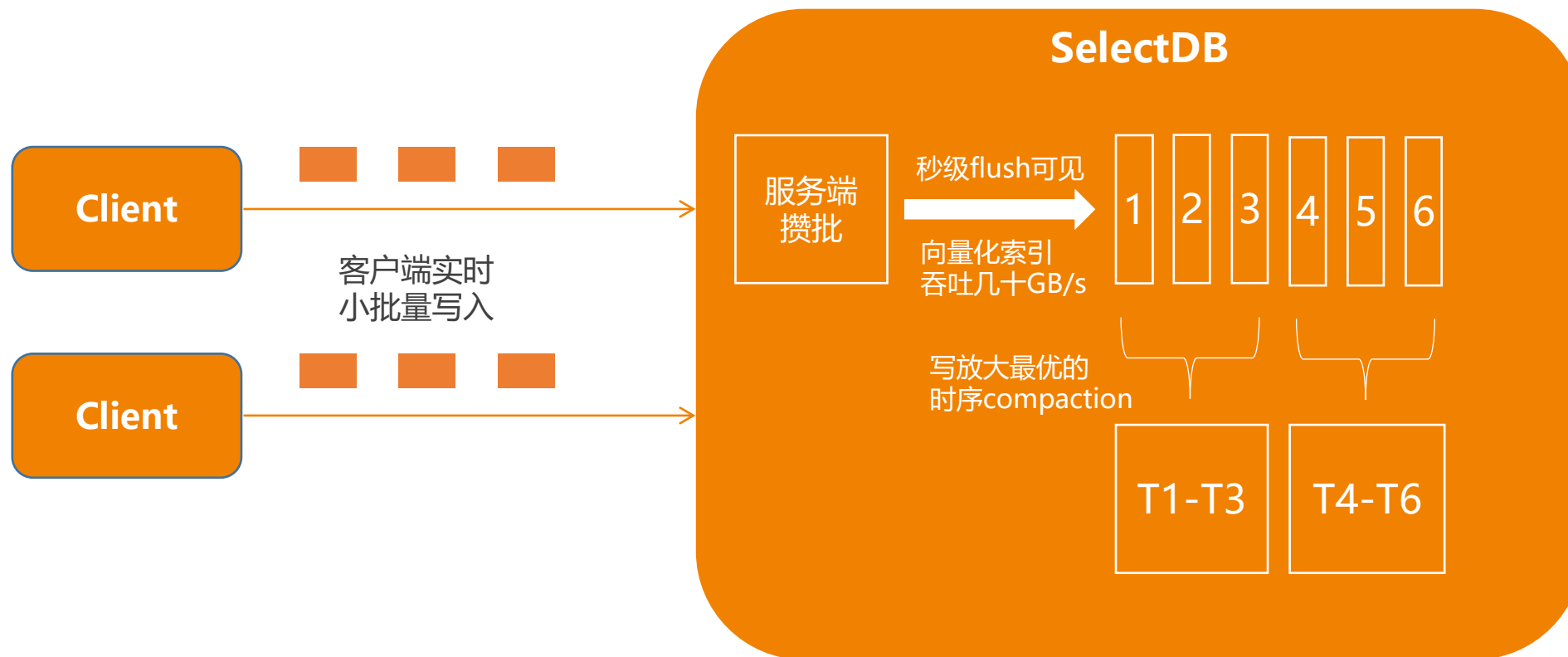
负载隔离，避免业务互相影响



存储成本大幅降低



关键技术：高吞吐实时写入



关键技术：快速交互式查询



```
SELECT * FROM log  
WHERE ts >= t1 AND ts <=t2 AND message MATCH 'error'  
ORDER BY ts DESC LIMIT 100
```

挑战：从海量日志中全文检索关键词



基于分区、主键的时间范围快速跳过
基于倒排索引的全文检索精确定位

挑战：从时间排序取满足条件的最新N条日志



按时间排序的时序存储模型
动态剪枝的TopN查询算法

百亿日志检索秒级响应



关于我们

Contact US

飞轮科技：专注于开源技术创新的云原生实时数据库厂商



开源数据仓库技术创新

秉持开源开放核心理念，大力投入研发力量，加强Apache Doris 在数据分析技术上的持续创新力，使其成为世界领先的开源分析数据库。



云端数据仓库商业服务

基于 Apache Doris，构建运行于多云之上的新一代云原生实时数仓 SelectDB，为客户提供极简运维和极致性价比的数仓服务。

联系我们



欢迎关注SelectDB微信公众号

公司邮箱: support@selectdb.com

SelectDB 官网: www.selectdb.com

Apache Doris 官网: <https://doris.apache.org/>

Apache Doris GitHub: <https://github.com/apache/doris>

THANKS

SQL Server
vertica
D B 2
G B a s e
O r a c l e
达梦数据库
神舟通用
KingbaseES

2010

2014

2018

openGauss
OceanBase
ArkDB
RASESQL
HotDB
StellarDB
QianBase xTP
云树Shard
GoldenDB
DolphinDB
MatrixDB
DynamoDB
SinoDB
FastData
Galaxybase
KunDB
GDB
GaussDB
PolarDB
KunDB
Spacture
Sequoiadb
OushuDB
ArgoDB
开务数据库
GreatDB
MongoDB
TDSQL
TiDB
Tapdata
UbiSQL
StarRocks