

Architect

SACC

2022 中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2022

· 激发架构性能 点亮业务活力

云上会议 网络直播 | 2022年10月20-22日

IT168.com

ChinaUnix.net

ITPUB

从 Apache Doris 到 SelectDB 下一代云原生实时数仓的演进之路

SelectDB 联合创始人& 产品VP 杨勇强

个人介绍



SelectDB 产品 VP

曾任职百度智能云存储部主任架构师

具有 10 年云存储产品与架构经验

本科毕业于南开大学 计算机科学与技术专业

硕士毕业于中国科学院计算技术研究所

杨勇强

Apache Doris Committer



目录 CONTENTS

01

Apache Doris 介绍

02

云原生时代的数据分析需求

03

下一代云原生实时数仓的演进之路

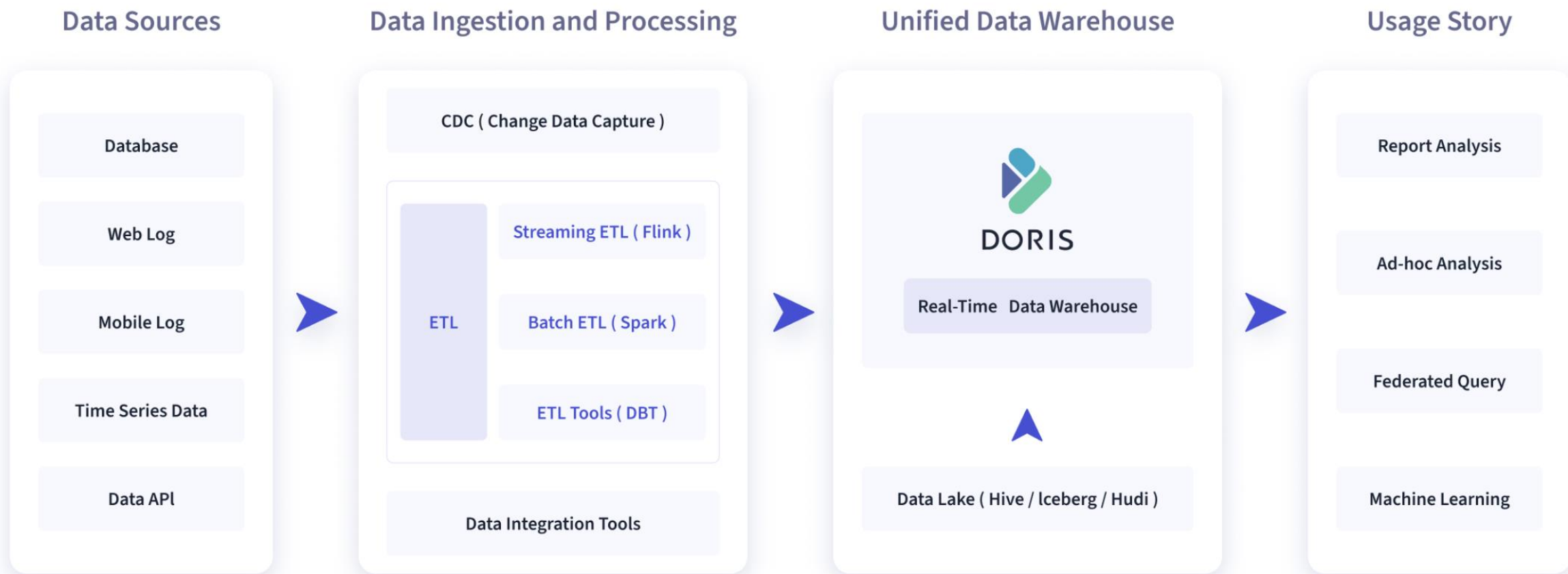
Apache Doris 介绍

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近百家企业的**近 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 **已经在全球接近 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。

易用、极速、实时、统一的湖仓分析引擎



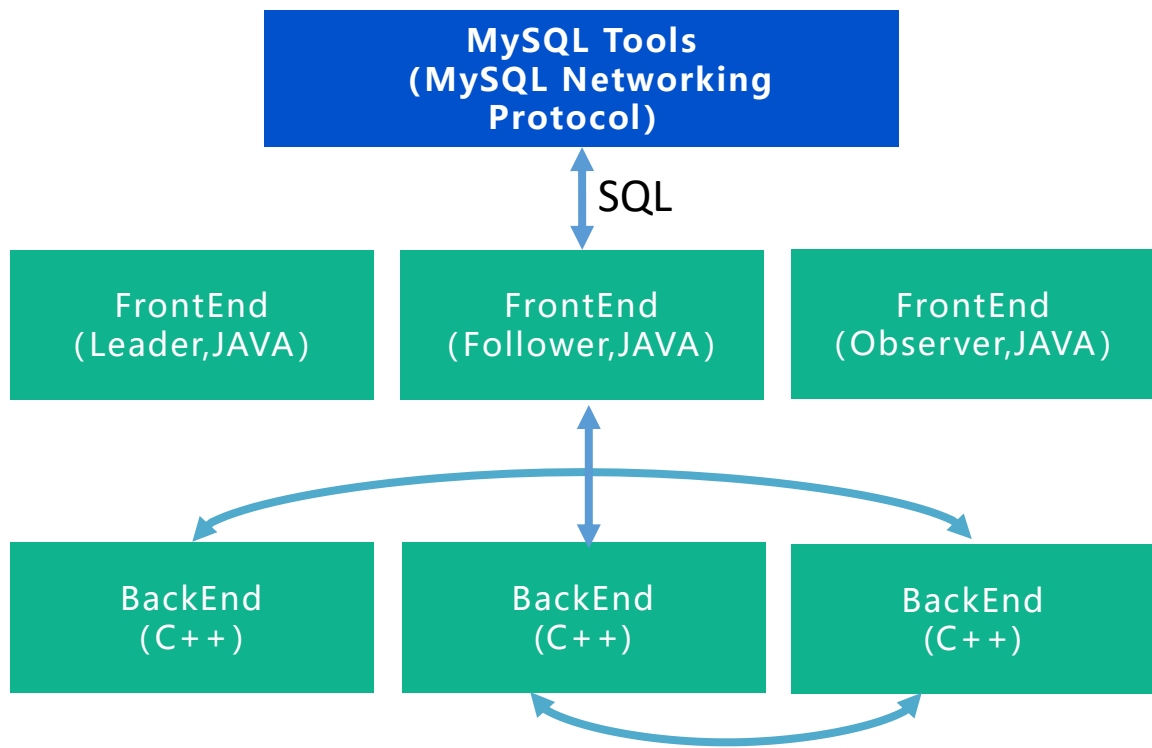
极简易用

业务侧

- 采用MySQL协议和标准SQL，各类客户端工具和BI可以无缝对接
- 支持多表大表join，不同场景有多种优化方案
- 支持子查询、窗口函数、udf/udaf, Grouping Set等高级功能
- 丰富的连接器，flink kafka等生态无缝对接，支持离线高效批量导入以及流式实时导入

运维侧

- 架构简洁，部署简单，只有FE和BE两类进程，无外部依赖
- 在线扩缩容，自动副本恢复
- 在线schema change



系统架构

极致性能（存储引擎）

列式存储

- 数据按照列存储，编码压缩高效
- 丰富的索引结构，减少数据扫描
 - ✓ sorted short key
 - ✓ min/max（等值、范围查询过滤）
 - ✓ bloom filter（高基数列等值过滤）
 - ✓ invert index（快速检索）

场景优化的存储模型

- aggregate key模型：相同key列value列合并，通过提前聚合大幅提升性能
- unique key模型：key唯一，相同key的数据覆盖，实现行级别数据更新
- duplicate key模型：明细模型
- 强一致物化视图，智能选择，加速查询

aggregate key模型

基础数据

时间	门店	销售额
2022/02/14	1	2000
2022/02/15	2	1500
2022/02/15	3	3000

新增数据

时间	门店	销售额
2022/02/14	1	200
2022/02/15	2	100
2022/02/15	4	4000

时间	门店	销售额
2022/02/14	1	2200
2022/02/15	2	1600
2022/02/15	3	3000
2022/02/15	4	4000

unique key模型

基础数据

订单	时间	状态
1	2022/02/14-22	待支付
2	2022/02/14-22	支付完

新增数据

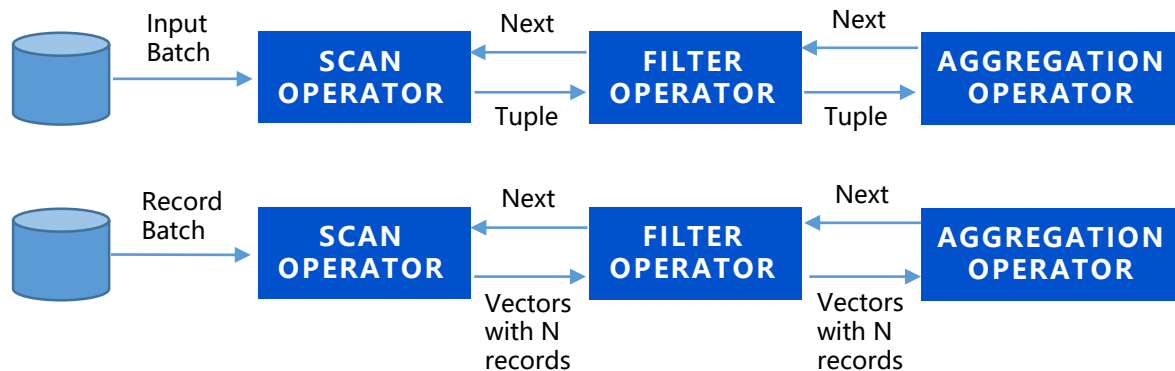
订单	时间	状态
1	2022/02/14-23	支付完

订单	时间	状态
1	2022/02/14-23	支付完
2	2022/02/14-22	支付完

极致性能（查询引擎）

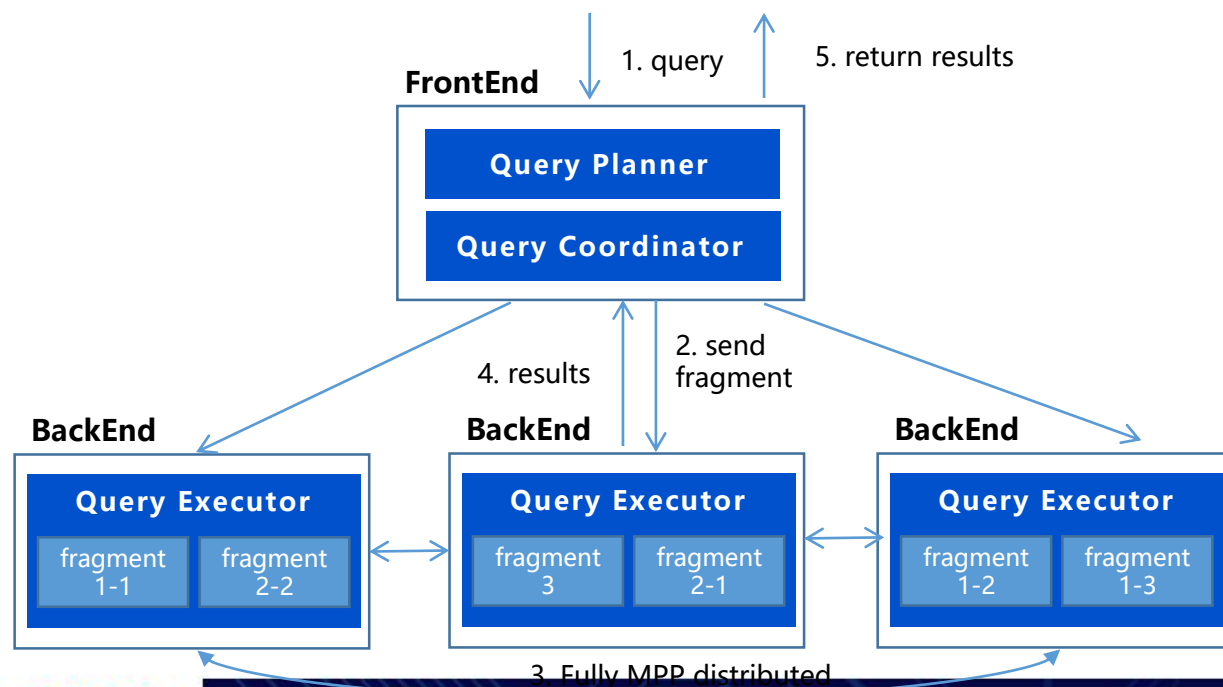
向量化

- 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



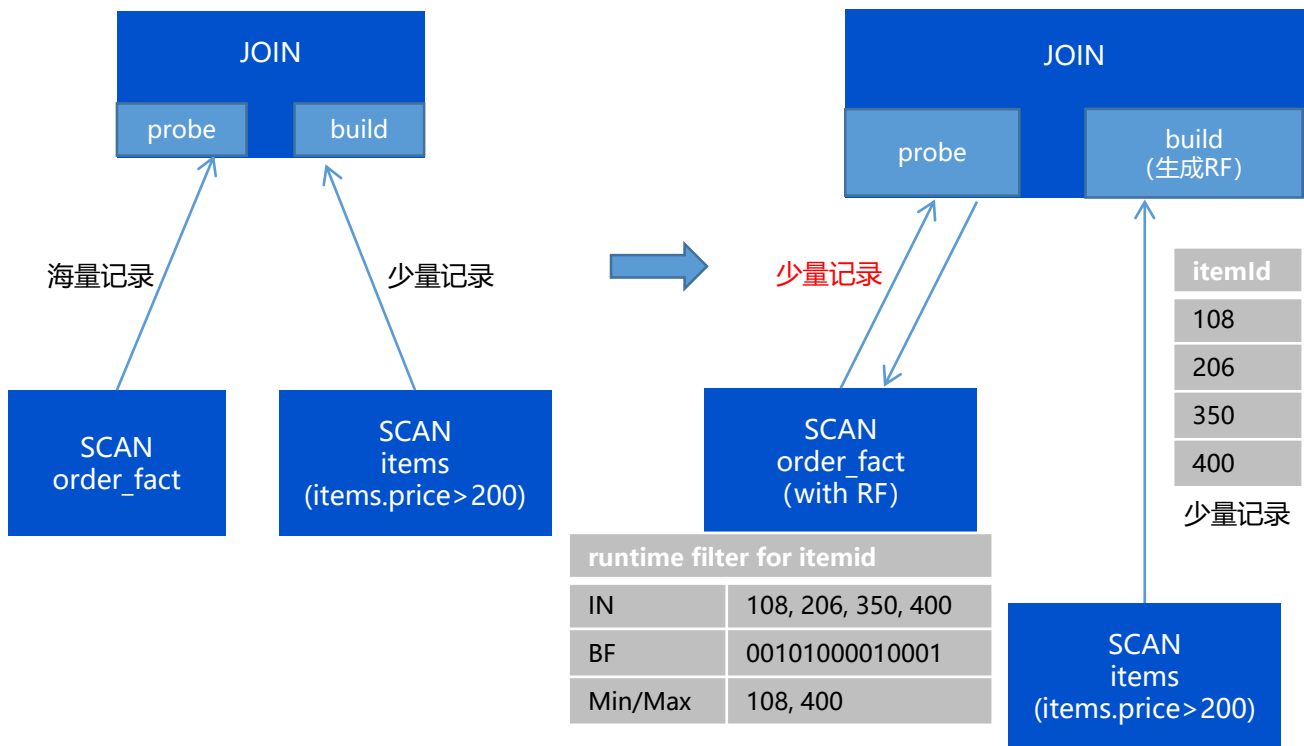
极致性能（查询优化器）

算子优化

- 自适应两阶段聚合算子优化
- Join的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

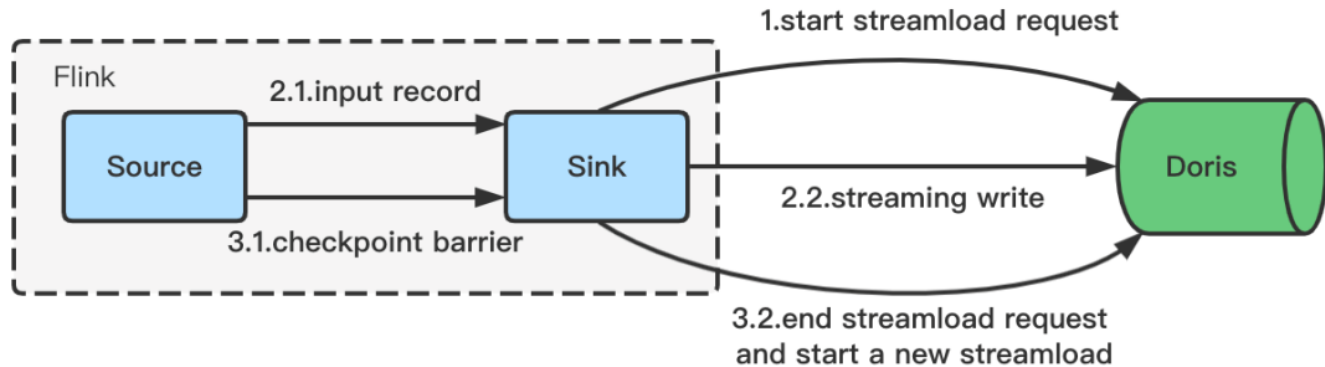
- CBO和RBO结合的优化器
- RBO常见规则常量折叠、子查询改写、谓词下推等
- CBO支持Join Reorder
- 新一代智能优化器（Nereids）



实时性保证

导入

- 两阶段提交，数据不重不丢
- 高并发微批，数据快速聚合
- flink / kafka 无缝集成
- online schema change

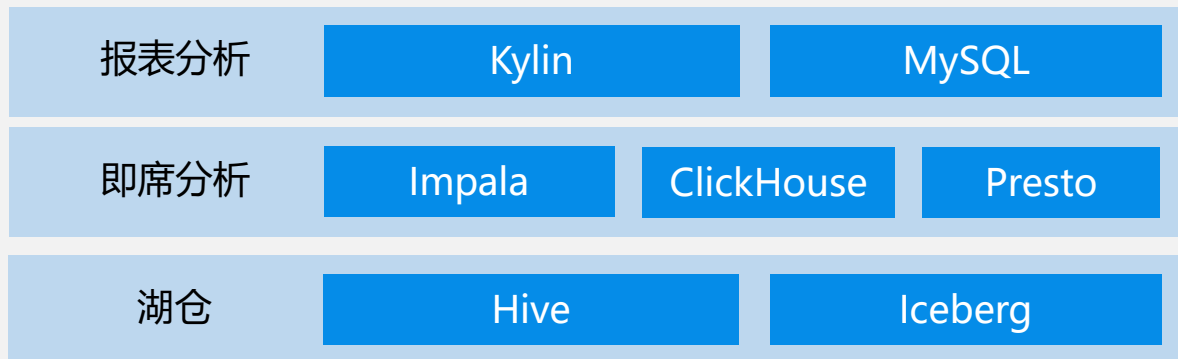


更新删除

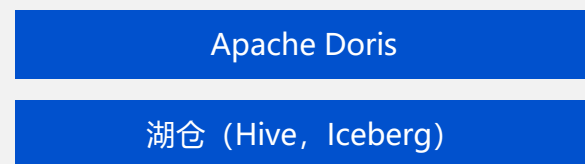
- 强一致
- 条件更新 (sequence column)
- 无冲突删除 (batch delete)

统一分析架构

过去的分析架构



基于 Doris 的统一分析架构



数据加工



统一 OLAP 分析引擎，从过去每个场景使用不同组件，到最终使用 Apache Doris 构建统一的分析架构，降低复杂架构带来的运维压力与技术融合承办。

使用 Apache Doris 作为Hive和Iceberg的**湖仓查询加速引擎**，性能相比Presto有 **3倍** 提升，相比Hive有 **数10倍** 以上提升。

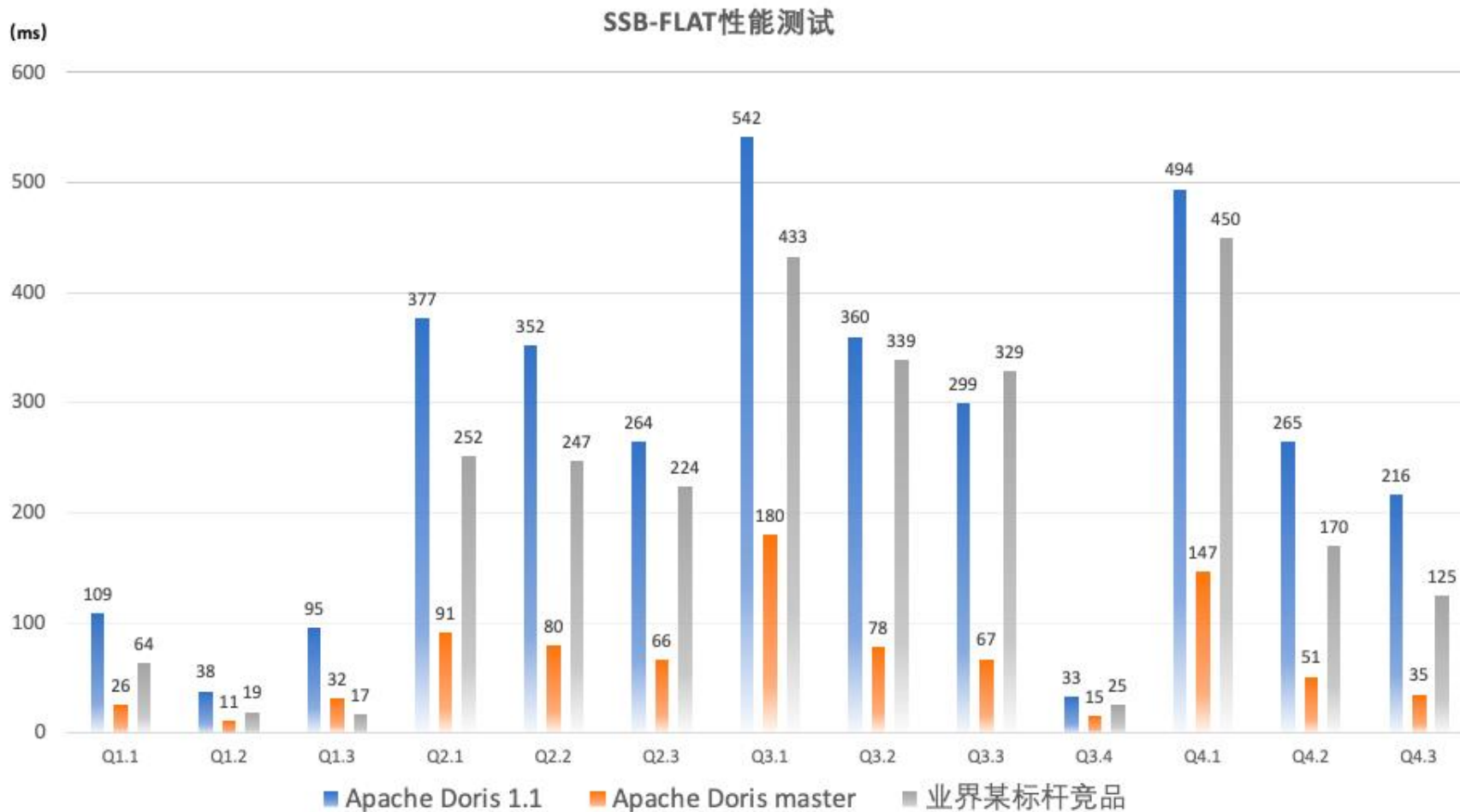
| Apache Doris 更多新功能

■ 1.2 新版本功能

- **极速数据湖分析**: 增加Multi Catalog, 无缝对接多种数据源, 性能超 Trino **3倍**
- **Light Schema Change**: 毫秒级在线元数据更新, 结合 Flink CDC 可实现 DDL 实时同步
- **支持实时更新的主键模型**: 全新 Merge on Write 主键模型, 实时更新场景性能提升 **10倍** 以上
- **JDBC数据源**: 提供便捷的外部数据源访问方式
- **Array 类型**: 支持嵌套、支持行列转换
- **New Decimal**: 更大精度、更高计算效率、更准确的精度信息
- **New Date/Datetime**: 更高计算效率、支持微妙存储
- **Java UDF**: 兼容 Hive UDF
-

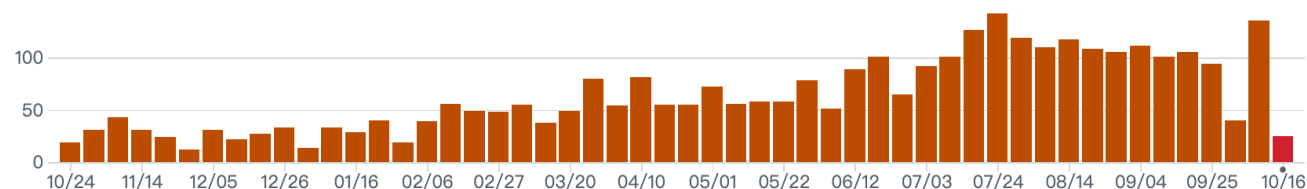
■ **发布时间 : 2022 年 10 月底**

Apache Doris 性能提升

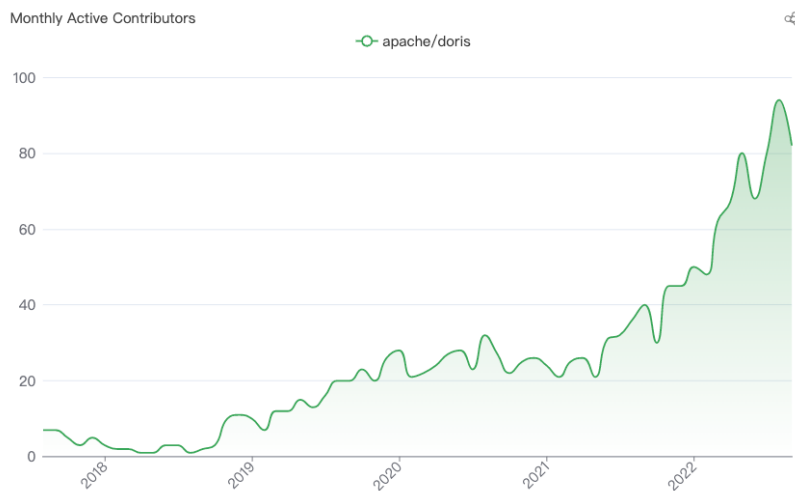


最新 1.2 版本较 1.1 版本提升近 **4** 倍，是业内标杆竞品 **3** 倍以上

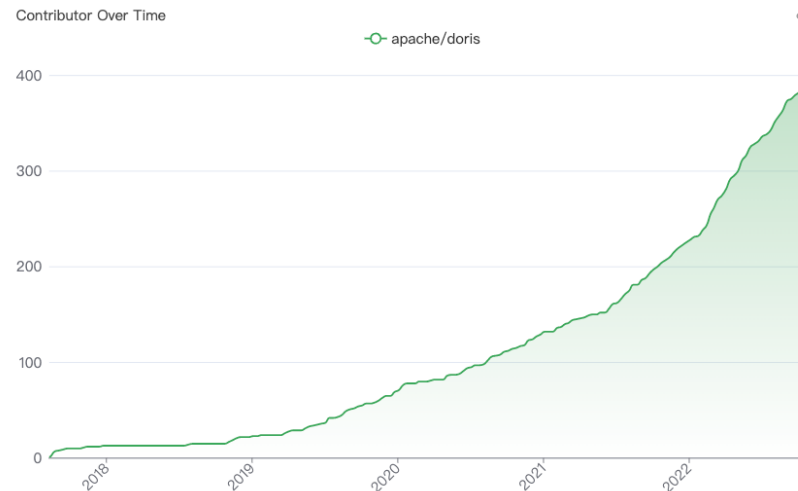
Apache Doris 社区发展 – 期待更多开发者的加入



每周超过100个Commits



活跃贡献者



累计贡献者



目录

CONTENTS

01

Apache Doris 介绍

02

云原生时代的数据分析需求

03

下一代云原生实时数仓的演进之路

大数据技术发展

- 宏观经营分析、业务决策
- 固定报表、批处理报告、BI

- 互联网业务大规模、多场景的数据计算和分析
- 批量计算、流式计算、报表和自助式分析

- 数据共享与联动分析、业务的数据驱动与敏捷创新
- BI应用、实时分析、湖仓联邦查询、云服务化

传统数据仓库时代

互联网大数据时代

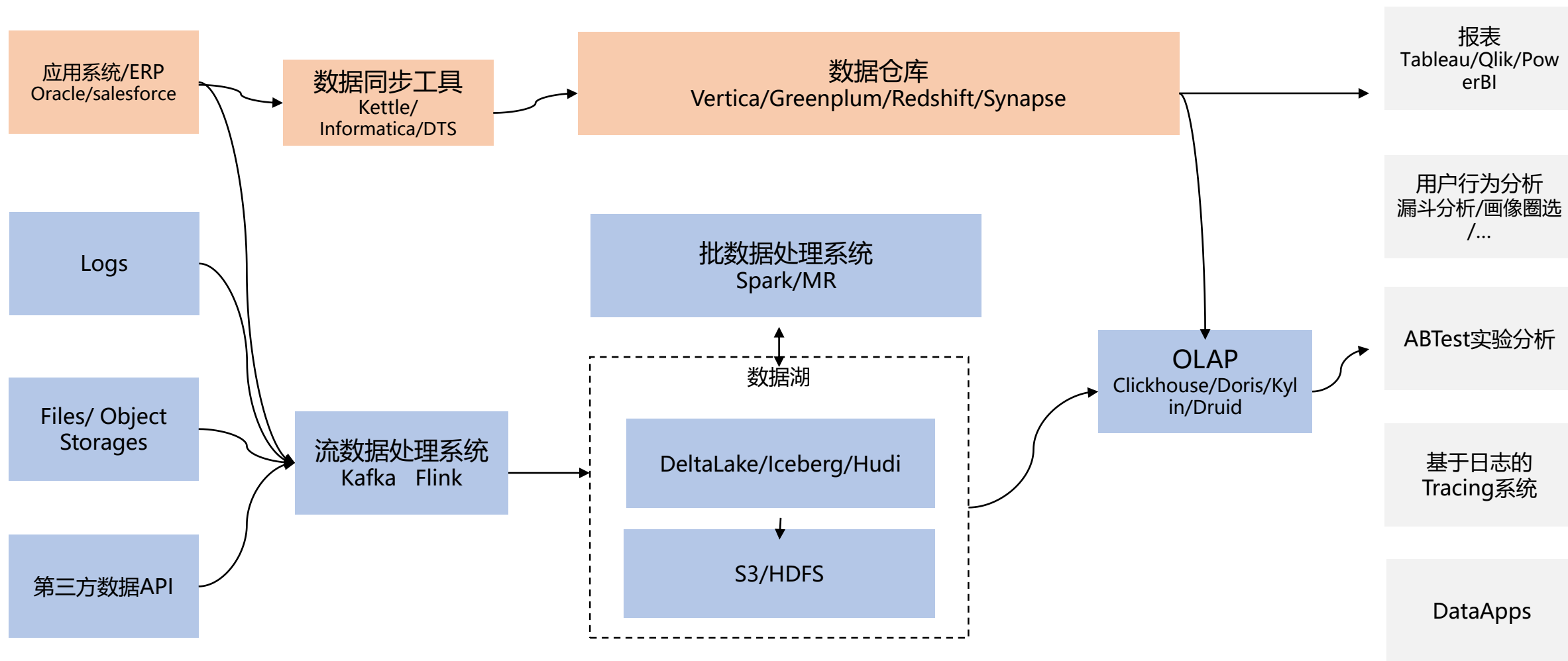
云原生数据分析时代

- Teradata
- Greenplum
- Vertica

- Hadoop
- Spark/Flink
- Hive
- Presto/Impala
- Doris/Clickhouse

- Snowflake
- Databricks
- AWS Redshift
- SelectDB

当前数据分析技术栈



云原生时代数据分析需求

日益增长的实时数据分析诉求

- Increasing Real-time Requirement -

复杂性高	性价比低	灵活性差	开放性弱
<ul style="list-style-type: none">■ 软件维护的复杂性■ 使用的复杂性■ 业务系统维护的复杂性	<ul style="list-style-type: none">■ 从on-premise到cloud-native 需要新设计;■ 新硬件的发展突飞猛进, 需要适配;	<ul style="list-style-type: none">■ 分析的数据类型与分析 work-load的变化没有被很好满足;■ 对半结构化和非结构化数据没有native高效支持■ 数据科学、机器学习等深度分析场景满足	<ul style="list-style-type: none">■ 客户出于防止锁定、容灾和降本的多云诉求■ 客户期待有第三方厂商推出的多云中立的, 统一接口和体验的, 开源开放的 PaaS层产品

用户满意度痛点

云原生时代数据分析技术

性能 - 实时 realtime

- 数据服务实时化 (HTAP、HSAP)
- 数据分析实时化 (向量化、代码生成、CBO)
- 数据处理实时化 (流计算Flink、表格式Iceberg)

功能 - 统一 unified

- 湖仓一体
- 在离线一体
- 流批一体
- 结构化、半结构化、日志等多数据类型支持

开源技术

全新
数据分析
基础设施

云端服务

云原生 cloud-native

- 存算分离、弹性使用
- 计算隔离，多租户
- 极致性价比
- 极简使用和运维

多云 multi-cloud

- 多云中立
- 统一接口
- 多云复制



目录

CONTENTS

01

Apache Doris 介绍

02

云原生时代的数据分析需求

03

下一代云原生实时数仓的演进之路

SELECTDB Cloud 概览

selectdb_dw

集群

连接

查询

监控

用量

设置

集群

load

运行中

运行时长 0 天 0 小时 1 分钟

内核版本 2.0

创建时间 2022-10-26 17:23:58

重启 详情

etl

运行中

运行时长 0 天 0 小时 0 分钟

内核版本 2.0

创建时间 2022-10-26 17:25:07

重启 详情

serving

运行中

运行时长 0 天 0 小时 0 分钟

内核版本 2.0

创建时间 2022-10-26 17:25:02

重启 详情

+ 新建集群

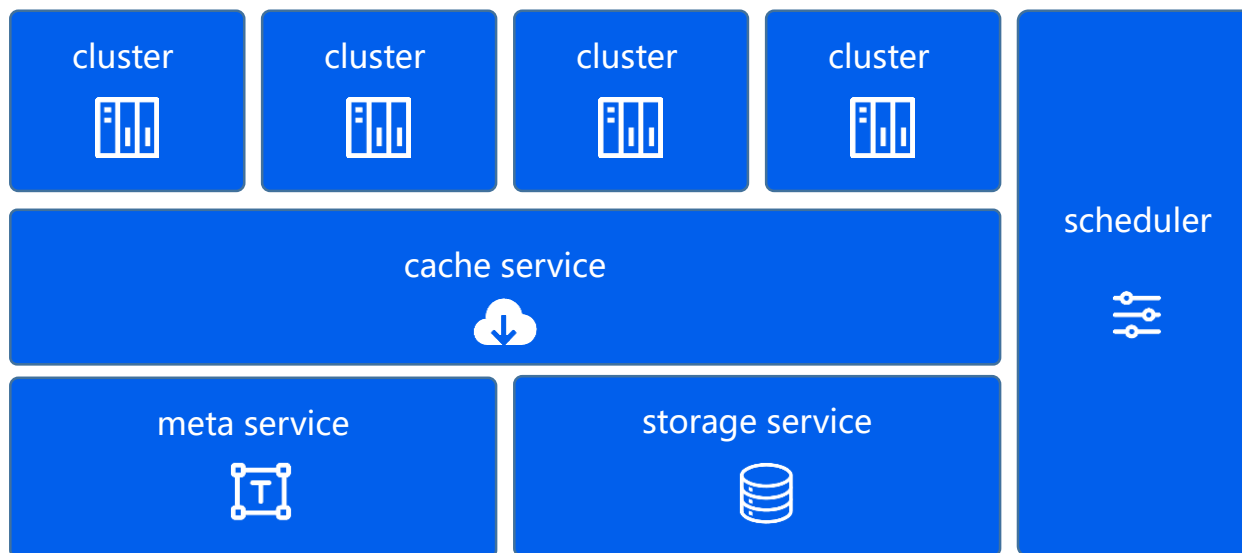
SELECTDB Core 架构

存算分离，以对象存储为主存储

共享缓存，确保在线高可用

timetravel 与 snapshot

数据共享与跨区域复制



SELECTDB 核心功能



极致性能

- 高效列式存储与索引
- 现代 MPP 计算架构
- 适配 X64、ARM64 的向量化执行引擎



融合统一

- 实时/交互/批量数据处理
- 结构化/半结构化数据
- 对数据湖和其他数据库进行联合查询



简单易用

- 兼容 MySQL 连接协议
- 强大的管理控制台
- Spark/Flink/DBT/Kafka 等丰富的连接器



高性价比

- 存算分离
- 按需自动扩缩容
- 冷热数据分级存储



开源开放

- 基于开源 Apache Doris
- 与 Doris 数据自由迁移
- 全云可用

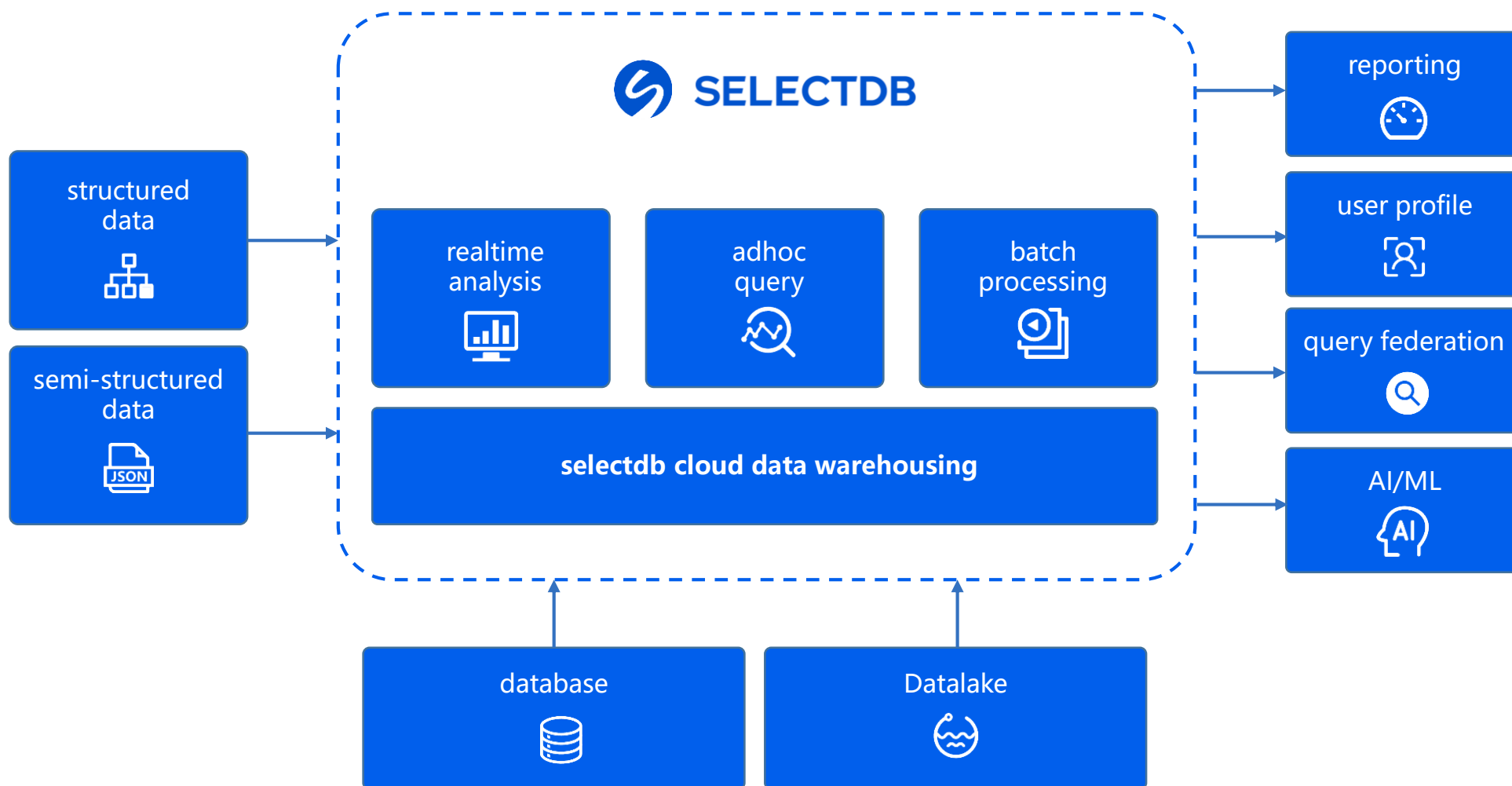


企业特性

- 用户管理与访问控制
- 数据保护与备份
- 数据治理促进数据共享使用

查询总耗时远低于行业竞品

融合统一



融合统一 | 混合负载

高并发报表

- 支持数千 ~ 数十万的QPS
 - ✓ 分区裁剪
 - ✓ sorted short key index
 - ✓ 物化视图
 - ✓ SQL, Partition Cache

湖仓 联邦分析

- 用SelectDB作为统一SQL查询层查询已经存在于离线湖仓中的数据
- 相对presto和hive性能提升数倍
 - ✓ MPP和向量化
 - ✓ 智能优化器（下推和Join）
 - ✓ 数据和元数据cache

即席分析

- 支持即席查询，秒级响应
 - ✓ MPP和向量化
 - ✓ 智能索引skip
 - ✓ 智能优化器

统一数仓

- Data Engineering：ETL
- Data Analytics：即席分析 + 湖仓联邦分析
- Data Serving：报表分析

ETL

- DAG执行模型增强
- Native支持半结构化和非结构化数据

混合负载 管理

- 资源隔离（硬限, 软限）
- 落盘与容错模型
- WorkLoad Manager (WLM)

融合统一 | 半结构化分析

JSON数据

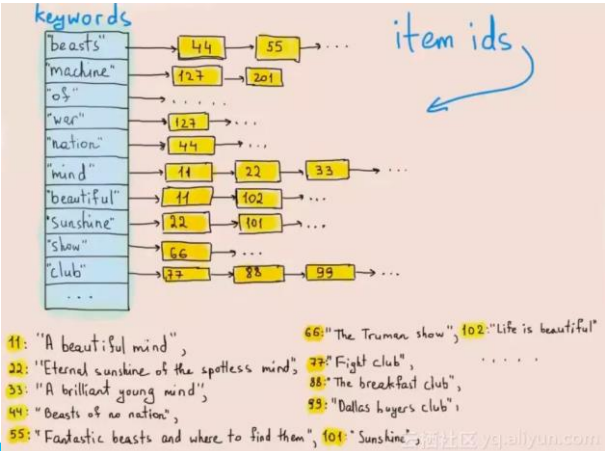
高效点查索引

Free Text

对非数值型字段提供Inverted Index索引

对数值型字段提供BKD Index索引

```
1  -- 建表
2  CREATE TABLE IF NOT EXISTS t_dynamic
3  (
4      name string,
5      age int,
6      ...
7  )
8  DUPLICATE KEY(`name`)
9  DISTRIBUTED BY HASH(`name`) BUCKETS 10
```



```
1 SELECT UNIX_TIMESTAMP(log_time) as time, host, type, level, msg
2 FROM dorislog
3 WHERE log_time >= $__timeFrom() AND log_time <= $__timeTo()
4   AND type = '$log_type'
5   AND msg MATCH_ANY '$keywords'
6 ORDER BY time DESC
7 LIMIT $limit
```

简单易用

兼容 MySQL 链接协议，标准SQL语法

基于 WebUI 的数据库管理和开发工具

Spark/Flink/Kafka/DBT 等丰富的连接器

The screenshot displays the SACC 2022 web interface. On the left is a sidebar with navigation options: 集群 (Cluster), 连接 (Connection), 查询 (Query), 监控 (Monitoring), 用量 (Usage), and 设置 (Settings). The '查询' (Query) option is selected. The main area shows a SQL query editor with the following code:

```
1 SELECT c_nation, s_nation, d_year,
2 SUM(lo_revenue) AS REVENUE
3 FROM customer, lineorder, supplier, dates
4 WHERE lo_custkey = c_custkey
5 AND lo_suppkey = s_suppkey
6 AND lo_orderdate = d_datekey
7 AND c_region = 'ASIA'
8 AND s_region = 'ASIA'
9 AND d_year >= 1992 AND d_year <= 1997
10 GROUP BY c_nation, s_nation, d_year
11 ORDER BY d_year ASC, REVENUE DESC;
12
```

Below the query editor, the '查询结果' (Query Results) tab is active, showing a table with the following data:

c_nation	s_nation	d_year	REVENUE
INDONESIA	INDIA	1992	272807600708
CHINA	INDIA	1992	271490331577
VIETNAM	INDIA	1992	270829651288
INDIA	INDIA	1992	269209915240
JAPAN	INDIA	1992	269067447506
CHINA	JAPAN	1992	268236710029
CHINA	CHINA	1992	267210920993
CHINA	VIETNAM	1992	266772531078
JAPAN	JAPAN	1992	266502358467

At the bottom of the interface, a status bar indicates: 执行成功 当前返回 150 行, 耗时 2s202ms, 单次查询结果最多返回 1000 条.

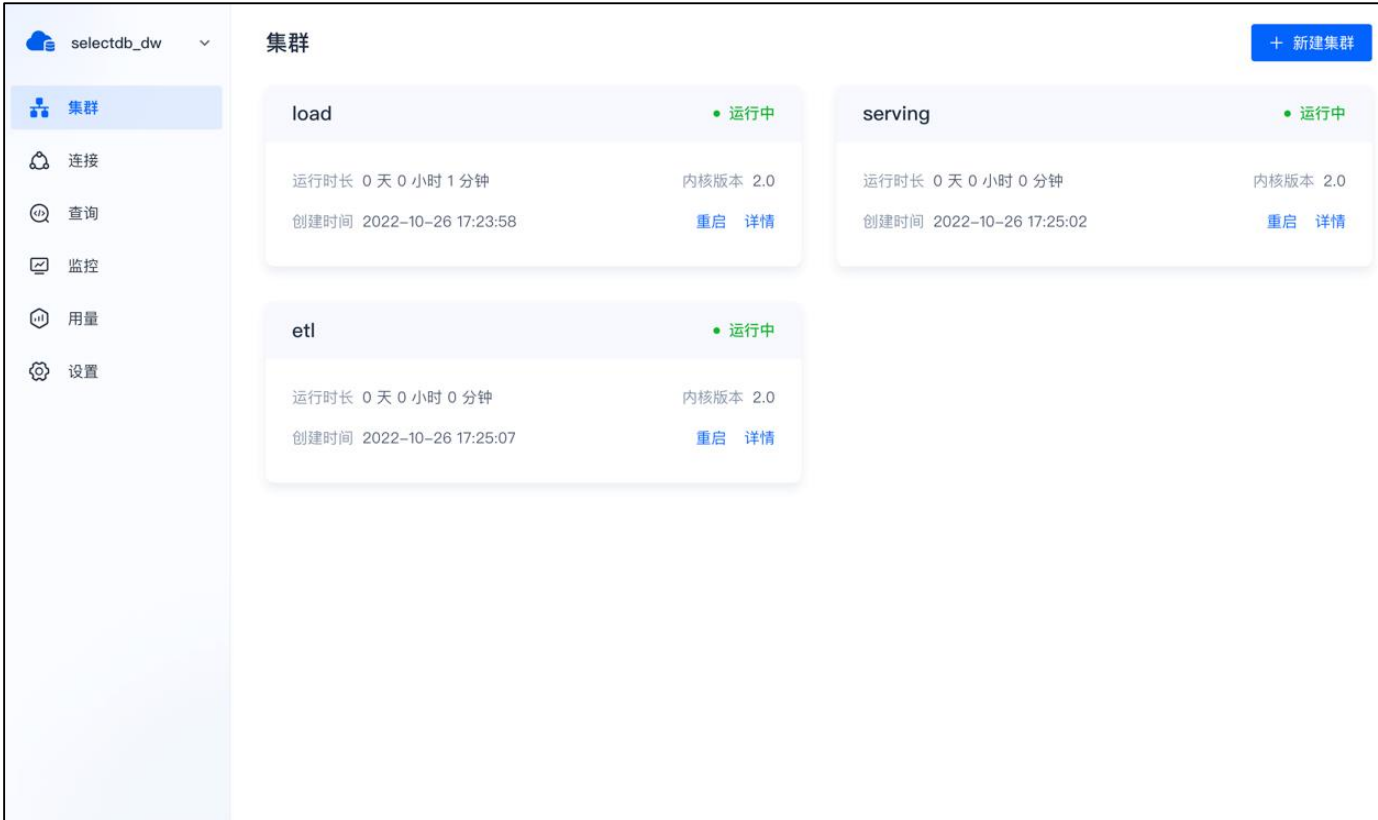
高性价比

存储按使用付费，计算按需弹性

数据共享，消除冗余数据

负载隔离，避免负载互相影响需要的额外资源

突发或者后台任务，云上竞价资源灵活调度



开源开放

- 基于Apache Doris研发，积极贡献Doris，充分利用社区力量和应用场景打磨引擎
- 多云中立，运行在多云之上，提供统一的产品体验，防止被锁定
- 兼容Doris数据格式，方便开源用户到商业产品的迁移



| 企业特性 | 安全

数据保护：落盘和传输加密，密钥管理，数据脱敏

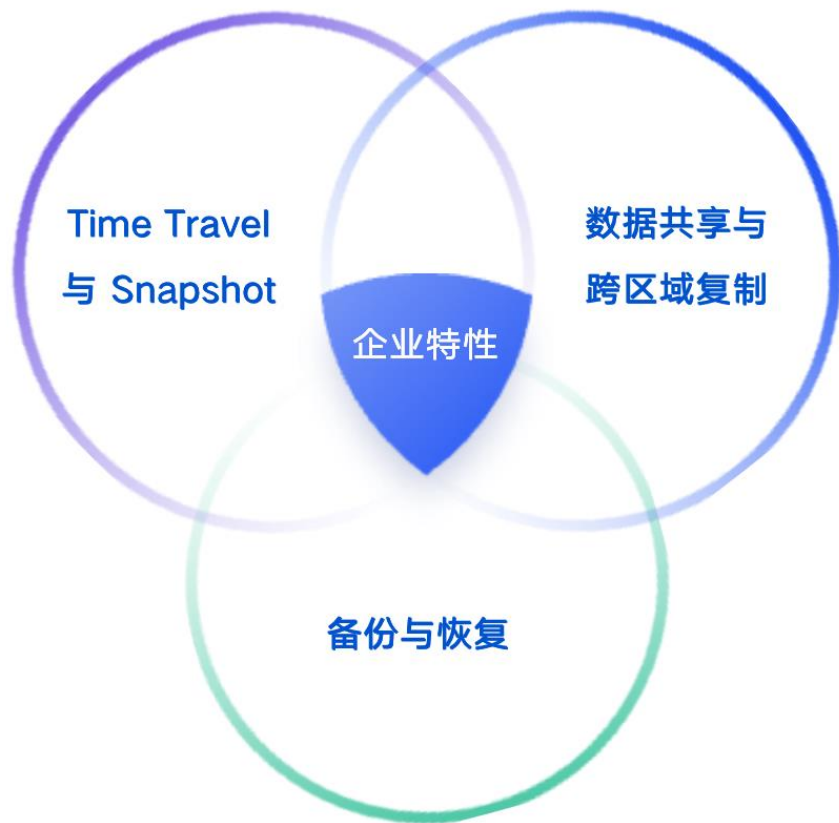
访问控制：RBAC，细粒度访问控制（行级别和列级别）

认证：多因子认证，IAM，ID Federation

网络安全：网络隔离

企业特性

- 存储表或者数据库的各个导入版本，并且支持对表或者数据库建立快照
- 同一地域多个数据库实例支持数据共享，支持跨区域数据实时复制（CCR）
- 依托对象存储的增量备份及数据恢复机制



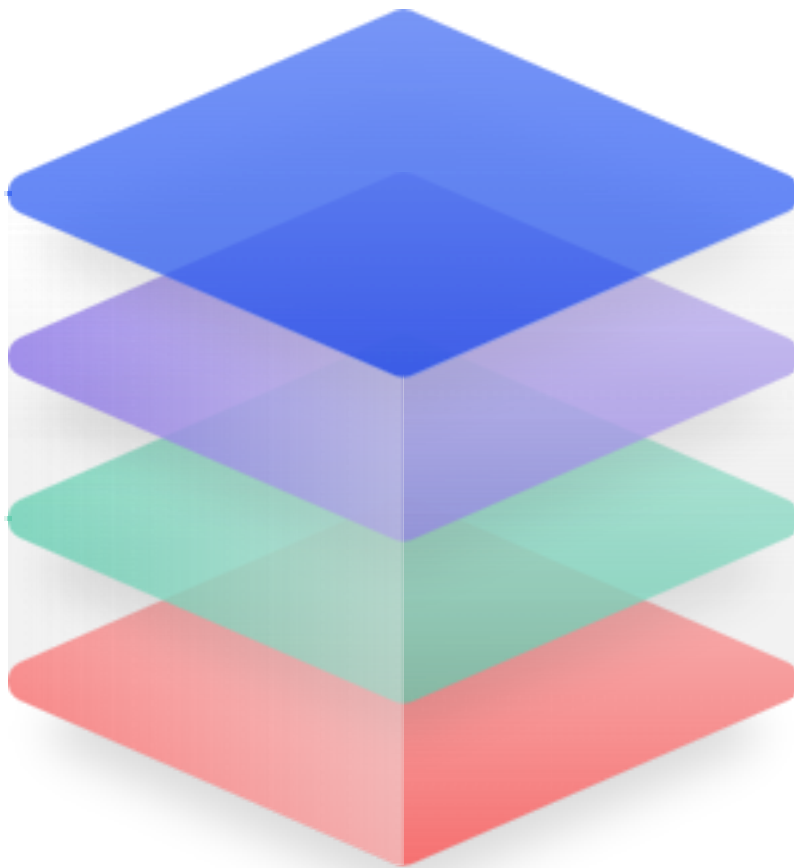
SELECTDB Cloud 应用场景

面向客户的报表与分析

面向管理者的驾驶舱Dashboard;
面向业务分析人员的Reporting;
面向C端/B端用户的高并发报表分析

用户画像与行为分析

收集用户属性与行为数据, 构建用户数据平台, 进行用户参与、留存和转化等行为分析, 以及人群洞察和人群圈选等画像分析



湖仓一体的现代数据平台

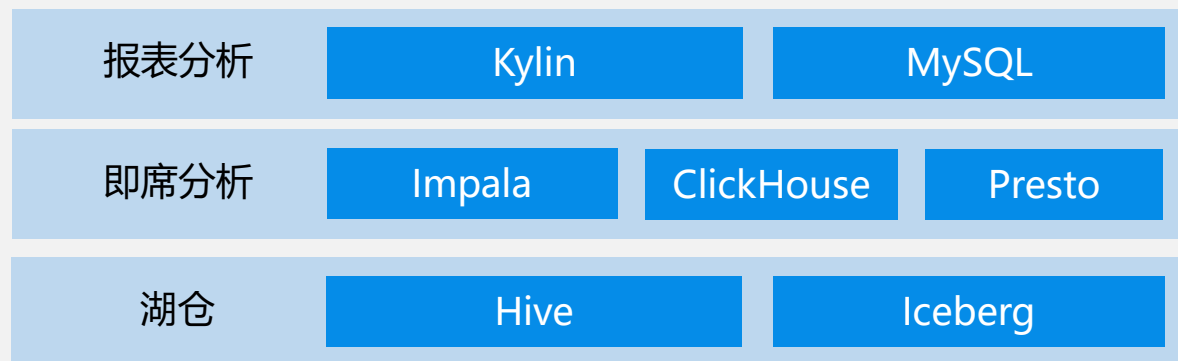
统一数据仓库和数据湖到统一平台, 提供面向企业内部的Bi报表和adhoc分析, 批量以及增量ELT

日志存储与分析

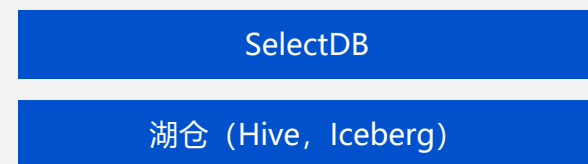
将业务、系统或者物联网相关的日志数据存储为结构化、半结构化或者原始文本, 构建统一的日志存储与分析平台。

典型应用场景 – 湖仓一体的现代化分析平台

之前的分析架构



基于SelectDB的统一分析架构



数据加工



1. 公司每天新增明细数据100亿，10TB单副本，大数据平台有数千台机器。
2. 旧的架构报表分析使用Kylin和MySQL，即席分析使用Impala, ClickHouse, Presto，将**报表分析**和**即席分析**场景用**SelectDB**，大大降低架构的复杂性，性价比提升**3倍**。
3. 使用SelectDB作为Hive和Iceberg的**湖仓查询加速引擎**，性能相比Presto有**3倍**提升，相比Hive有**10倍**以上提升。

典型应用场景 – 半结构化数据分析

业务应用场景

日志存储分析



用户行为分析
CND流量分析
服务故障分析

可观测性



系统资源监控
服务质量监测
APM性能监测

安全分析



网络流量安全监测
网络安全态势感知
安全事件追踪溯源

分析需求特点

写入吞吐大
存储周期长
快速检索

时效性要求高
历史数据聚合
监控指标多且扩展

写入吞吐大
快速检索
多种数据关联分析

业务价值

性能提升**2倍**以上

资源成本节省**60%+**

融合ES能力在统一数仓，**简化系统架构、减少维护成本**

典型应用场景 – 半结构化数据分析

业务价值

性能提升**2倍**以上

资源成本节省**60%+**

融合ES能力在统一数仓, **简化系统架构、减少维护成本**

快速
检索

复杂
分析

高性
价比

半结构化增强

倒排索引

动态表

冷热分离

Doris引擎

执行优化器

标准SQL

场景化存储模型

MPP架构

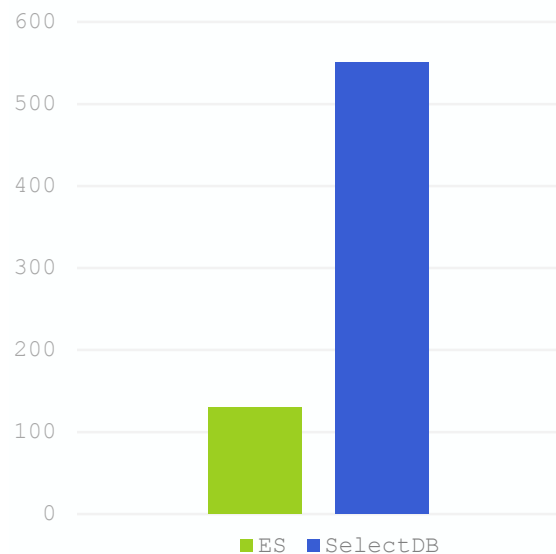
向量化执行

列式存储

典型应用场景 – 半结构化数据分析

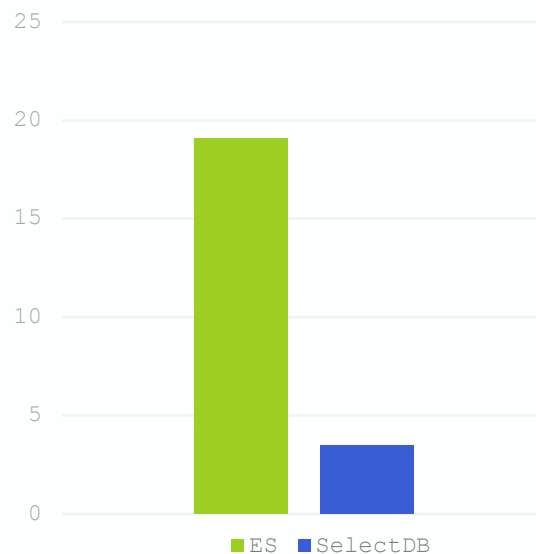
4.2倍

写入速度(MB/s)



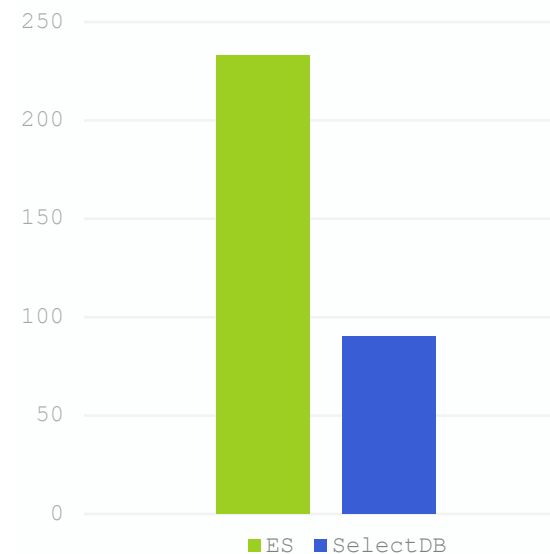
5.4倍

磁盘空间(GB)



2.3倍

查询时间(秒)



说明:

1. 以上数据为ES官方性能benchmark中http_logs场景测试结果
2. 写入速度越高越好, 磁盘空间越低越好, 查询时间越低越好

欢迎关注



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

Architect