

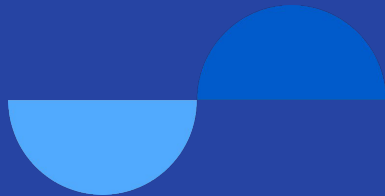


Databend

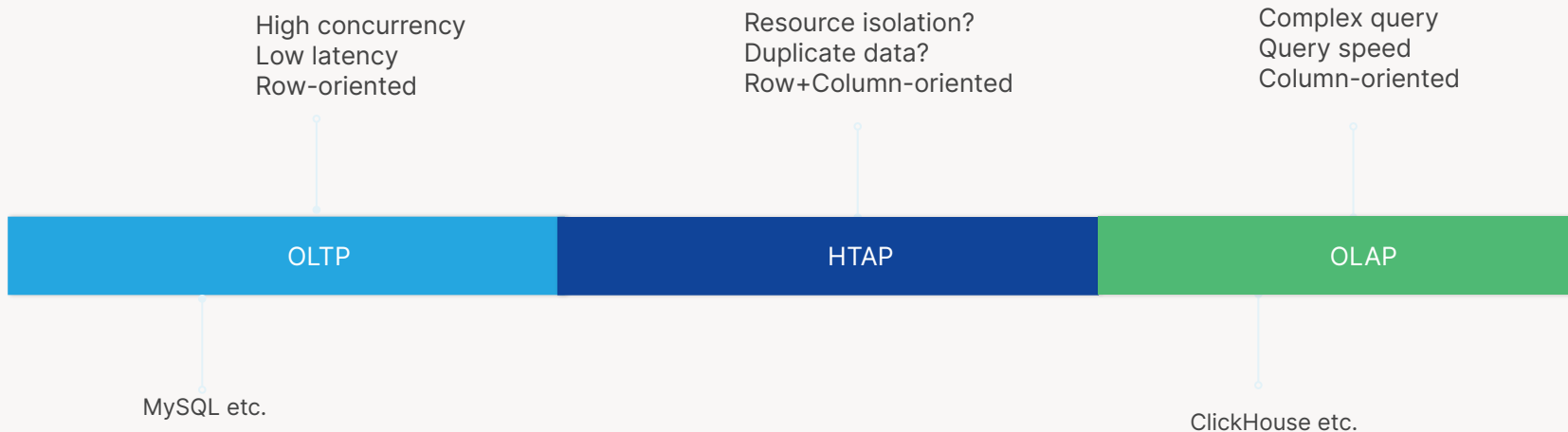
Databend

A modern warehouse with **Rust** for your massive-scale analytics

<https://github.com/datafuselabs/databend>



Database and Data Warehouse



大纲

- 大数据分析遇到了什么“新”问题？
- 传统数仓为什么无法解决这些“新”问题？
- 新一代实时弹性数仓如何设计？
- 使用 Rust 从 0 到 1 研发一款数仓是种什么体验？



Bohu TANG (张雁飞)

Co-Creator of Databend: <https://github.com/datafuselabs/databend>

ClickHouse and MySQL(TokuDB) 重度贡献者

Database Kernel | Distributed Database | Data Warehouse

<https://bohutang.me/>



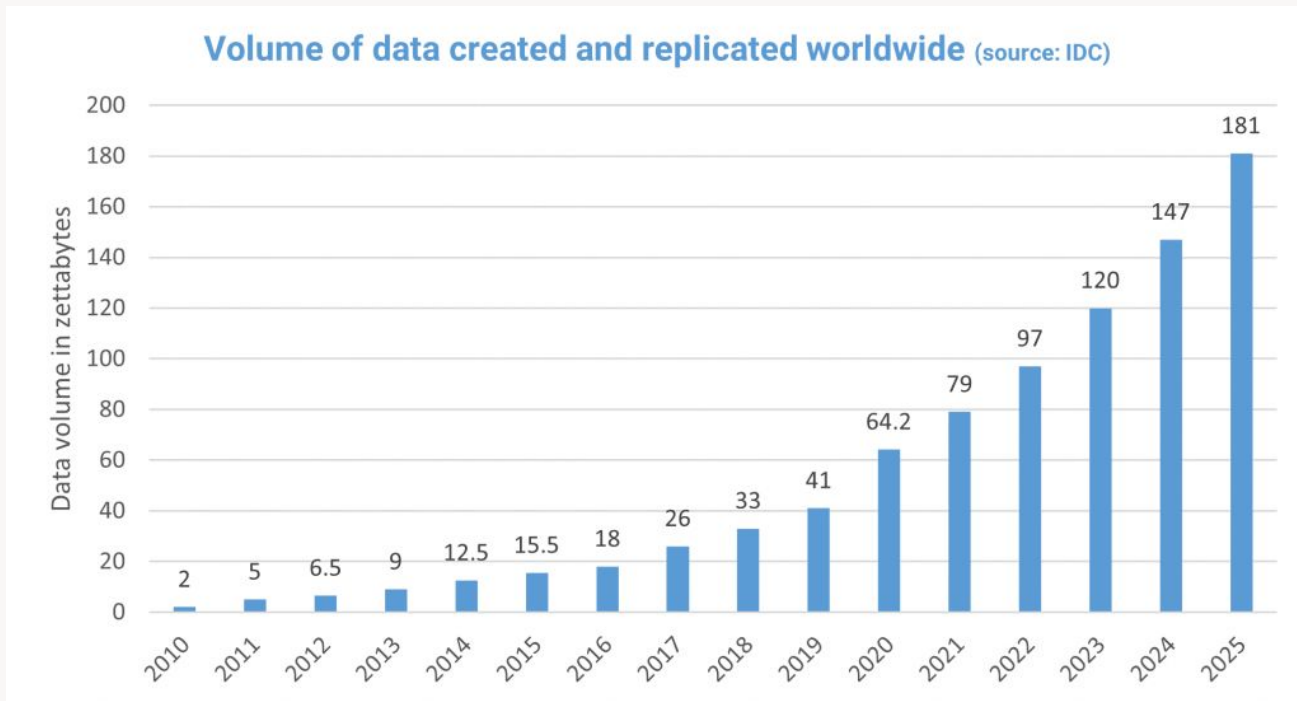


Databend

01

当今(2022)大数据新问题

全球数据指数级增长



1024PB = 1EB, 1024EB = 1ZB

当今大数据新问题 1

- 大数据量下的资源利用率问题, $< 50\%$
- 物理资源常驻问题
- 大数据分析, 波峰、波谷问题



当今大数据新问题 2

- 大数据量下的存储成本问题
- PB 级数据, 每月存储成本百万美金 !

▼ S3 Standard [Info](#)

The calculations below exclude Free Tier discounts.

S3 Standard storage

50000

TB per month

How will data be moved into S3 Standard?

Automatically calculates PUT, COPY, POST costs for moving data into S3 Standard initially. To compare the cost of current storage in S3 Standard to lifecycle in S3 Standard while selecting Lifecycle under the new storage class to capture the upfront cost of moving your data.

The specified amount of data is already stored in S3 Standard

PUT. COPY. POST. LIST requests to S3 Standard

Total Upfront cost: 0.00 USD

Total Monthly cost: 1,075,763.20 USD

Show Details ▼

当今大数据新问题 3

- 大数据量下的计算成本问题
- 对扫描数据量要求非常高, 容易破产

Configure Amazon Athena [Info](#)

Queries

Total number of queries

10 per day

Data amount scanned per query

1 TB

Spark

Total number of spark sessions

per day

Code execution per session

DPU-hour

▼ Show calculations

Unit conversions

Total number of queries: 10 per day * (730 hours in a month / 24 hours in a day) = 304.17 queries per month

Pricing calculations

Total Upfront cost: 0.00 USD

Total Monthly cost: 1,520.00 USD

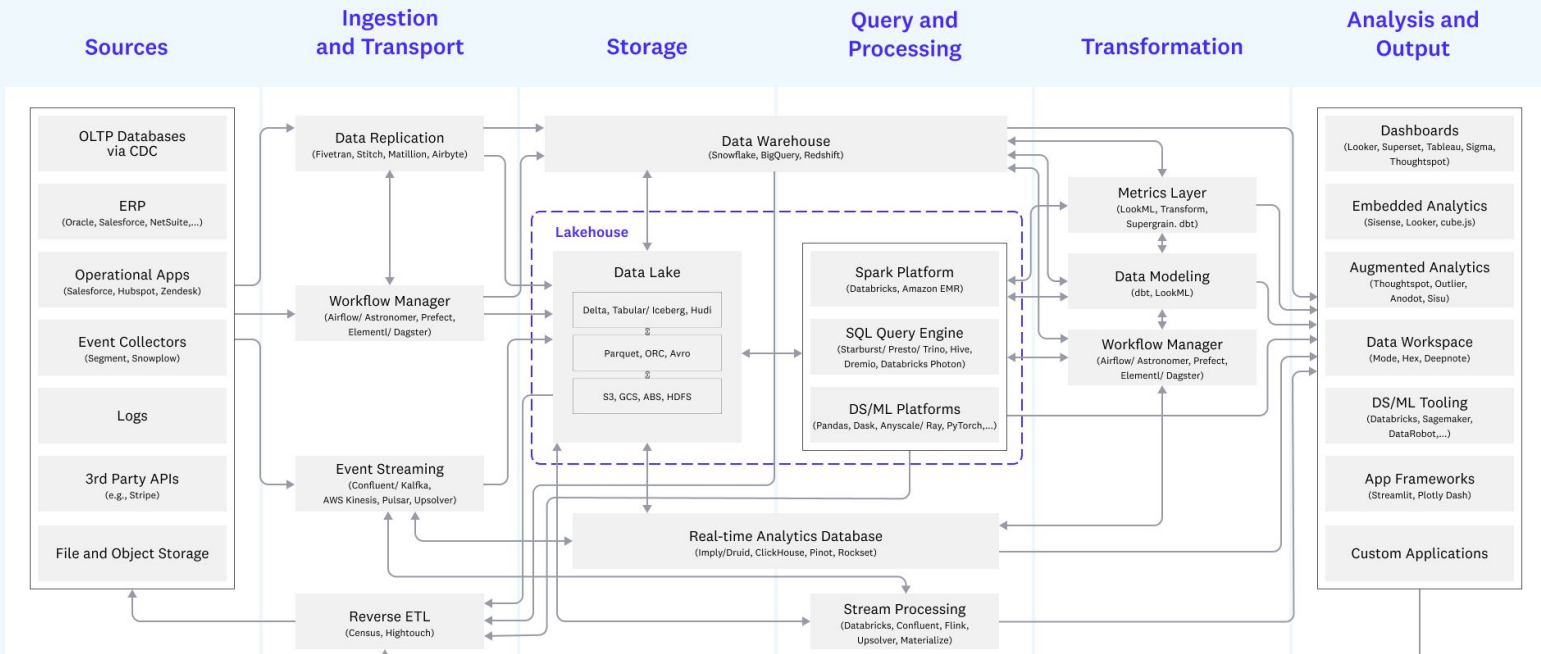
[Show Details ▼](#)

[Save and view summary](#)

当今大数据新问题 4

- 大数据量下的数据平台复杂度越来越高

Unified Data Infrastructure (2.0) (From a16z)





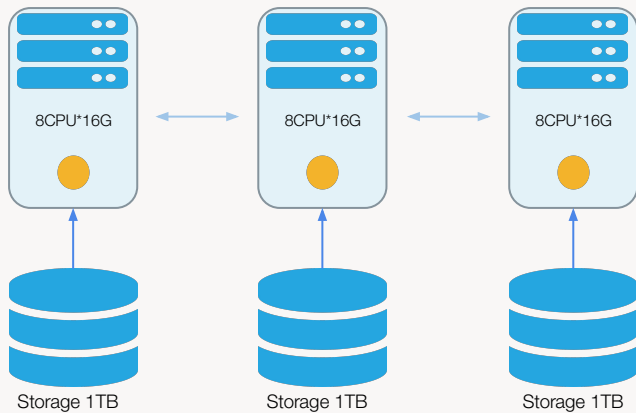
Databend

02

传统数仓架构 vs. 弹性数仓架构

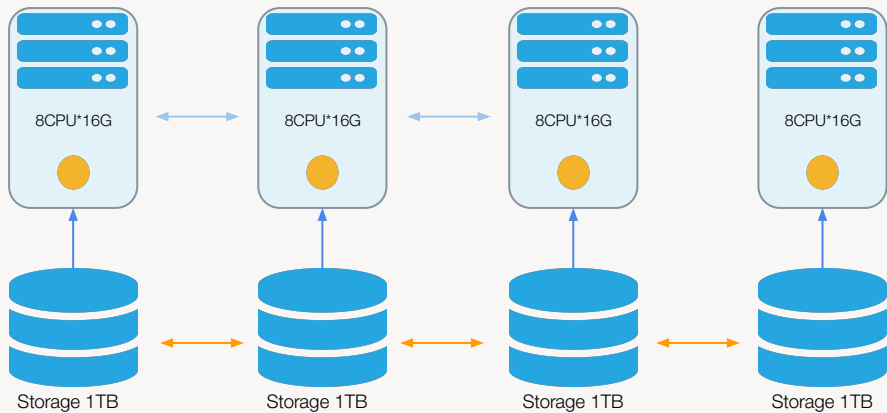
传统数仓架构

- Shared-Nothing
- 存储、计算一体
- 资源控制粒度粗



传统数仓架构

- Shared-Nothing
- 存储、计算一体
- 资源控制粒度粗



传统数仓架构

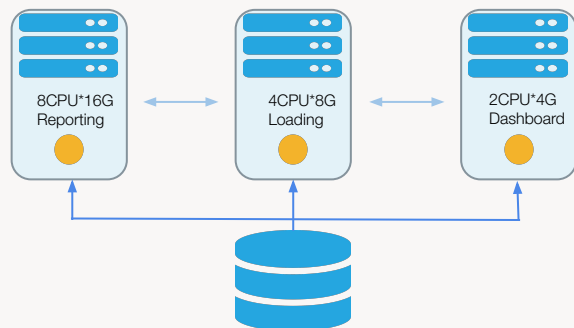
- Shared-Nothing - 弱弹性
- 存储、计算一体 - 弱弹性
- 资源控制粒度粗 - 成本高

$$\text{成本(高)} = \text{Resource} * \text{Time}$$



新一代弹性数仓架构

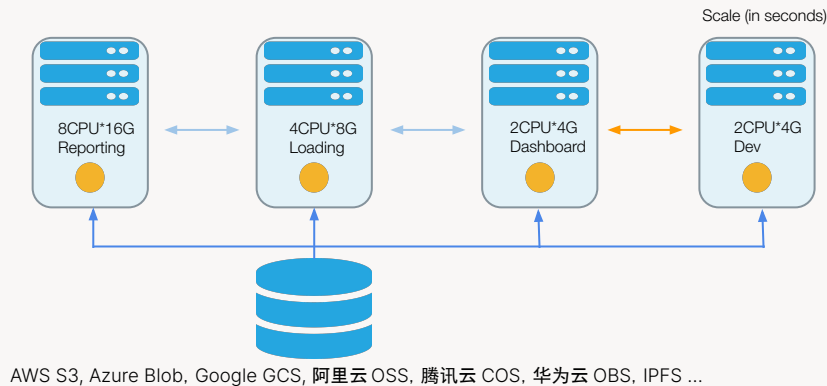
- Shared-Storage (Amazon S3, Azure Blob ...)
- 真正存储、计算分离
- 实时弹性扩容和缩容
- 资源控制粒度细



AWS S3, Azure Blob, Google GCS, 阿里云 OSS, 腾讯云 COS, 华为云 OBS, IPFS ...

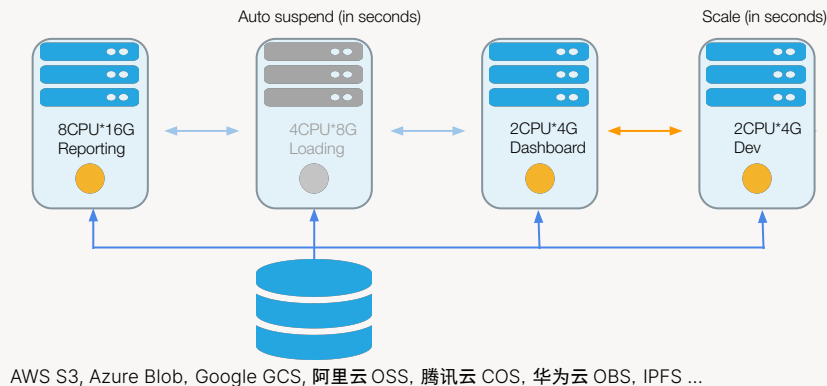
新一代弹性数仓架构

- Shared-Storage (Amazon S3, Azure Blob ...)
- 真正存储、计算分离
- 实时弹性扩容和缩容
- 资源控制粒度细



新一代弹性数仓架构

- Shared-Storage (Amazon S3, Azure Blob ...)
- 真正存储、计算分离
- 实时弹性扩容和缩容
- 资源控制粒度细



新一代弹性数仓架构

- Shared-Storage (Amazon S3, Azure Blob ...) - 高弹性
- 真正存储、计算分离 - 高弹性
- 实时弹性扩容和缩容 - 高弹性
- 资源控制粒度细 - 成本低

成本(低) = Resource * Time





Databend

03

Databend 新一代实时弹性架构设计

ClickHouse

- OS Warehouse
- 向量化计算, 细节优化到位
- Pipeline 处理器和调度器
- MergeTree 列式存储引擎
- 单机性能非常强悍
- 缺点: 分布式能力较弱, 运维复杂度高

[ClickHouse Group By 为什么这么快]: <https://bohutang.me/2021/01/21/clickhouse-and-friends-groupby/>

[ClickHouse Pipeline 处理器和调度器]: <https://bohutang.me/2020/06/11/clickhouse-and-friends-processor/>

[ClickHouse 存储引擎技术进化与MergeTree]: <https://bohutang.me/2020/06/20/clickhouse-and-friends-merge-tree-algo/>



Snowflake

- Cloud Warehouse
- 多租户, 存储、计算分离
- 基于对象存储便宜介质
- 弹性能力非常强悍
- 缺点: 单机性能一般, 重度依赖分布式



Databend = ClickHouse + Snowflake + Rust



Databend

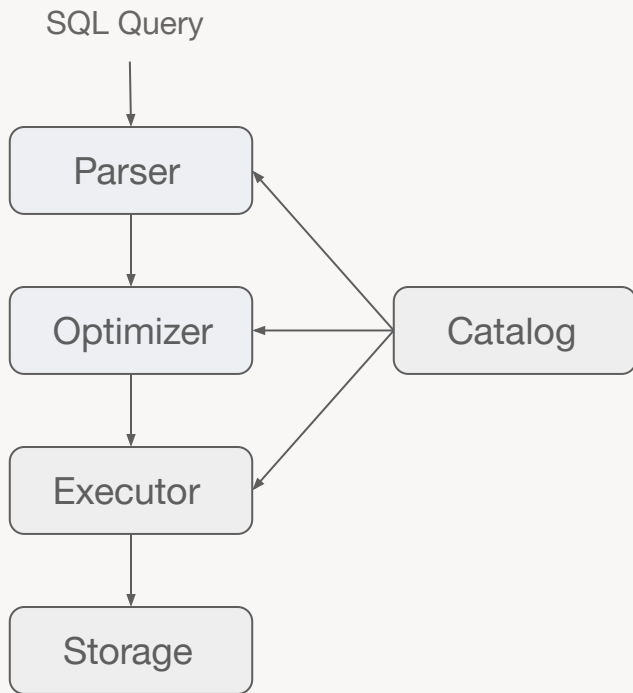
- 借鉴 ClickHouse 向量化计算, 提升单机计算性能
- 借鉴 Snowflake 存储、计算分离思想, 提升分布式计算能力
- 借鉴 Git, MVCC 列式存储引擎, Insert/Read/Delete/Update(WIP)
- 高弹性 + 强分布式, 致力于解决大数据分析**成本**和**复杂度**问题
- 基于便宜的对象存储也能方便的做实时性分析
- 完全使用 Rust 研发(30w+ loc), Day1 在 Github 开源

挑战

- 计算层做到高弹性，计算的状态如何管理？
- 对象存储不是为数据库而设计，高延迟和高性能如何平衡？
- 如何让系统更加智能，根据查询模式自动创建索引？
- 如何面向 Warehouse + Datalake 需求设计？

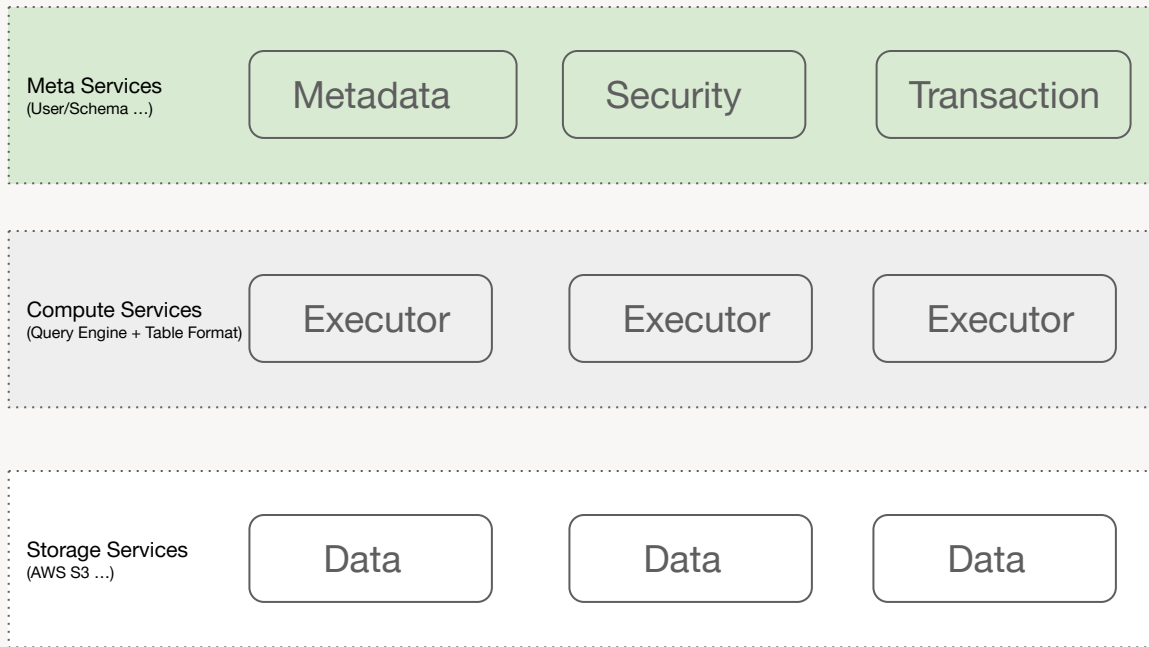


单机功能模块

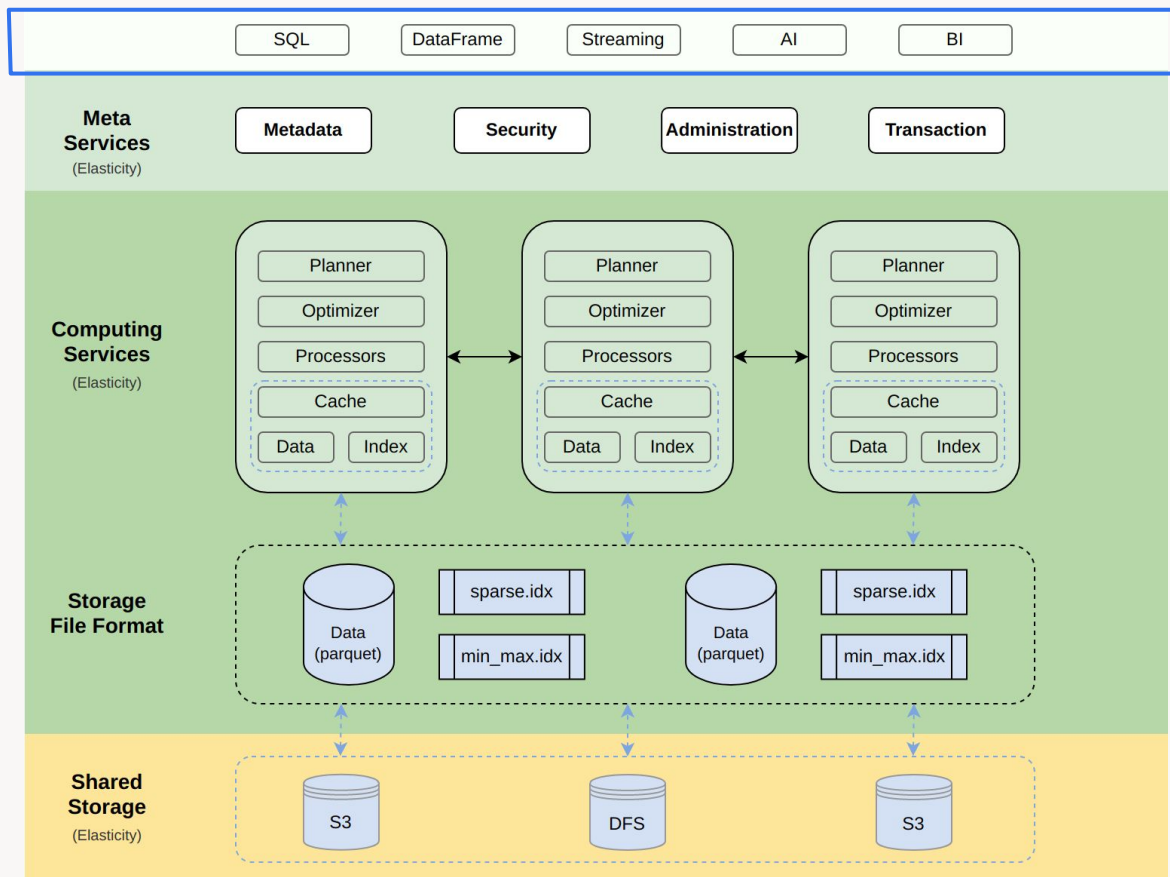


模块微服务化

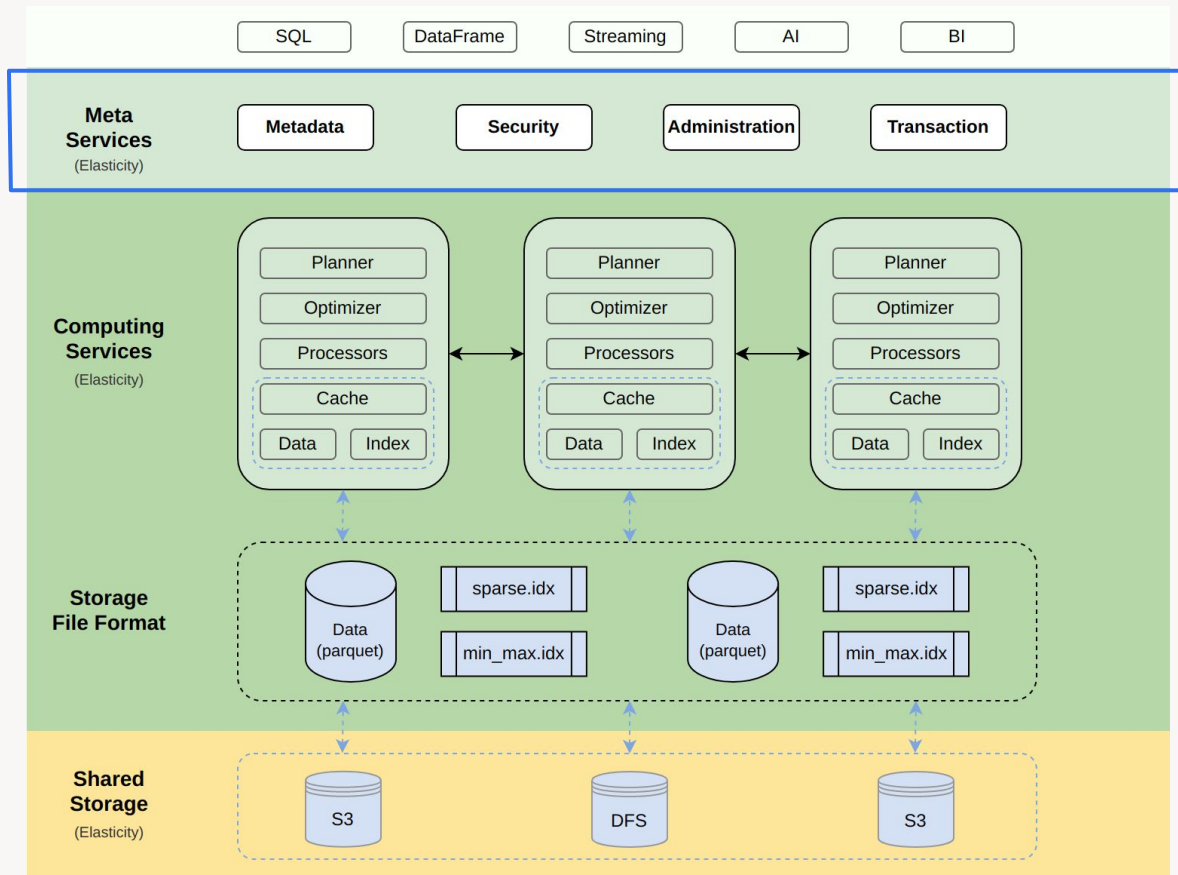
SQL Query



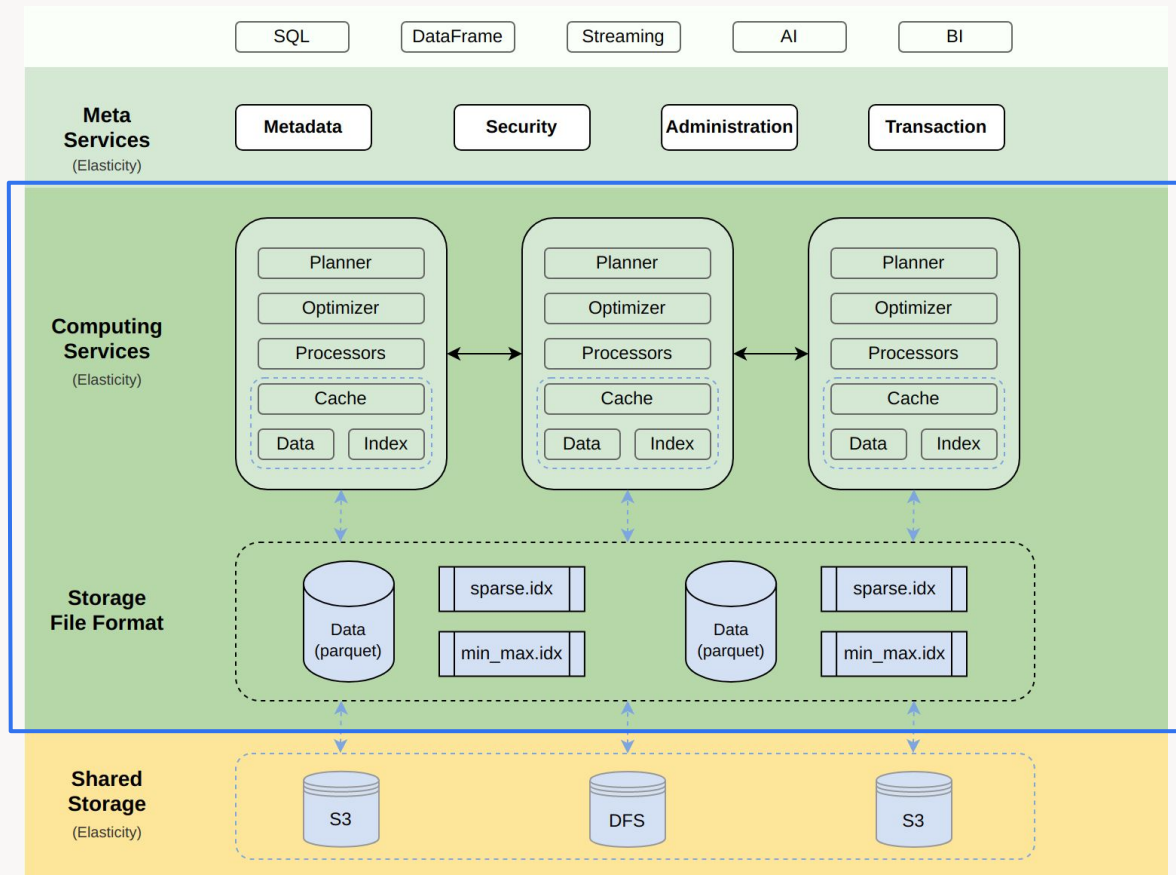
Databend 架构



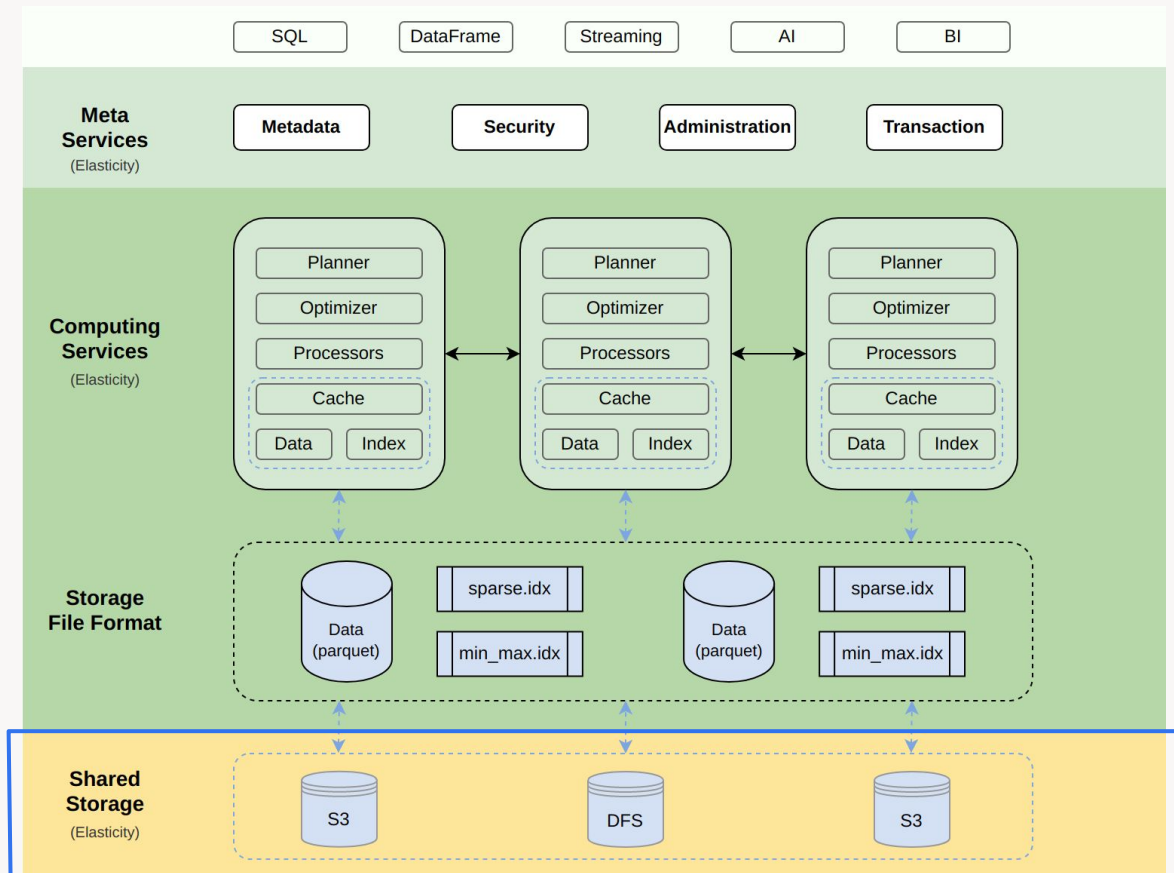
Databend 架构



Databend 架构

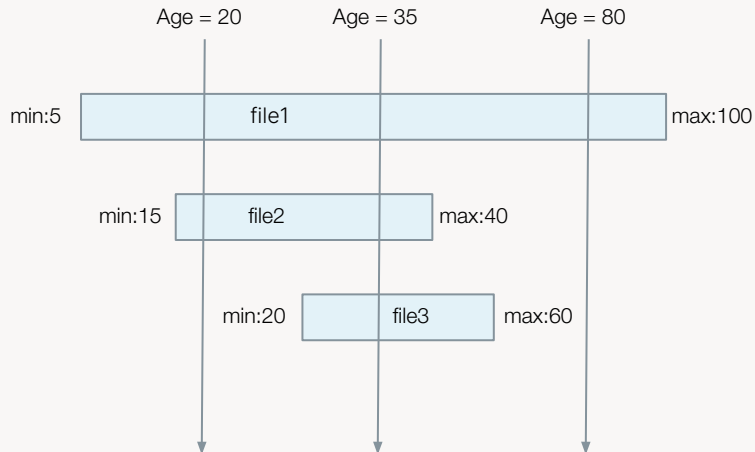


Databend 架构

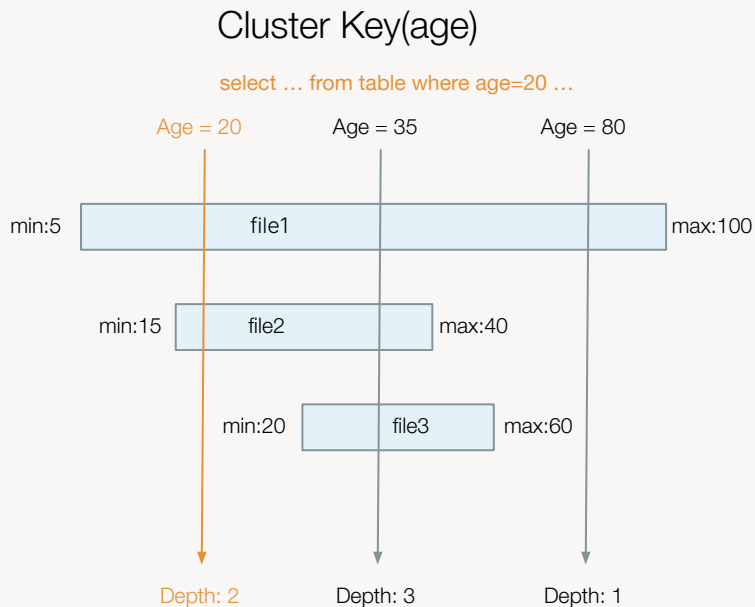


Automatic Tuning

Cluster Key(age)



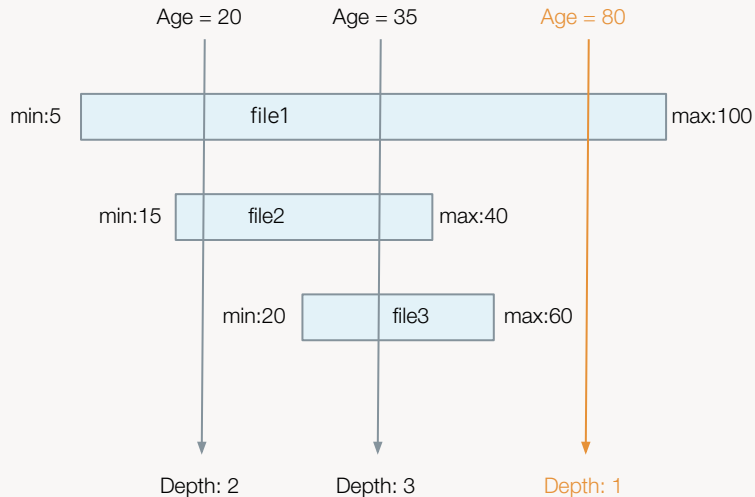
Automatic Tuning



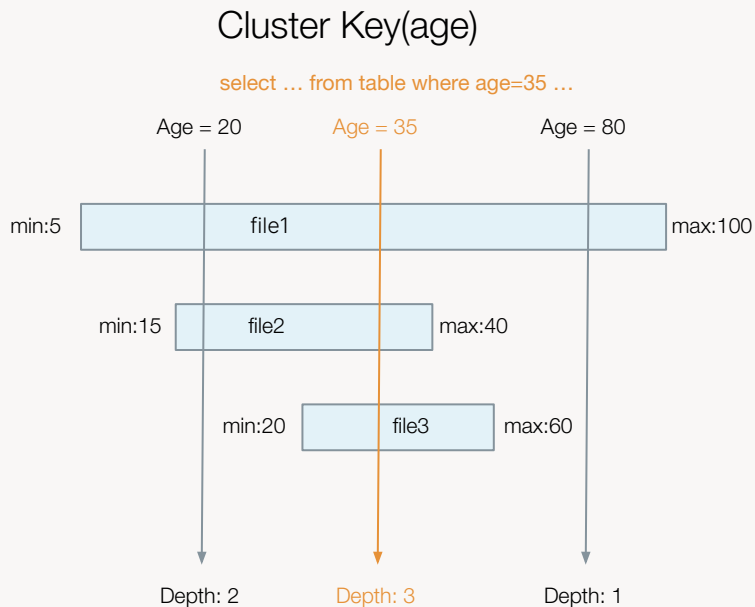
Automatic Tuning

Cluster Key(age)

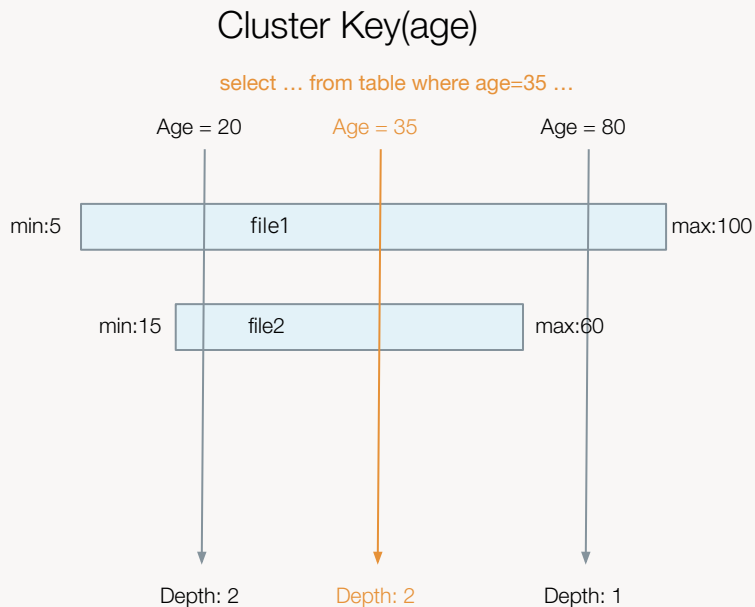
`select ... from table where age=80 ...`



Automatic Tuning



Automatic Tuning



Databend + Hive

```
CREATE CATALOG my_hive  
  TYPE=HIVE  
  CONNECTION = (URL='<hive-meta-store>'  
THRIFT_PROTOCOL=BINARY);  
SELECT * FROM my_hive.db1.table;
```

[Multiple Catalog RFC]: <https://databend.rs/doc/contributing/rfcs/multiple-catalog>



Databend + Iceberg

```
CREATE CATALOG my_iceberg  
  TYPE=ICEBERG  
  CONNECTION = (URL='s3://my_bucket/path/to/iceberg');  
SELECT * FROM my_iceberg.db1.table;
```

[Multiple Catalog RFC]: <https://databend.rs/doc/contributing/rfcs/multiple-catalog>





Databend

04

Databend 开源社区

Databend 开源社区



4.9K Stars

140+ Contributors

迭代非常快

Overview

350 Active pull requests

279 Active issues

341

Merged pull requests

9

Open pull requests

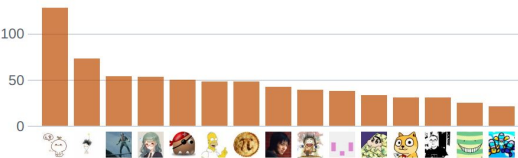
188

Closed issues

91

New issues

Excluding merges, **39 authors** have pushed **845 commits** to main and **845 commits** to all branches. On main, **1,552 files** have changed and there have been **73,585 additions** and **38,579 deletions**.



36 Releases published by 1 person

<https://github.com/datafuselabs/databend>

Databend 开源社区



~40 月度活跃开发者:

SAP

Yahoo

Fortinet

Shopee

PingCAP

Alibaba

Tencent

ByteDance

EMQ

快手 (湖仓一体共建)

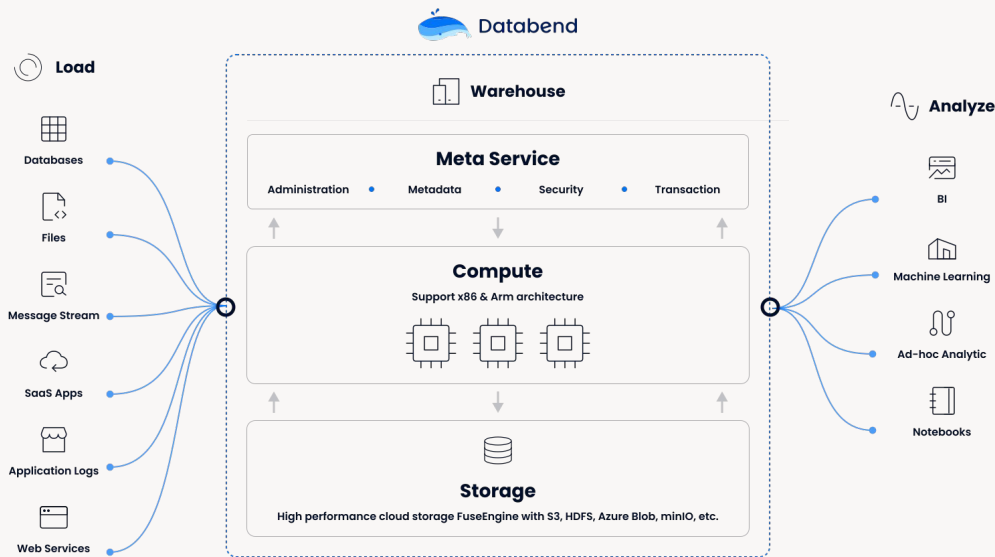
... ..



Databend 体验: On-Premises, Serverless



- On-Premises
社区版: <https://databend.rs>
- Serverless Cloud
海外(AWS) <https://app.databend.com>
国内(阿里云) <https://app.databend.cn>



Databend 用户



More ...





Databend

Thanks

