# Unbundling of the DBMS stack

*Mosha Pasumansky & Benjamin Wagner*

# DBMSs created per Year

| | | |
|---|---|---|
| **CeresDB** | Time series database | Query Engine - DataFusion<br>WAL - RocksDB, OceanBase, Memtable - AgateDB<br>SST - derived from Parquet |
| **LingoDB** | Data processing system that leverages compiler technology | Parser - libpg_query |
| **CnosDB** | Time series database | RPC - ArrowFlight, Query Engine - DataFusion |
| **RisingWave** | Distributed SQL for stream processing | DataFusion |
| **MonographDB** | Multi model database | Compute - MariaDB, Storage - Cassandra |
| **nucliadb** | AI search / generative answers / vector database | LMBD and/or TiKV |
| **spicedb** | Database for managing security permissions checking | CockroachDB |
| **Dragonfly** | Redis replacement | |
| **Oriole** data base | Next gen storage engine for PostgreSQL | PostgreSQL extension |
| **Edgeless**DB | Database for confidential computing (inside SGX enclave) | Forked MariaDB, Storage engine - RocksDB |
| **NEON** | Serverless PostgreSQL | PostgreSQL |
| **FerretDB** | MongoDB alternative | PostgreSQL |

**Admission Control**

**Dispatch and Scheduling**

**Process Manager (Section 2)**

**Local Client Protocols**

**Remote Client Protocols**

**Client Communications Manager**

**Query Parsing and Authorization**

**Query Rewrite**

**Query Optimizer**

**Plan Executor**

**DDL and Utility Processing**

**Relational Query Processor (Section 4)**

**Access Methods**

**Buffer Manager**

**Lock Manager**

**Log Manager**

**Transactional Storage Manager (Sections 5 & 6)**

**Catalog Manager**

**Memory Manager**

**Administration, Monitoring & Utilities**

**Replication and Loading Services**

**Batch Utilities**

**Shared Components and Utilities (Section 7)**

Fig. 1.1 Main components of a DBMS.

Foundations and Trends® in Databases 1:2 (2007)

**Architecture of a Database System**

Joseph M. Hellerstein, Michael Stonebraker, and James Hamilton

now

the essence of knowledge

Fig. 1.1 Main components of a DBMS.

**Client Protocols**

**Process Manager**

**Admission Control**

**Dispatch And Scheduling**

**Query Processor**

**Language Frontend**

**Planner / Optimizer**

**Local Execution**

**Distributed Execution**

**DDL and Utility Processing**

**Storage Engine**

**Data Formats, Indexes**

**Buffer Manager**

**Distributed Storage**

**Log Manager**

**Cloud Services**

**Authorization**

**Catalog Manager**

**Transaction Manager**

Query Interfaces
Optimizers
Execution Engines
Storage Engines

Credit: Towards a Modular Data Management System Framework
https://cdmsworkshop.github.io/2022/Slides/Fri_C2.5_HaralamposGavriilidis.pdf

**Client Protocols**

**Language Frontend**

**Planner / Optimizer**

**Local Execution**

**Distributed Execution**

**DDL and Utility Processing**

**Query Processor**

**Admission Control**

**Dispatch And Scheduling**

**Process Manager**

**Cloud Services**

**Authorization**

**Catalog Manager**

**Transaction Manager**

**Data Formats, Indexes**

**Buffer Manager**

**Distributed Storage**

**Log Manager**

**Storage Engine**

**Process Manager**

Admission Control

Dispatch And Scheduling

**Client Protocols**

**Query Processor**

Language Frontend

Planner / Optimizer

Local Execution

Distributed Execution

DDL and Utility Processing

**Storage Engine**

Data Formats, Indexes

Buffer Manager

Distributed Storage

Log Manager

Cloud Services

Authorization

Catalog Manager

Transaction Manager

| Process Manager | Query Processor / Storage Engine | | |
|---|---|---|---|

**Process Manager**

- Admission Control
- Dispatch And Scheduling

**Client Protocols**

**Query Processor**

- Language Frontend
- Planner / Optimizer
- Local Execution
- Distributed Execution
- DDL and Utility Processing

**Storage Engine**

- Data Formats, Indexes
- Buffer Manager
- Distributed Storage
- Log Manager

- Cloud Services
- Authorization
- Catalog Manager
- Transaction Manager

1. What language ?

2. What language ?

1. What language ?



2. What language ?

1. What language ?



2. What language ?

**Process Manager**

Admission Control

Dispatch And Scheduling

**Client Protocols**

**Query Processor**

Language Frontend

Planner / Optimizer

Local Execution

Distributed Execution

DDL and Utility Processing

**Storage Engine**

Data Formats, Indexes

Buffer Manager

Distributed Storage

Log Manager

Cloud Services

Authorization

Catalog Manager

Transaction Manager

**Query Processor**

| Client Protocols | |
|---|---|
| Language Frontend | |
| Planner / Optimizer | DDL and Utility Processing |
| Local Execution | |
| Distributed Execution | |

**Admission Control**

**Dispatch And Scheduling**

**Process Manager**

**Storage Engine**

| Data Formats, Indexes | Buffer Manager |
|---|---|
| Distributed Storage | Log Manager |

**Cloud Services**

**Authorization**

**Catalog Manager**

**Transaction Manager**

**Client Protocols**

**Process Manager**

**Admission Control**

**Dispatch And Scheduling**

**Query Processor**

**Language Frontend**

**Planner / Optimizer**

**Local Execution**

**Distributed Execution**

**DDL and Utility Processing**

**Storage Engine**

**Data Formats, Indexes**

**Buffer Manager**

**Distributed Storage**

**Log Manager**

**Cloud Services**

**Authorization**

**Catalog Manager**

**Transaction Manager**

SQLite

DuckDB

DATA FUSION

Velox

ClickHouse

| Process Manager | Query Processor | | | Storage Engine | |
|---|---|---|---|---|---|
| **Admission Control** | Client Protocols | | | **Cloud Services** | |
| | Language Frontend | | | **Authorization** | |
| | Planner / Optimizer | DDL and Utility Processing | | **Catalog Manager** | |
| **Dispatch And Scheduling** | **Local Execution** | | | | |
| | Distributed Execution | | | **Transaction Manager** | |
| | Data Formats, Indexes | Buffer Manager | | | |
| | Distributed Storage | Log Manager | | | |

SQLite

DuckDB

DATA FUSION

Velox

ClickHouse

ARROW    Unbundled Execution Engine ?

| Admission Control | Client Protocols | | Cloud Services |
| | Language Frontend | | |
| | Planner / Optimizer | DDL and Utility Processing | Authorization |
| | Local Execution | | Catalog Manager |
| Dispatch And Scheduling | Distributed Execution | | |
| | Query Processor | | Transaction Manager |
| | Data Formats, Indexes | Buffer Manager | |
| | Distributed Storage | Log Manager | |
| Process Manager | Storage Engine | | |

| Process Manager | Client Protocols | Cloud Services |
|---|---|---|

**Client Protocols**

**Admission Control**

**Dispatch And Scheduling**

**Query Processor**

**Language Frontend**

**Planner / Optimizer**

**Local Execution**

**DDL and Utility Processing**

**Distributed Execution**

**Storage Engine**

**Data Formats, Indexes**

**Buffer Manager**

**Distributed Storage**

**Log Manager**

**Cloud Services**

**Authorization**

**Catalog Manager**

**Transaction Manager**

**Process Manager**

**Storage Engine**

**Figure 1: Open source modular data stack outline.**

The Composable Data Management System Manifesto
https://dl.acm.org/doi/pdf/10.14778/3603581.3603604

**Figure 1: Open source modular data stack outline.**

The Composable Data Management System Manifesto
https://dl.acm.org/doi/pdf/10.14778/3603581.3603604



Substrait: Rethinking DBMS Composability
https://cdmsworkshop.github.io/2022/Proceedings/Keynotes/Abstract_JacquesNadeau.pdf

## Layered Query Compilation

SQL

↓ Parser

Logical Plan

↓ Query Opt.

Physical Plan

↓ Translation

**Scope of Prior Work**

Opt. ↻ IR 1

↓ (N-1) Lowerings

Opt. ↻ IR N

↓ Translation

Backend

## Our Open Query Compilation Stack

SQL

↓ Parser & Semantic Analysis

Query Opt. ↻ [ IR 1 | IR 2 | IR X ] ↺ Cross Opt.

↓ Lower IR 1

Opt. ↻ [ IR 3 | IR 2 | IR Y ]

↓ Progressive Lowerings

Low-Level IR

↓ Translation

Backend

■ relational ■ DB-specific □ general-purpose ■ other domains

**Figure 1: Our proposal of an *open* query compilation stack. It enhances prior work on layered query compilation with two major ideas: 1) Introducing open IRs, designed to be combined with other IRs, and 2) implementing query optimization as compiler passes.**
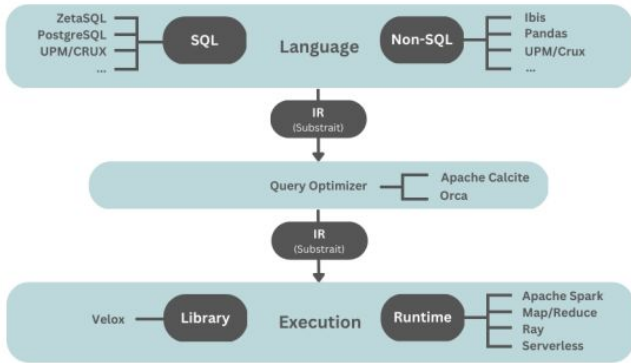
Designing an Open Framework for Query Optimization and Compilation
https://db.in.tum.de/~jungmair/papers/p2485-jungmair.pdf

| | Client Protocols | | |
|---|---|---|---|
| Admission Control | Language Frontend | | Cloud Services |
| | Planner / Optimizer | DDL and Utility Processing | Authorization |
| Dispatch And Scheduling | Local Execution | | Catalog Manager |
| | Distributed Execution | | Transaction Manager |
| | Query Processor | | |
| | Data Formats, Indexes | Buffer Manager | |
| | Distributed Storage | Log Manager | |
| Process Manager | Storage Engine | | |

**Client Protocols**

**Query Processor**

Language Frontend

Planner / Optimizer

Local Execution

Distributed Execution

DDL and Utility Processing

**Process Manager**

Admission Control

Dispatch And Scheduling

**Storage Engine**

Data Formats, Indexes

Buffer Manager

Distributed Storage

Log Manager

Cloud Services

Authorization

Catalog Manager

Transaction Manager

PostgreSQL Tools

ODBC / JDBC

PostgreSQL queries

CitusDB Metadata

Distributed Query Planner

PostgreSQL Query Planner

Distributed Query Executors

CitusDB Metadata (Master) Node

PostgreSQL extensions for full cooperation

Distributed queries

CitusDB Worker Node

Distributed Join Funcs (optional)

PostgreSQL Query Planner

PostgreSQL Query Executor

Storage Engine

PostgreSQL extensions for full cooperation

Distributed Join Funcs (optional)

PostgreSQL Query Planner

PostgreSQL Query Executor

Storage Engine

Distributed Join Funcs (optional)

PostgreSQL Query Planner

PostgreSQL Query Executor

Storage Engine

Citus Query Processing
https://docs.citusdata.com/en/v7.0/performance/query_processing.html

Admission Control

Dispatch And Scheduling

Process Manager

Client Protocols

Language Frontend

Planner / Optimizer

Local Execution

Distributed Execution

DDL and Utility Processing

Query Processor

Data Formats, Indexes

Buffer Manager

Distributed Storage

Log Manager

Storage Engine

Cloud Services

Authorization

Catalog Manager

Transaction Manager

Generic pushdown mechanism in YugabyteDB

1. Build PostgreSQL execution tree

2. Identify subtree to pushdown

3. Identify nodes to send pushdown to

4. Execute subtree using PostgreSQL libraries on DocDB storage layer

SQL Query — CLIENT

YSQL Query Layer

DocDB Storage Layer

Node #1    Node #2    Node #3

5 Query Pushdowns for Distributed SQL and How They Differ from a Traditional RDBMS
https://www.yugabyte.com/blog/5-query-pushdowns-for-distributed-sql-and-how-they-differ-from-a-traditional-rdbms/

Admission Control

Dispatch And Scheduling

Process Manager

Client Protocols

Language Frontend

Planner / Optimizer

Local Execution

Distributed Execution

DDL and Utility Processing

Query Processor

Data Formats, Indexes

Buffer Manager

Distributed Storage

Log Manager

Storage Engine

Cloud Services

Authorization

Catalog Manager

Transaction Manager

Spark Cluster Mode Overview
https://spark.apache.org/docs/latest/cluster-overview.html

**Worker Node**
Executor
Cache
Task
Task

**Driver Program**
SparkContext

**Cluster Manager**

**Worker Node**
Executor
Cache
Task
Task



Admission Control

Dispatch And Scheduling

Process Manager

Client Protocols

Language Frontend

Planner / Optimizer

Local Execution

Distributed Execution

DDL and Utility Processing

Query Processor

Data Formats, Indexes

Buffer Manager

Distributed Storage

Log Manager

Storage Engine

Cloud Services

Authorization

Catalog Manager

Transaction Manager

Spark Cluster Mode Overview
https://spark.apache.org/docs/latest/cluster-overview.html



**Figure 8: Presto on Spark architecture**

Presto: A Decade of SQL Analytics at Meta
https://research.facebook.com/publications/presto-a-decade-of-sql-analytics-at-meta/

Spark shuffle recap

Mappers | Map Output Files | Reducers



Client Protocols

Admission Control

Dispatch And Scheduling

Language Frontend

Planner / Optimizer

Local Execution

Distributed Execution

DDL and Utility Processing

Query Processor

Data Formats, Indexes

Buffer Manager

Distributed Storage

Log Manager

Storage Engine

Process Manager

Cloud Services

Authorization

Catalog Manager

Transaction Manager

# Spark shuffle recap

**Mappers**  **Map Output Files**  **Reducers**

Partition

Map 0 → Reduce 0

Map 1 → Reduce 1

Map *m* → Reduce *r*

- InBuilt Shuffle service
- Standalone Shuffle Service
- YARN Shuffle Service
- Mesos Shuffle Service
- Kubernetes Shuffle Service
- Cosco (Meta)
- Magnet (Linkedin)
- Riffle (Meta)
- Zeus (Uber)
- EMR Remote Shuffle Service (Alibaba)

| Admission Control | Client Protocols | | Cloud Services |
| | Language Frontend | | |
| | Planner / Optimizer | DDL and Utility Processing | Authorization |
| Dispatch And Scheduling | Local Execution | | |
| | Distributed Execution | | Catalog Manager |
| | Query Processor | | |
| | Data Formats, Indexes | Buffer Manager | Transaction Manager |
| | Distributed Storage | Log Manager | |
| Process Manager | Storage Engine | | |

**(a)** Monolithic shuffle systems.

**(b)** Exoshuffle.

**Figure 1:** Exoshuffle builds on an extensible architecture. Shuffle as a library is easier to develop and more flexible to integrate with applications. The data plane ensures performance and reliability.

**Exoshuffle: An Extensible Shuffle Architecture**

https://arxiv.org/abs/2203.05072

**Process Manager**

Admission Control

Dispatch And Scheduling

**Query Processor**

Client Protocols

Language Frontend

Planner / Optimizer

Local Execution

Distributed Execution

DDL and Utility Processing

**Storage Engine**

Data Formats, Indexes

Buffer Manager

Distributed Storage

Log Manager

Cloud Services

Authorization

Catalog Manager

Transaction Manager

**Delta Lake**

**Iceberg**

**Apache Hudi**



| | Client Protocols | |
|---|---|---|
| **Admission Control** | **Language Frontend** | **Cloud Services** |
| | Planner / Optimizer | **Authorization** |
| **Dispatch And Scheduling** | Local Execution / DDL and Utility Processing | **Catalog Manager** |
| | Distributed Execution | |
| | Query Processor | **Transaction Manager** |
| | Data Formats, Indexes / Buffer Manager | |
| | Distributed Storage / Log Manager | |
| Process Manager | Storage Engine | |

DELTA LAKE

ICEBERG

Apache hudi

Overview

Iceberg Catalog

db1.table1
current metadata pointer

metadata layer

metadata file
s0

metadata file
s0  s1

manifest list

manifest list

manifest file

manifest file

manifest file

data layer

data files

data files

data files

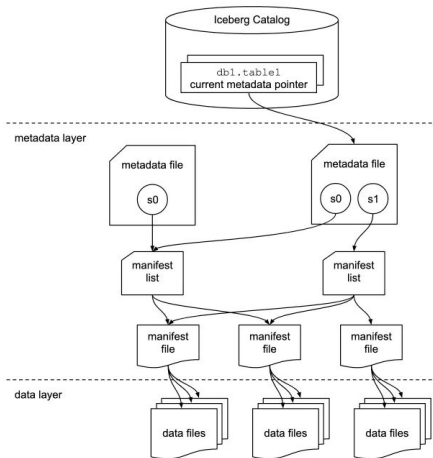mytable/date=2020-01-01/1b8a32d2ad.parquet
                    /a2dc5244f7.parquet
        /date=2020-01-02/f52312dfae.parquet
                    /ba68f6bd4f.parquet

Data objects
(partitioned
by date field)

/_delta_log/000001.json
        /000002.json
        /000003.json
        /000003.parquet
        /000004.json
        /000005.json
        /_last_checkpoint
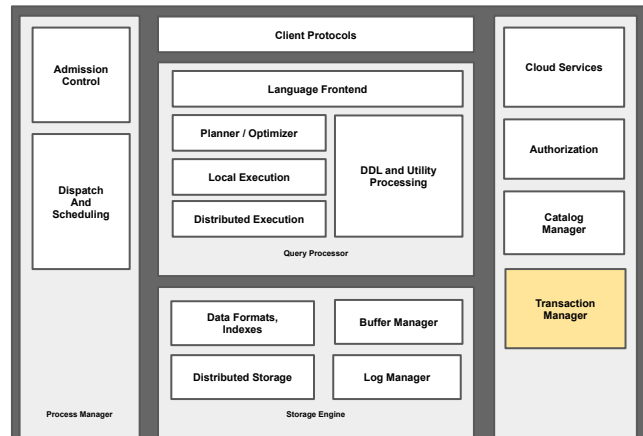
Log records
& checkpoints

Contains {version: "000003"}

Combines log
records 1 to 3

Transaction's operations, e.g.,
add date=2020-01-01/a2dc5244f7f7.parquet
add date=2020-01-02/ba68f6bd4f1e.parquet

| Client Protocols | | | |
|---|---|---|---|
| Admission Control | Language Frontend | | Cloud Services |
| | Planner / Optimizer | DDL and Utility Processing | Authorization |
| Dispatch And Scheduling | Local Execution | | Catalog Manager |
| | Distributed Execution | | |
| | Query Processor | | Transaction Manager |
| | Data Formats, Indexes | Buffer Manager | |
| | Distributed Storage | Log Manager | |
| Process Manager | Storage Engine | | |

# Snowflake Architecture



Authentication & access control

Infrastructure manager | Optimizer | Transaction Manager | Security

Cloud services
Compilation and Management

Metadata Key-Value Store
FoundationDB

Virtual warehouse — Cache

Compute
Virtual warehouses

Storage
Databases

**FOUNDATIONDB**



| Admission Control | Client Protocols | | Cloud Services |
|---|---|---|---|
| | Language Frontend | | Authorization |
| | Planner / Optimizer | DDL and Utility Processing | Catalog Manager |
| Dispatch And Scheduling | Local Execution | | Transaction Manager |
| | Distributed Execution | | |
| | Query Processor | | |
| | Data Formats, Indexes | Buffer Manager | |
| | Distributed Storage | Log Manager | |
| Process Manager | Storage Engine | | |

ScalarDB: Universal Transaction Manager for Polystores
https://github.com/scalar-labs/scalardb

**Process Manager**

- Admission Control
- Dispatch And Scheduling

**Query Processor**

- Client Protocols
- Language Frontend
- Planner / Optimizer
- Local Execution
- Distributed Execution
- DDL and Utility Processing

**Storage Engine**

- Data Formats, Indexes
- Buffer Manager
- Distributed Storage
- Log Manager

- Cloud Services
- Authorization
- Catalog Manager
- Transaction Manager

RocksDB 26

levelDB 10

Pebble 1

SQLite 13

LMDB 12



| | Client Protocols | |
|---|---|---|
| Admission Control | Language Frontend | Cloud Services |
| | Planner / Optimizer    DDL and Utility Processing | Authorization |
| Dispatch And Scheduling | Local Execution | Catalog Manager |
| | Distributed Execution | Transaction Manager |
| | Query Processor | |
| | Data Formats, Indexes    Buffer Manager | |
| | Distributed Storage    Log Manager | |
| Process Manager | Storage Engine | |

RocksDB  26

levelDB  10

Pebble  1

SQLite  13

LMDB  12



**Figure 1: MySQL and MyRocks Storage Engine**

MyRocks: LSM-Tree Database Storage Engine Serving Facebook's Social Graph
https://www.vldb.org/pvldb/vol13/p3217-matsunobu.pdf

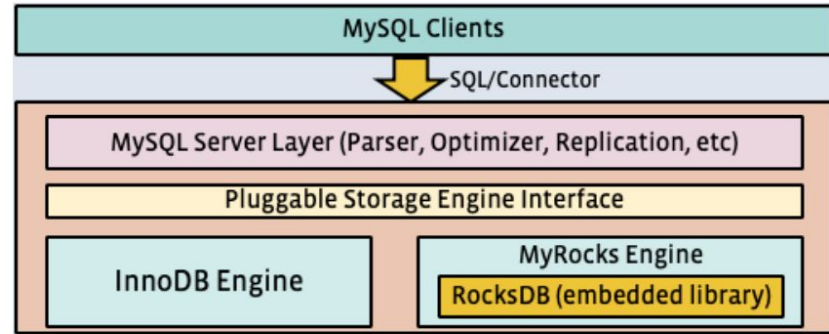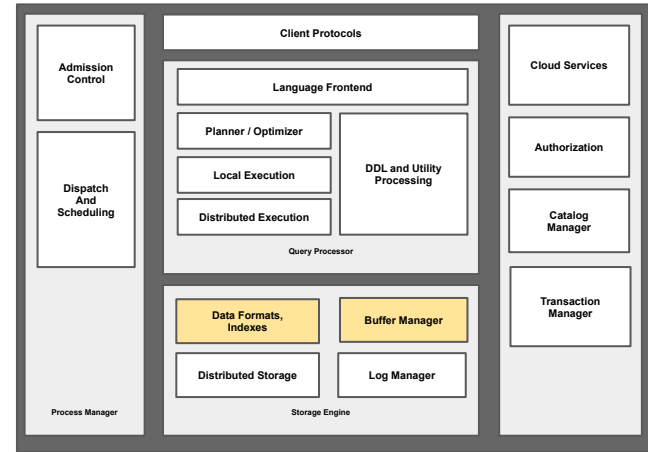| | | Parquet | ORC |
|---|---|---|---|
| **FEATURES** | Internal Layout (§3.1) | PAX | PAX |
| | Encoding Variants (§3.2) | plain, RLE_DICT, RLE, Delta, Bitpacking | plain, RLE_DICT, RLE, Delta, Bitpacking, FOR |
| | Compression (§3.3) | Snappy, gzip, LZO, zstd, LZ4, Brotli | Snappy, zlib, LZO, zstd, LZ4 |
| | Type System (§3.4) | Separate logical and physical type system | One unified type system |
| | Zone Map / Index (§3.5) | Min-max per smallest zone map/row group/file | Min-max per smallest zone map/row group/file |
| | Bloom Filter (§3.5) | Supported per column chunk | Supported per smallest zone map |
| | Nested Data Encoding (§3.6) | Dremel Model | Length and presence |
| **CONCEPTS** | Row Group | Row Group | Stripe |
| | Smallest Zone Map | Page Index (a Page) | Row Index (10k rows) |
| | Encoding Unit | Page | Stream |
| | Compression Unit | Page | Compression Chunk |

**Table 1:** Feature Taxonomy and Concepts Mapping



(a) Parquet layout.

(b) ORC layout.

**Figure 1:** Illustration of file layout — Parquet (a) and ORC (b). Blocks in  gray  are optional depending on configurations/data.

An Empirical Evaluation of Columnar Storage Formats
https://arxiv.org/pdf/2304.05028.pdf

**Client Protocols**

**Query Processor**

- **Language Frontend**
- **Planner / Optimizer**
- **Local Execution**
- **Distributed Execution**
- **DDL and Utility Processing**

**Storage Engine**

- **Data Formats, Indexes**
- **Buffer Manager**
- **Distributed Storage**
- **Log Manager**

**Process Manager**

- **Admission Control**
- **Dispatch And Scheduling**

- **Cloud Services**
- **Authorization**
- **Catalog Manager**
- **Transaction Manager**

**NATIVE DRIVERS**  **HTTP ACCESS**  **MONGODB DRIVERS**

**Tigris**

**gRPC and REST APIs**  **MONGODB COMPATIBILITY**

QUERY ENGINE

**AUTHENTICATION**  **ACCESS CONTROL**  **AUDIT LOGGING**  **MONITORING**

**DOCUMENT MANAGEMENT**  **QUERY PARSING**  **METADATA MANAGEMENT**

**QUERY PLANNING AND EXECUTION**

**TRANSACTION MANAGEMENT AND COORDINATION**

**IN-MEMORY SEARCH INDEXING**

**SHARD 1**  **SHARD 2**  .....  **SHARD N**

**FOUNDATIONDB**

STORAGE

**SHARD 1**  **SHARD 2**  .....  **SHARD N**

Durability  Transactions  Sharding  Replication

**FOUNDATIONDB**

**INFRASTRUCTURE PROVIDERS**  aws  Google Cloud  Azure  DATA CENTER

Tigris - Document Database built on FoundationDB

https://www.tigrisdata.com/docs/concepts/architecture/

**Client Protocols**

**Admission Control**

**Dispatch And Scheduling**

**Language Frontend**

**Planner / Optimizer**

**Local Execution**

**DDL and Utility Processing**

**Distributed Execution**

Query Processor

**Data Formats, Indexes**  **Buffer Manager**

**Distributed Storage**  **Log Manager**

Storage Engine

Process Manager

**Cloud Services**

**Authorization**

**Catalog Manager**

**Transaction Manager**

Fig. 1.  Different applications are able to run on disaggregated storage with RocksDB.

## Disaggregating RocksDB: A Production Experience

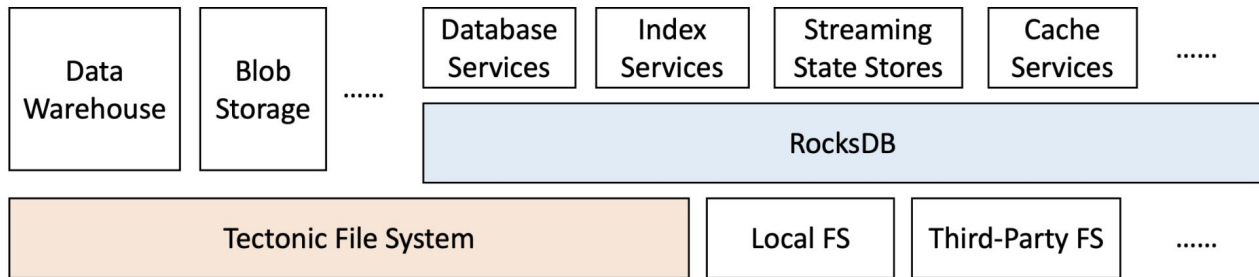https://research.facebook.com/publications/disaggregating-rocksdb-a-production-experience/

| Process Manager | Query Processor | | | Cloud Services |
|---|---|---|---|---|
| Admission Control | Client Protocols | | | Cloud Services |
| | Language Frontend | | | Authorization |
| Dispatch And Scheduling | Planner / Optimizer | DDL and Utility Processing | | Catalog Manager |
| | Local Execution | | | |
| | Distributed Execution | | | Transaction Manager |
| | Storage Engine | | | |
| | Data Formats, Indexes | Buffer Manager | | |
| | Distributed Storage | Log Manager | | |

Neon Architecture
https://neon.tech/docs/introduction/architecture-overview

(1) **Disaggregated Compute-Storage OLAP Architecture**

# Computing

## Firebolt Engine

CPU optimized
Query engine

**Ad hoc**

## Local SSD Cache

▤|| Sparse Index A

▤|| Sparse Index B

---

▤|| Indexes

🗄 Data

---

# Storage (S3)

F3  F3  F3

F3  F3

## File A

Range A.1
Range A.2
Range A.3
Range A.4
Range A.5

▤|| Sparse Index A

## File B

Range B.1
Range B.2
Range B.3
Range B.4
Range B.5

▤|| Sparse Index B

---

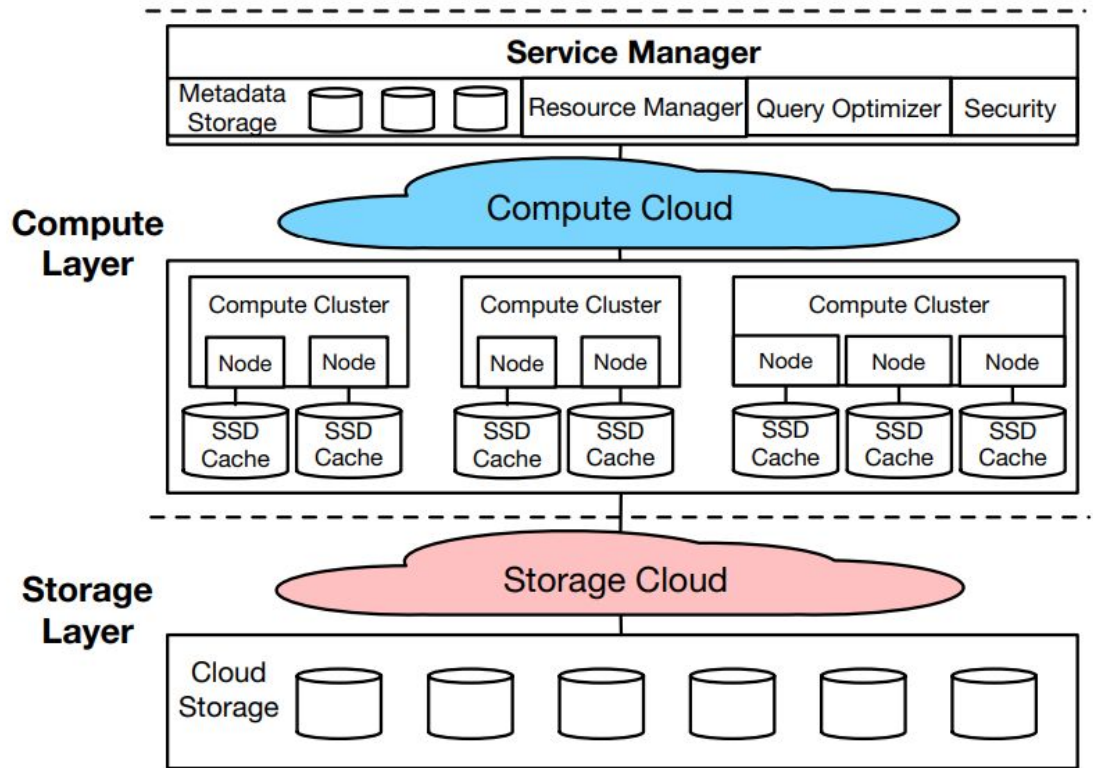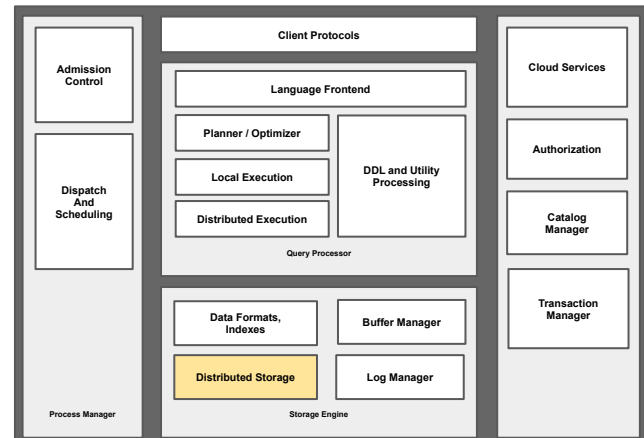| Process Manager | Query Processor | | | Cloud Services |
|---|---|---|---|---|
| **Admission Control** | **Client Protocols** | | | **Cloud Services** |
| | **Language Frontend** | | | **Authorization** |
| **Dispatch And Scheduling** | **Planner / Optimizer** | **DDL and Utility Processing** | | **Catalog Manager** |
| | **Local Execution** | | | |
| | **Distributed Execution** | | | **Transaction Manager** |

### Storage Engine

| Data Formats, Indexes | Buffer Manager |
|---|---|
| **Distributed Storage** | **Log Manager** |

**Computing**

SQL

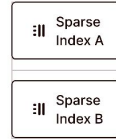Optimize → Logical query plan

Physical query plan

What data do I need?

**Firebolt Engine**

CPU optimized
Query engine
Ad hoc

**Local SSD Cache**

Sparse Index A

Sparse Index B

Indexes          Data

**Storage (S3)**

F3  F3  F3

F3  F3

Sparse Index A

Sparse Index B

File A
Range A.1
Range A.2
Range A.3
Range A.4
Range A.5

File B
Range B.1
Range B.2
Range B.3
Range B.4
Range B.5

Admission Control

Dispatch And Scheduling

Process Manager

Client Protocols

Language Frontend

Planner / Optimizer          DDL and Utility Processing

Local Execution

Distributed Execution

Query Processor

Data Formats, Indexes          Buffer Manager

Distributed Storage          Log Manager

Storage Engine

Cloud Services

Authorization

Catalog Manager

Transaction Manager

# Computing

**Firebolt Engine**

CPU optimized
Query engine

**Ad hoc**

SQL

Optimize → Logical query plan

Physical query plan

What data do I need?

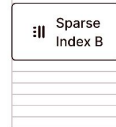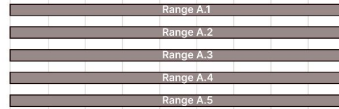## Local SSD Cache

Sparse Index A

Sparse Index B

≡ll Indexes

🗄 Data

# Storage (S3)

F3  F3  F3

F3  F3

Sparse Index A

Sparse Index B

**File A**

Range A.1
Range A.2
Range A.3
Range A.4
Range A.5

**File B**

Range B.1
Range B.2
Range B.3
Range B.4
Range B.5

---

Client Protocols

Admission Control

Dispatch And Scheduling

Language Frontend

Planner / Optimizer

Local Execution

Distributed Execution

DDL and Utility Processing

Query Processor

Data Formats, Indexes

Buffer Manager

Distributed Storage

Log Manager

Storage Engine

Process Manager

Cloud Services

Authorization

Catalog Manager

Transaction Manager

**Computing**

Firebolt
Engine

CPU optimized
Query engine

Ad hoc

SQL

Optimize → Logical query plan

Physical query plan

What data do I need?

**Local SSD Cache**

Sparse Index A

Sparse Index B

Range A.2
Range B.4

| Indexes

Data

**Storage (S3)**

F3  F3  F3

F3  F3

Sparse Index A

Sparse Index B

**File A**

Range A.1
Range A.2
Range A.3
Range A.4
Range A.5

**File B**

Range B.1
Range B.2
Range B.3
Range B.4
Range B.5

Admission Control

Dispatch And Scheduling

Process Manager

Client Protocols

Language Frontend

Planner / Optimizer

Local Execution

Distributed Execution

DDL and Utility Processing

Query Processor

Data Formats, Indexes

Distributed Storage

Buffer Manager

Log Manager

Storage Engine

Cloud Services

Authorization

Catalog Manager

Transaction Manager

## Left diagram

**Computing**

**Firebolt Engine**

CPU optimized
Query engine

**Ad hoc**

SQL
→ Optimize → Logical query plan
→ Physical query plan
→ What data do I need?
→ Execute the query

**Local SSD Cache**

Sparse Index A

Sparse Index B

Scan

Range A.2
Range B.4

**Indexes** | **Data**

**Storage (S3)**

F3  F3  F3
F3  F3

Sparse Index A

**File A**
Range A.1
Range A.2
Range A.3
Range A.4
Range A.5

Sparse Index B

**File B**
Range B.1
Range B.2
Range B.3
Range B.4
Range B.5

## Right diagram

**Admission Control**

**Dispatch And Scheduling**

Process Manager

**Client Protocols**

**Language Frontend**

**Planner / Optimizer**

**Local Execution**

**Distributed Execution**

**DDL and Utility Processing**

Query Processor

**Data Formats, Indexes**

**Buffer Manager**

**Distributed Storage**

**Log Manager**

Storage Engine

**Cloud Services**

**Authorization**

**Catalog Manager**

**Transaction Manager**

| | SELECT | DDL/DML/DCL/TCL | Total | |
|---|---|---|---|---|
| SQLite | 7.2M | 225K | 7.4M | 12K |

|  | SELECT | DDL/DML/DCL/TCL | Total |  |
|---|---|---|---|---|
|  | 7.2M | 225K | 7.4M | 12K |
|  | 24K | 22K | 46K | 4K |

|  | SELECT | DDL/DML/DCL/TCL | Total |  |
|---|---|---|---|---|
|  | 7.2M | 225K | 7.4M | 12K |
|  | 24K | 22K | 46K | 4K |
|  | 17K | 28K | 35K | 3K |

|  | SELECT | DDL/DML/DCL/TCL | Total | |
|---|---|---|---|---|
| SQLite | 7.2M | 225K | 7.4M | 12K |
| | 24K | 22K | 46K | 4K |
| | 17K | 28K | 35K | 3K |
| | 45K | 4K | 49K | 3K |

| | SELECT | DDL/DML/DCL/TCL | Total |  |
|---|---|---|---|---|
|  SQLite | 7.2M | 225K | 7.4M | 12K |
|  | 24K | 22K | 46K | 4K |
|  | 17K | 28K | 35K | 3K |
|  | 45K | 4K | 49K | 3K |
|  | 30K | 7K | 37K | 22K |