

---



# StarRocks Technical Overview

A Linux Foundation Project

Albert Wong, albert.wong@celerdata.com



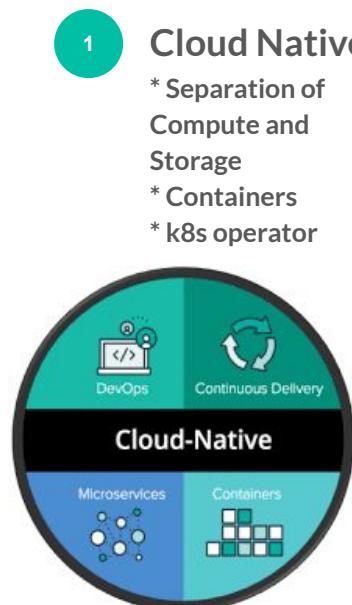
# Agenda

- State of the OLAP software landscape
- What is StarRocks
- StarRocks' Release Timeline
- Major Features in StarRocks

# State of the OLAP software landscape

# Trends in OLAP databases

Online analytical processing (OLAP) databases are evolving rapidly to meet the demands of modern data analytics. Here are some of the key trends in OLAP databases:



1

## Cloud Native

- \* Separation of Compute and Storage
- \* Containers
- \* k8s operator

2 A

## Sub-second vs. Second/Minute Query Response Time

2 B

## Streaming vs. Batch Data

2 C

## Mutable vs. Immutable Data

2 D

## Remote (Object) Storage vs. Local (SSD) Storage

2 E

## Open Table Format vs. Product Native Storage Format

3

## Data Warehouse vs. Data Lake vs. Data Lakehouse



Data Warehouse



Data Lake



Data Lakehouse

# Trends in OLAP databases

Open Lakehouse vs Proprietary / Hybrid Lakehouse.

	Proprietary / Hybrid	Open Storage	Open
Compute	 snowflake	 databricks photon	 StarRocks  trino
Data Catalog	 snowflake	 databricks  HIVE	 HIVE
Table Format	 snowflake  ICEBERG	 DELTA LAKE	 HIVE  ICEBERG  Apache Hudi
Storage Format	 snowflake  Parquet	 Parquet	 Parquet

StarRocks is an open-source  
query engine that delivers data  
warehouse performance on the  
data lake.

# StarRocks Community



7500+ Github Stars

350+ Contributors

18,000+ Community Members

As of Feb 2024



SHEIN

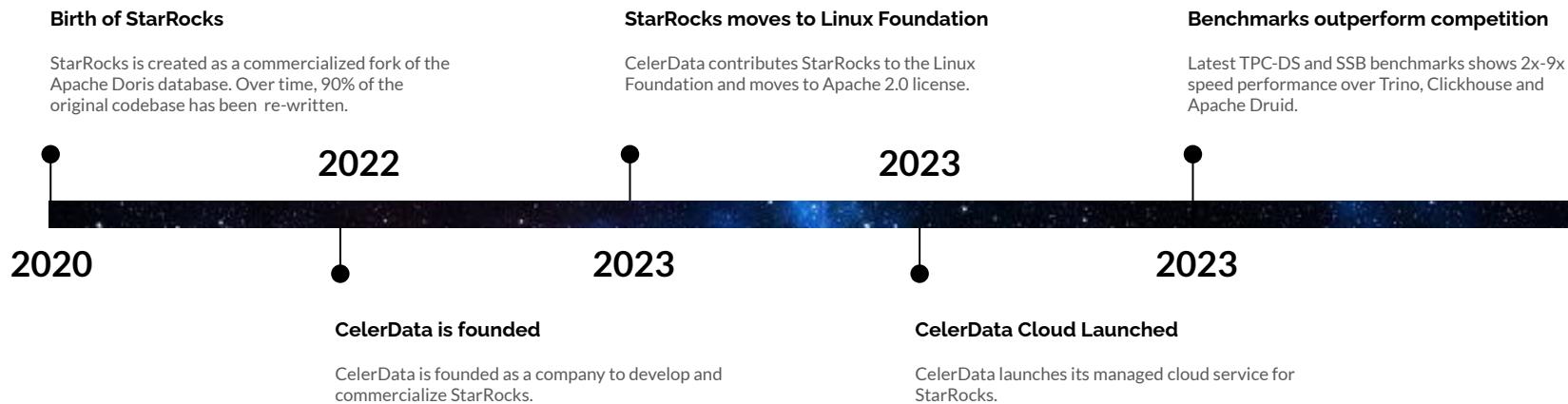


Tencent



# History of StarRocks and CelerData

StarRocks was designed to address the challenges of real-time analytics, including the need to support high concurrency, low latency, and a wide range of analytical workloads. StarRocks also offers a number of features that are not available in other real-time analytics databases, such as the ability to query data directly from data lakes.

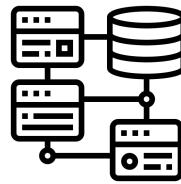


---

# StarRocks is an open-source query engine that delivers data warehouse performance on the data lake.



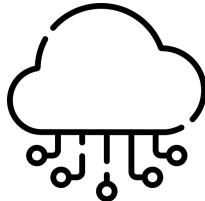
Directly query data on data lake



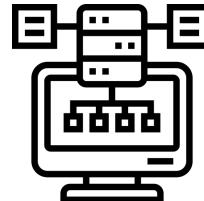
Sub-second joins and aggregations on billions of rows



Hundreds of thousands of concurrent end-user requests



Cloud Native w/ separating Compute and Storage tiers



JOIN performant at scale; denormalization optional



mysql protocol with Trino dialect

# StarRocks Use Cases: User-facing analytics

User-facing analytics (UFA) is a rapidly growing field that is transforming the way businesses deliver insights to their users. UFA empowers users to explore and analyze data for themselves, without the need for technical expertise. This can lead to a number of benefits, such as:

1

Improved decision making

2

Increased user engagement

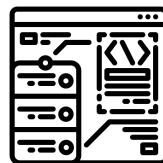
3

Reduced reliance on IT

Key trends in user-facing analytics:



Self-service Analytics



Embedded Analytics



Real-Time Analytics



Augmented Analytics



Conversational Analytics

# StarRocks Use Cases: Real-time analytics

Real-time analytics is the process of collecting, processing, and analyzing data as it is generated, in order to gain insights into the present state of a system or process. This can lead to a number of benefits, such as:

1

Make better decisions

2

Increased user engagement

3

Staying ahead of the competition

Key trends in real-time analytics:



Rise of streaming data



Growth of Edge Computing



Increasing use of machine learning



Democratization of real-time analytics

# StarRocks Use Cases: Data Lakehouse

A data lakehouse is a revolutionary data architecture that merges the best of both data lakes and data warehouses.

1    Democratized Data Access

2    Increased agility and insights

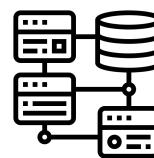
3    Reduced costs and complexity

3    Faster and more accurate analytics

Key trends in data lakehouse:



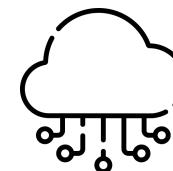
Directly query  
data on data lake



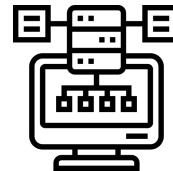
Sub-second joins and  
aggregations on  
billions of rows



Hundreds of thousands  
of concurrent end-user  
requests



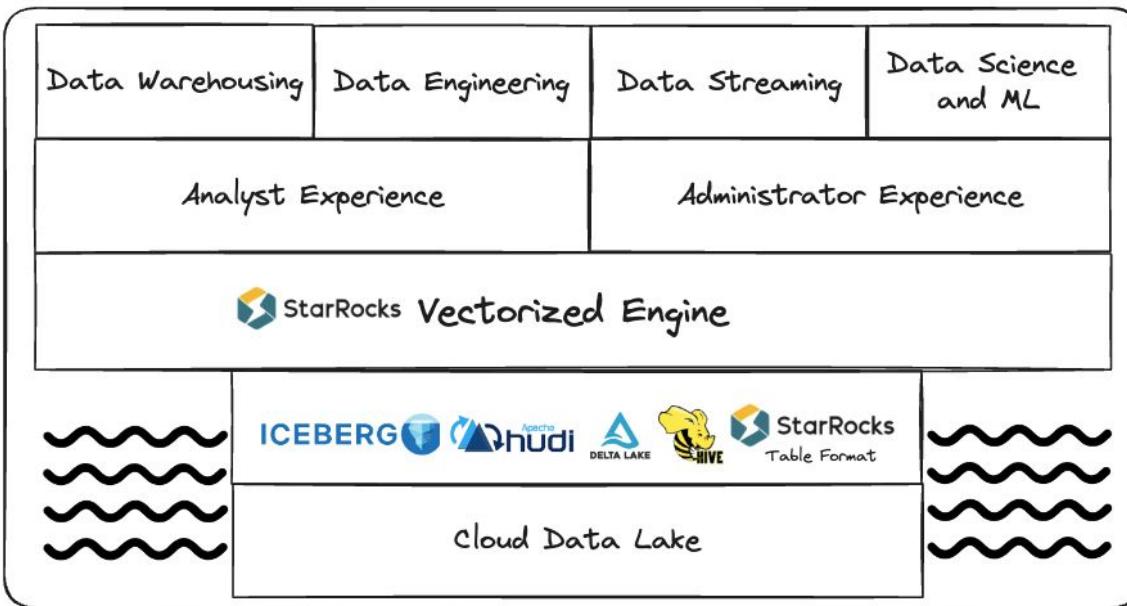
Cloud Native w/  
separating Compute  
and Storage tiers



JOIN performant at scale;  
denormalization optional

# StarRocks Architecture Overview

Linux Foundation project with Apache 2.0 license.



Seamless integration with the Ecosystem

Ease of Use

Real-world Performance

Open Source OLAP compute engine

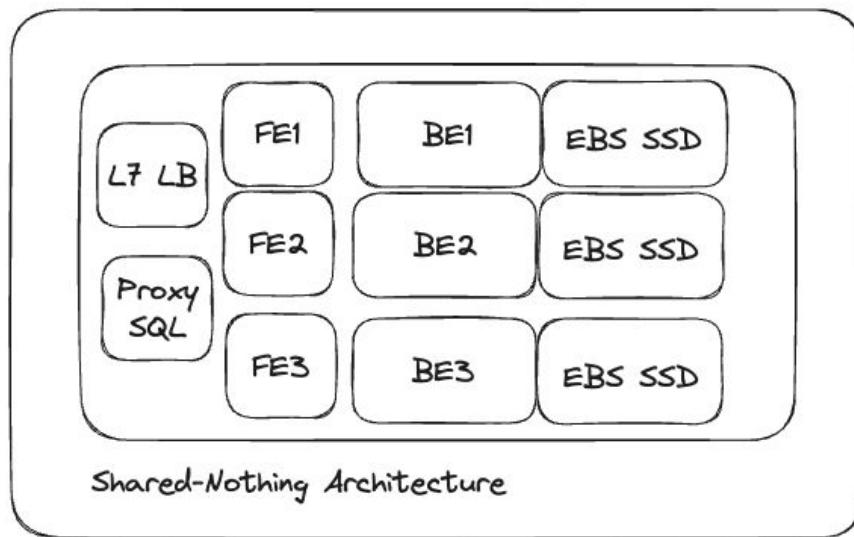
Open Table Formats as the Foundation

Support for Open Storage

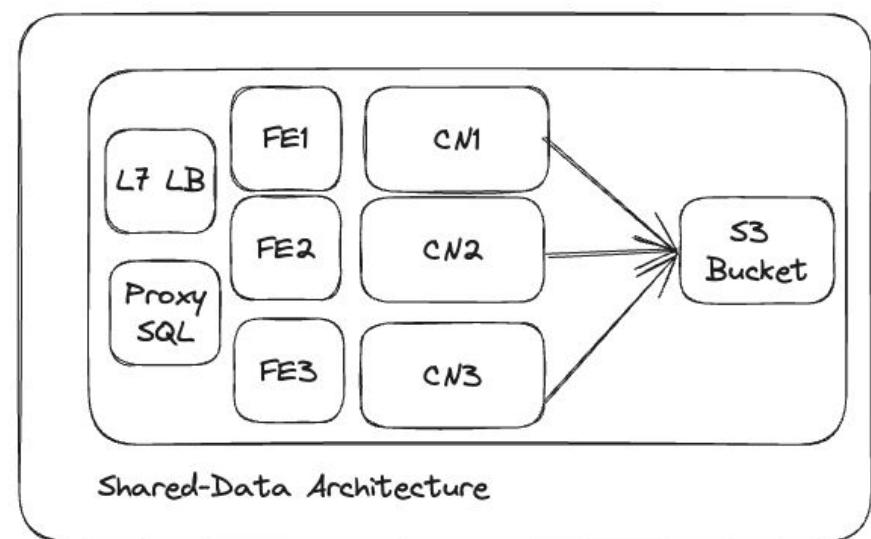
Separated compute and storage architecture

Cloud Native with k8s Operator

# Two Deployment Architectural Choices



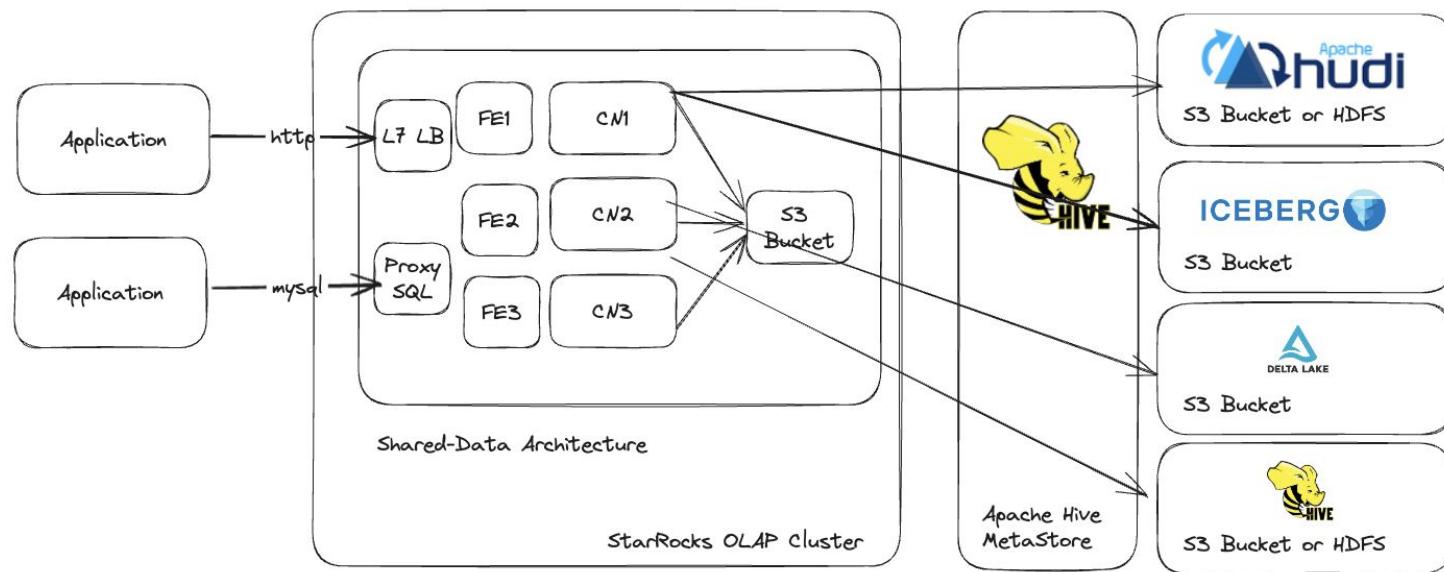
StarRocks OLAP Cluster



StarRocks OLAP Cluster

# StarRocks with Open Data Lake

StarRocks can access multiple open table formats at the same time and even be able to create a materialized view across all of them.



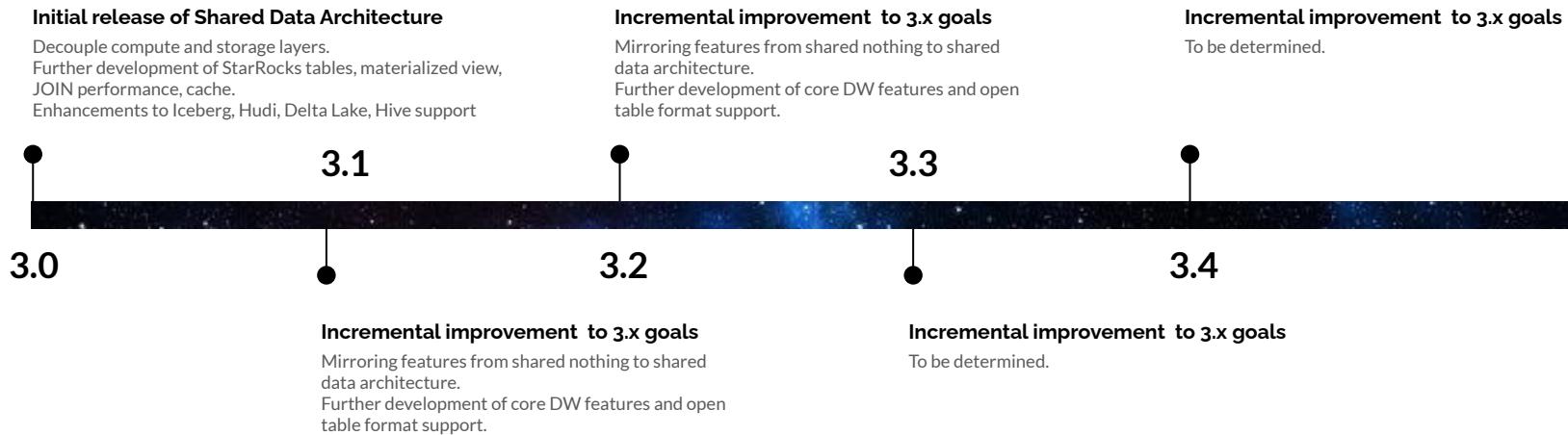
# StarRocks past and future

# StarRocks Technical Features

1.x (2020-2021)	2.x (2022-2023)	3.0 (2023)	3.1 (2023)	3.2 (2023)
<b>OLAP for Real Time Analytics</b> <ul style="list-style-type: none"> <li>• Vectorized Execution Engine</li> <li>• Cost Based Optimizer</li> <li>• Vectorized ingestion</li> <li>• Apache Hive support</li> <li>• Bitmap optimization</li> <li>• TopN optimization</li> <li>• Lateral JOIN</li> <li>• Fast Decimal support</li> <li>• Tableau compatibility</li> <li>• Global runtime filter</li> <li>• Primary Key Table</li> </ul>	<b>OLAP for Data Lake</b> <ul style="list-style-type: none"> <li>• Global low-cardinality dictionary</li> <li>• Pipeline Engine</li> <li>• Apache Iceberg Support</li> <li>• Resource Group</li> <li>• Java UDF</li> <li>• JSON data type support</li> <li>• Partial update feature</li> <li>• JDBC external catalog support</li> <li>• Primary key Index</li> <li>• Fully support delete and update operations</li> <li>• Multi-table materialized view</li> <li>• More table statistics including histogram</li> <li>• Compute node on k8s</li> <li>• Separation of storage and compute</li> <li>• Local cache for open table formats on data lake</li> <li>• Semantic cache</li> <li>• Fully support RBAC</li> <li>• Map/Struct data type</li> <li>• Lambda function</li> </ul>	<b>Shared Data Arch.</b> <ul style="list-style-type: none"> <li>• Shared Data Architecture</li> <li>• New RBAC privilege system</li> <li>• Spill to disk</li> <li>• Fully support for update</li> <li>• Support more complete UPDATE and DELETE syntax in primary key tables</li> <li>• Presto/Trino compatible dialect</li> <li>• Broadcast JOIN and Bucket Shuffle JOIN can use query cache</li> <li>• Global UDFs</li> </ul>	<b>Shared Data Arch. Optimization</b> <ul style="list-style-type: none"> <li>• Primary Key table support in Shared Data</li> <li>• Auto_Increment column attribute</li> <li>• Automatic partition creation during load</li> <li>• Support Apache Iceberg v2 tables</li> <li>• Random bucketing</li> <li>• FILES keyboard</li> <li>• Generated columns</li> <li>• Support loading data into MAP and STRUCT data types</li> <li>• Support nesting Fast Decimal values in ARRAY, MAP and STRUCT</li> <li>• Optimized creation of async materialized view</li> <li>• Optimized query rewrite with async materialized views</li> <li>• Optimized refreshing of aysn materialized views</li> <li>• Optimized caching, and query logic for StarRocks table format and Apache Iceberg</li> </ul>	<b>Shared Data Arch. Optimization v2</b> <ul style="list-style-type: none"> <li>• Persisting Primary Key table indexes to local disk</li> <li>• Spill to Disk enabled by default for async materialized views</li> <li>• Support creating, dropping database and managed tables in Apache Hive catalogs</li> <li>• Unified Catalog</li> <li>• Supports Information Schema for external tables</li> <li>• Enhanced Files()</li> <li>• Support unloading data from StarRocks to parquet</li> <li>• Supports manual optimization of table structure and data distribution strategy</li> <li>• Continuous data loading using PIPE</li> <li>• Support HTTP SQL API</li> <li>• Runtime Profile and text-based profile analysis commands</li> <li>• Support access control through Apache Ranger</li> <li>• Optimized open file format readers</li> <li>• Added data consistency features for async materialized view</li> <li>• Hot and warm storage support</li> <li>• Fast Schema Evolution</li> <li>• Dynamically adjusting number of tables</li> <li>• Data redistribution across local disks for primary key tables</li> </ul>

# StarRocks 3.x series roadmap

The goal of the 3.x series roadmap is to 1) Build more and optimize core data warehouse features, 2) have feature parity between the shared-nothing architecture and shared-data architecture and 3) be able to query the StarRocks table format and all the popular open table formats such as Apache Iceberg, Apache Hudi, Apache Hive, Delta Lake and Apache Paimon.

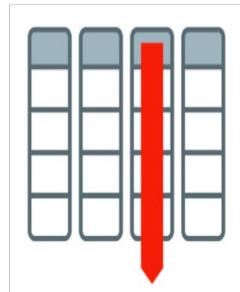


# Major Features in StarRocks

# Vectorized Query Engine with SIMD

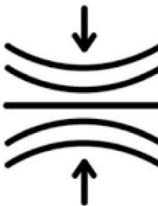
Modern CPUs have vectorized instruction sets, which can perform operations on multiple data elements simultaneously which means faster queries by 3x to 5x over non-SIMD databases.

## Scan



- Reducing IO
- SIMD filter

## Compression



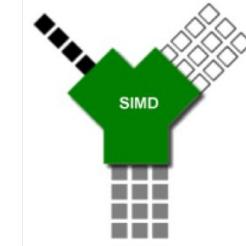
- 4-8x compression rate
- Adaptive codec selection

## Indexes



- Min/max
- Ordinal index
- Secondary index

## Vectorized



- SIMD Instructions
- CPU prefetch and cache friendly
- Branch prediction friendly

IO acceleration

CPU acceleration

# Table Type Support

Tables are units of data storage. Understanding the table structure in StarRocks and how to design an efficient table structure helps optimize data organization and enhance query efficiency.

## Types of Tables supported

- Duplicate Key
  - Analyze raw data, such as raw logs and raw operation records.
  - Query data by using a variety of methods without being limited by the pre-aggregation method.
  - Load log data or time-series data. New data is written in append-only mode, and existing data is not updated.
- Aggregate Key
  - Help website or app providers analyze the amount of traffic and time that their users spend on a specific website or app and the total number of visits to the website or app.
  - Help advertising agencies analyze the total clicks, total views, and consumption statistics of an advertisement that they provide for their customers.
  - Help e-commerce companies analyze their annual trading data to identify the geographic bestsellers within individual quarters or months.
- Primary Key
  - Stream data in real time from transaction processing systems into StarRocks.
  - Join multiple streams by performing update operations on individual columns.

Table Type	Duplicate Key	Append Only	<input checked="" type="checkbox"/>
	Aggregate Key	Append Only	<input checked="" type="checkbox"/>
	Primary Key	All CRUD	<input checked="" type="checkbox"/>

# JOIN performance at scale

Simply your data engineering pipeline and infrastructure by using JOINS; denormalization is optional.

## Types of JOINS supported

- CBO will do intelligent Join reorder and Join method selection
- Starrocks can join 100 million rows of data per second using only 1 CPU.

Details at

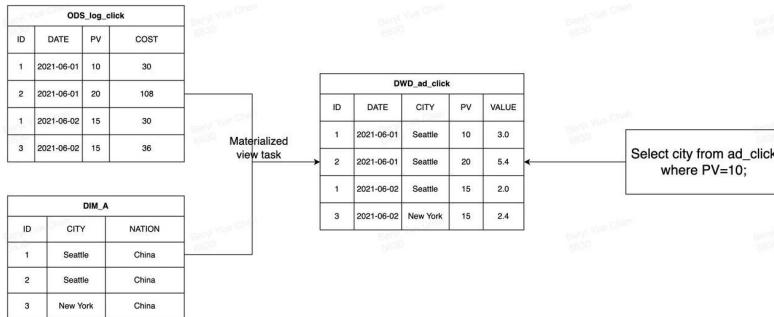
[https://www.starrocks.io/blog/bench  
mark-test](https://www.starrocks.io/blog/bench-mark-test)

SQL JOINS	
Inner Join	✓
Left Join	✓
Right Join	✓
Full Join	✓
Cross Join	✓
Semi Join	✓
Anti Join	✓
SQL JOINS Optimization Technique	
Broadcast Join	✓
Shuffle Join	✓
Bucket Shuffle Join	✓
Co-Located Join	✓
Replicated Join	✓
Local Join	✓

# Materialized View

Materialized views can significantly improve query performance by pre-computing common aggregations.

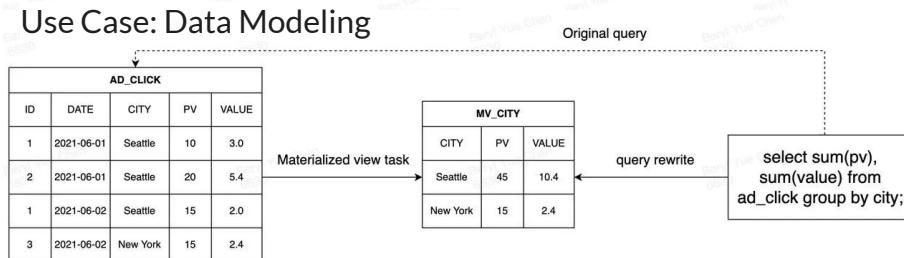
## Use Case: Query Acceleration



Select city from ad\_click  
where PV=10;

PROJECT	<input checked="" type="checkbox"/>
AGGREGATE	<input checked="" type="checkbox"/>
JOIN	<input checked="" type="checkbox"/>
Outer-Join	<input checked="" type="checkbox"/>
View-Delta-Join	<input checked="" type="checkbox"/>
PARTIAL-UNION	<input checked="" type="checkbox"/>
NESTED MV	<input checked="" type="checkbox"/>
View-Based	<input checked="" type="checkbox"/>
Auto Refresh	<input checked="" type="checkbox"/>
Scheduled Refresh	<input checked="" type="checkbox"/>
Partition-Wise	<input checked="" type="checkbox"/>

Transparent Speedup  
(Core Functionality)

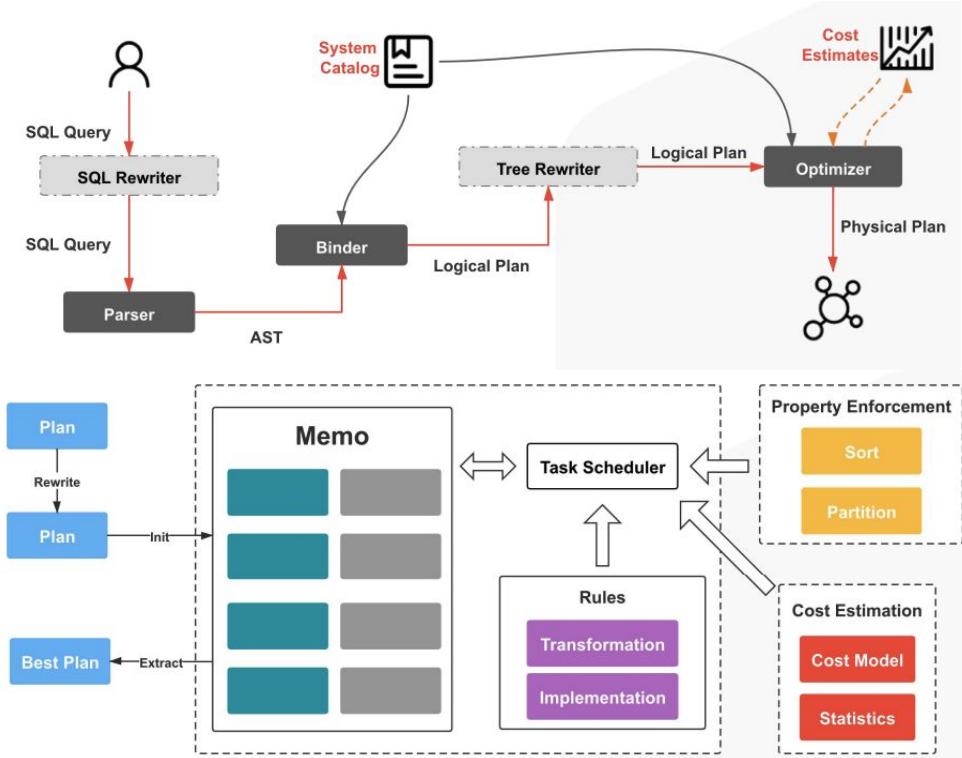
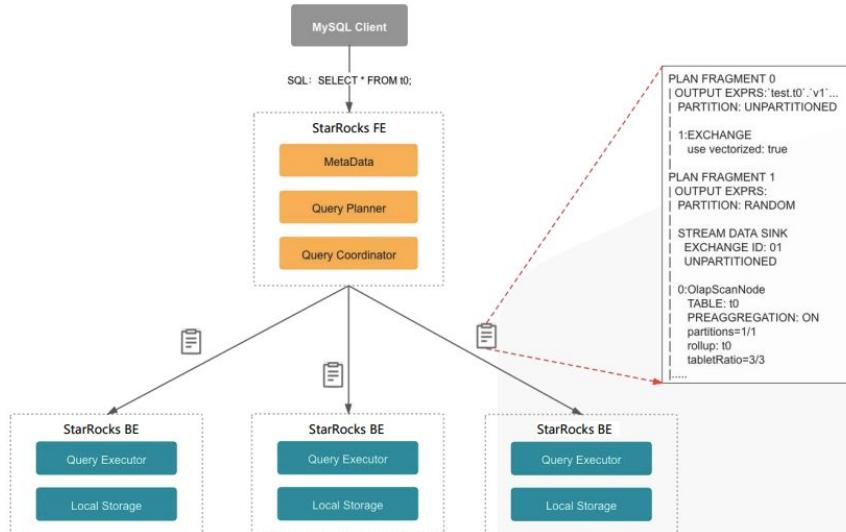


query rewrite  
select sum(pv),  
sum(value) from  
ad\_click group by city;

Incremental Refresh  
(Core Functionality)

# SQL Hybrid-Based Optimizer

Analyzes a SQL query and chooses the most efficient execution plan by estimating the cost of different potential plans

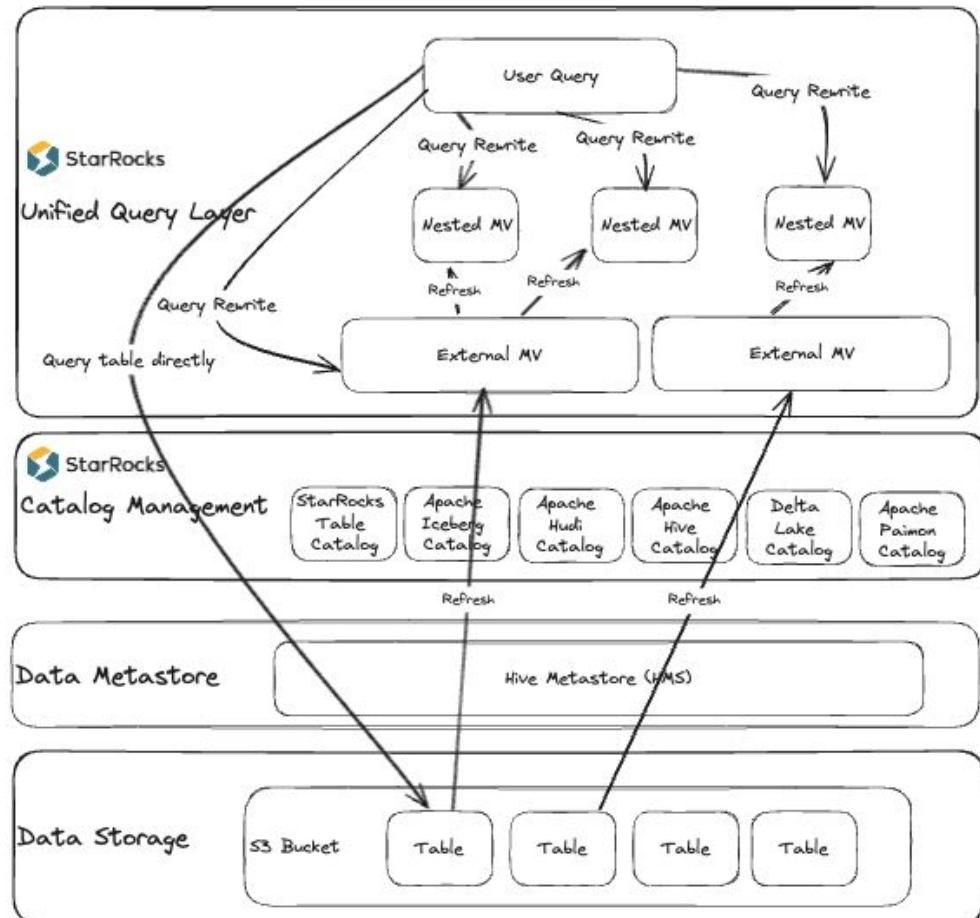


# Query Rewrite

Technique used to optimize database queries without the user needing to change their original query.

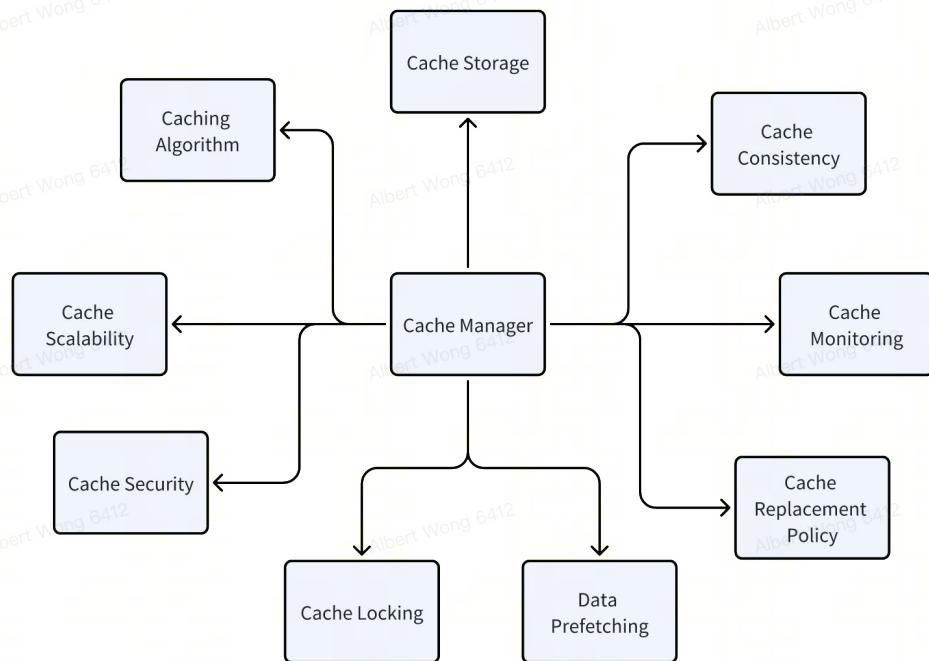
## Use Case: Semantic Layer

- Targeted at Select - Projection - Join - Aggregation (SPJA) query pattern
- Up to 10x performance increase



# Cache System

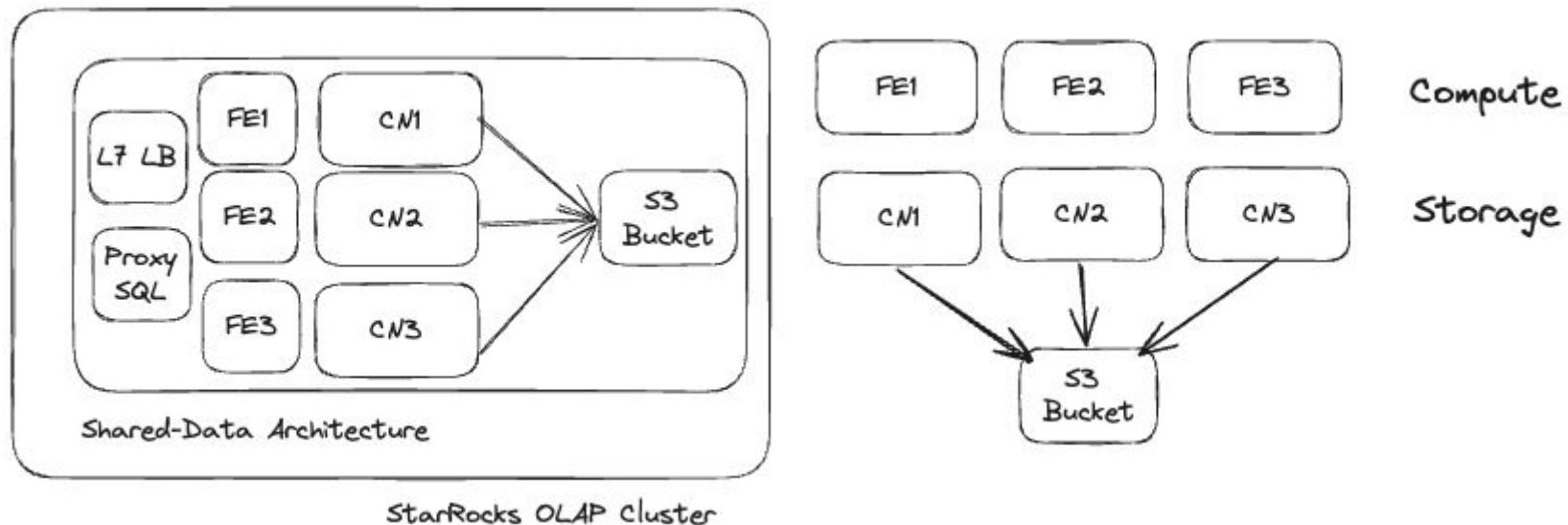
Cache allows you to pull the data from memory instead of storage which can improve query efficiency by 3x to 17x.



	Metadata	✓
Transparent Speedup (Cache Functionality)	Query	✓
	Page	✓
	Data	✓

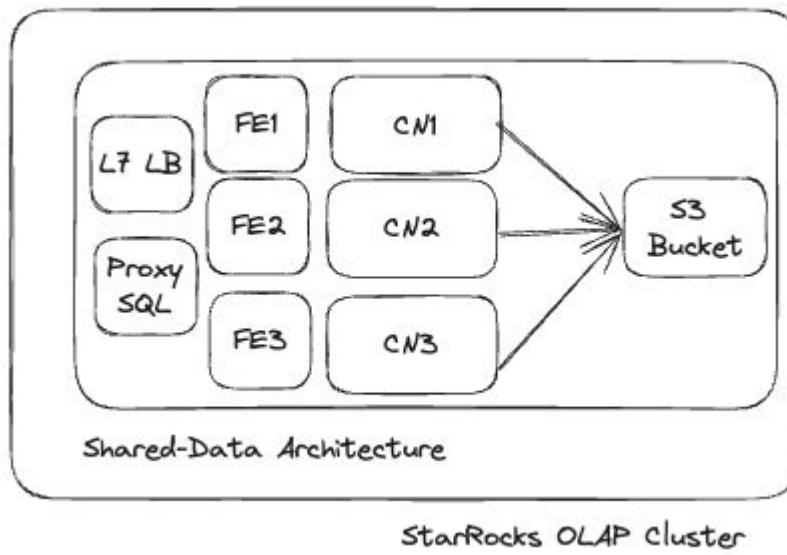
# Separated compute and storage architecture

Design approach for databases and data platforms that decouples the processing power (compute) from the data storage layer.



# High Availability

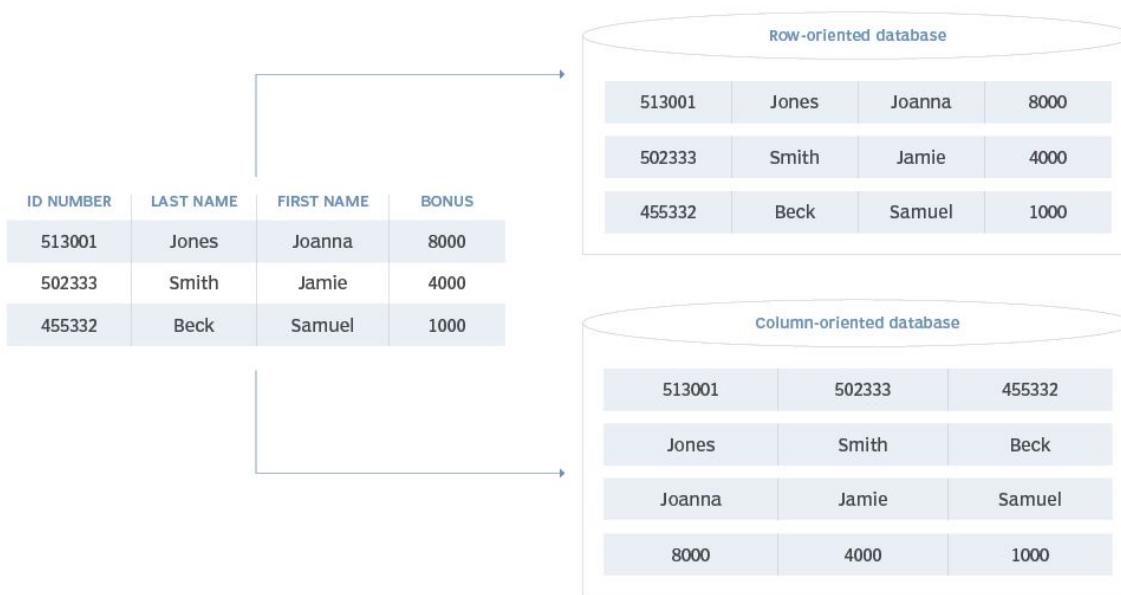
Redundant components and data allows the database to respond even when there is failure.



Service Availability	FE	Additional Nodes
	FE	✓ Additional Nodes
	CN	✓ Additional Nodes
	MySQL	✓ 3rd party ProxySQL
	HTTP Services	✓ 3rd party Load Balancer
	S3 Bucket	✓ 3rd party vendor
	Data Files	✓ 3rd party vendor S3 Bucket Vendor

# Columnar Storage

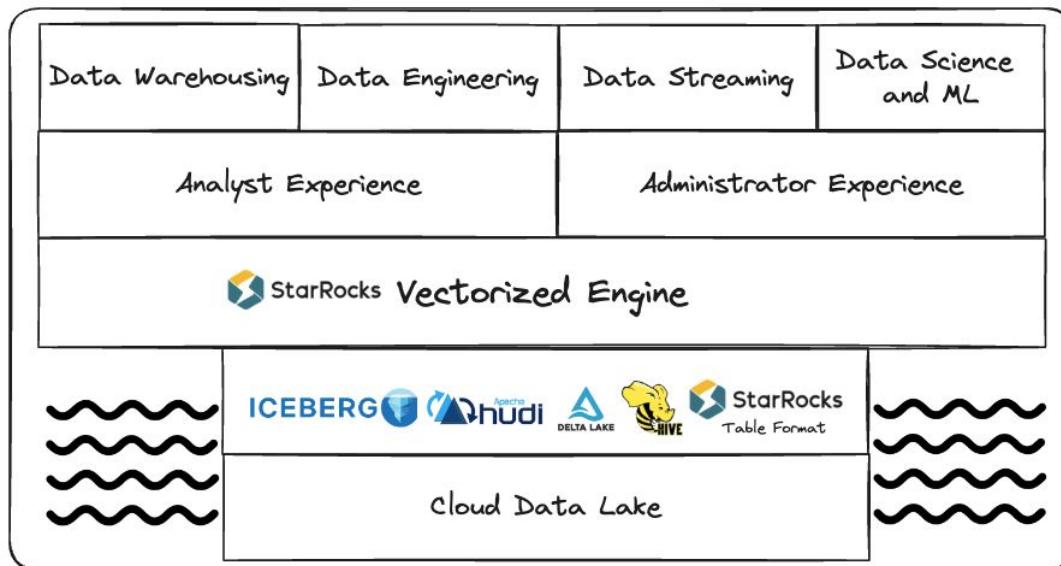
Stores data in a table by separating each column into its own continuous block instead of grouping entire rows together.



Columnar Storage Formats	StarRocks Table Format
	Apache Iceberg
	Apache Hudi
	Apache Hive
	Delta Lake
	Apache Paimon

# Support for Open Table Formats

Open Table Formats allow users to extract more value from their data while maintaining flexibility and control.



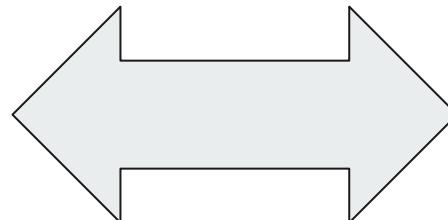
Open Table Formats	StarRocks Table Format	(Read/Write)
	Apache Iceberg	(Read/Write)
	Apache Hudi	(Read)
	Apache Hive	(Read/Write)
	Delta Lake	(Read)
	Apache Paimon	(Read)

# SQL Connectivity through MySQL wire protocol support with Trino dialect

Communicate with StarRocks through MySQL statements and utilities. Also understands the Trino SQL dialect.



Client



StarRocks

Server



Apache Flink



kafka



Airbyte

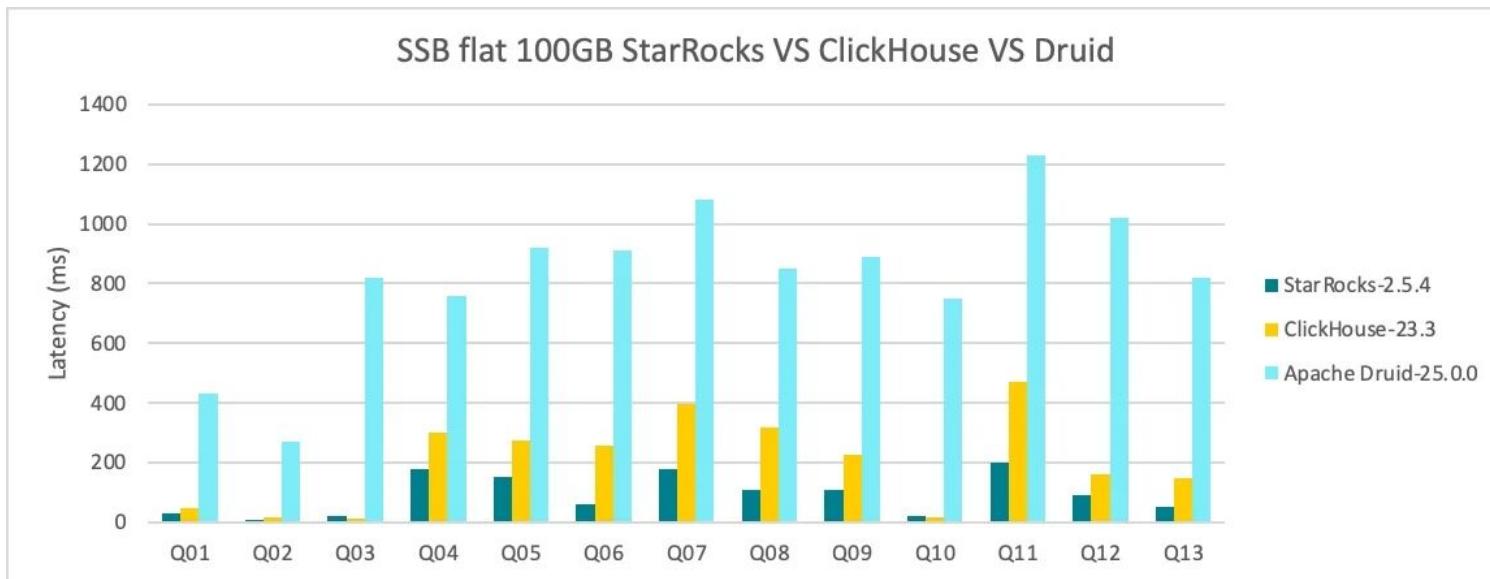


dbt™

# Benchmarks and Community References

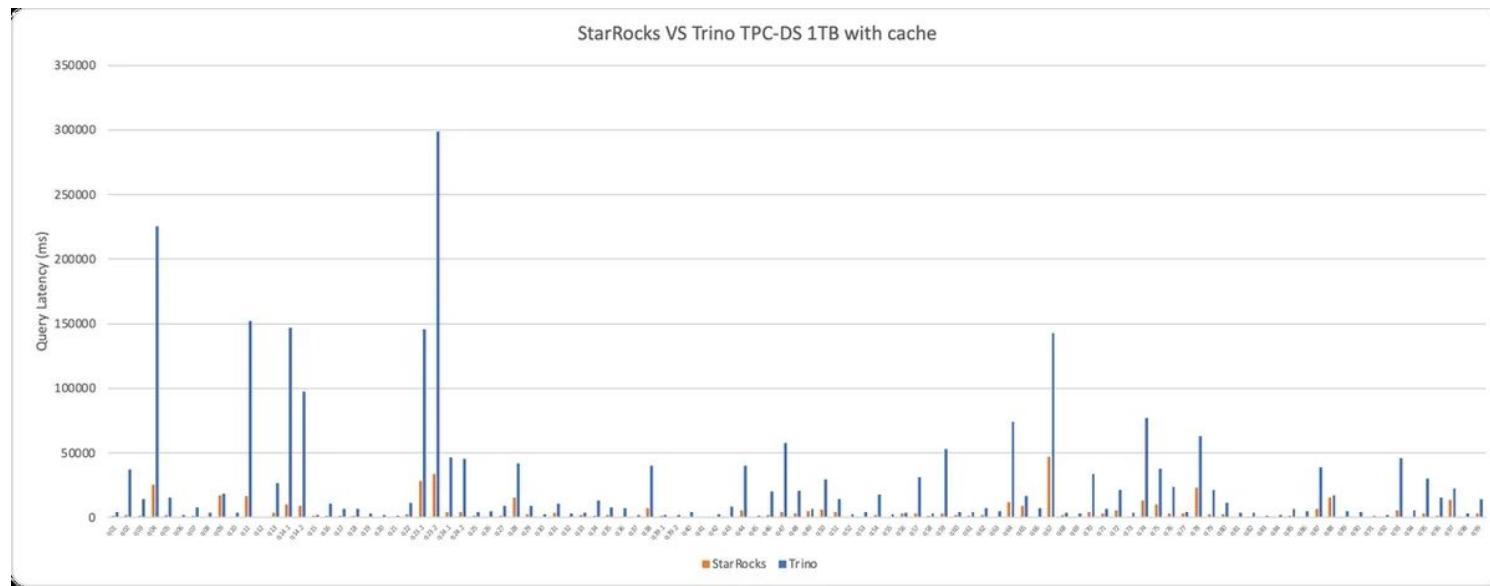
# Benchmark

StarRocks Offers 2.2x Performance over ClickHouse and 8.9x Performance over Apache Druid® in Wide-table Scenarios Out of the Box using product native table format.



# Benchmark

StarRocks Delivers 5.54x Query Performance over Trino in Multi-table Scenarios using Apache Iceberg table format with Parquet files.



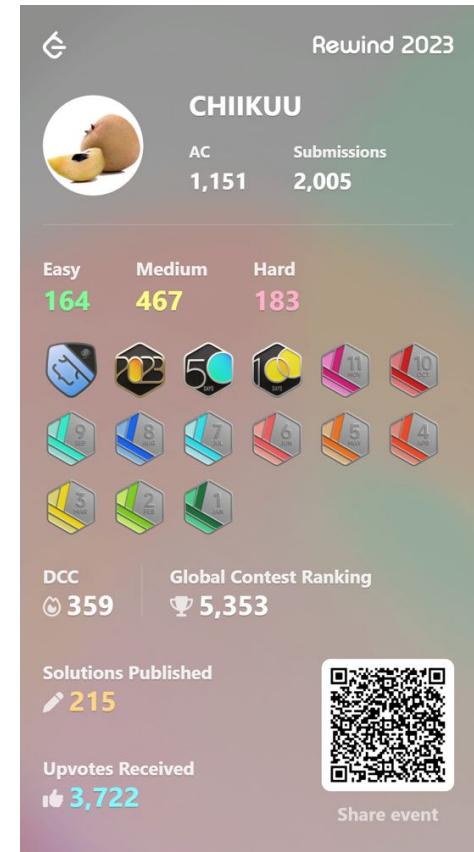
# Use Case: User Analytics at LeetCode



LeetCode's current data warehouse, built on an OLTP database, was struggling under the weight of terabytes of user activity data. Using this OLTP database, queries took ages, impacting user experience and hindering LeetCode's ability to analyze trends and optimize the platform. Scaling up the existing system proved costly and unsustainable.

## StarRocks Solution:

- Queries 100x Faster: Complex analytics that previously took hours now finished in seconds, empowering LeetCode to gain real-time insights into user behavior and platform performance. Additionally, some queries that couldn't run in the OLTP system were able to run successfully in StarRocks.
- Unlimited Scalability: StarRocks' horizontal scaling effortlessly accommodated LeetCode's growing data volume, eliminating concerns about future bottlenecks.
- Cost Savings of 80%: Compared to the a similar managed OLAP solution on GCP, StarRocks delivered significant cost savings, allowing LeetCode to reinvest in platform development and user experience.

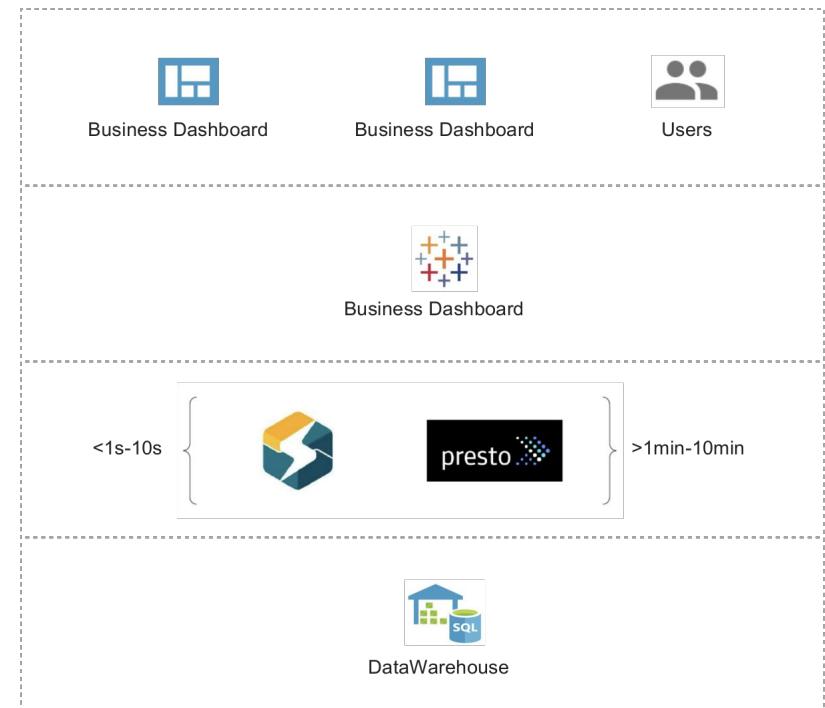


# Use Case: Tableau Dashboard at Airbnb

The Airbnb Tableau Dashboard project is designed to serve both internal and external users by providing interactive dashboards. It requires a quick response to user queries. However, the query latency of previous solutions is over 10 mins, which is not acceptable. This project was just suspended until StarRocks is adopted.

## StarRocks Solution:

- StarRocks can directly connect and works very well with Tableau.
- 3 tables (0.5B rows, 6B rows, 100M rows) + 4 joins + 3 distinct count + JSON functions and regex at same time, response time just 3.6s.
- Reduce the query response time from mins level to sub-seconds level.

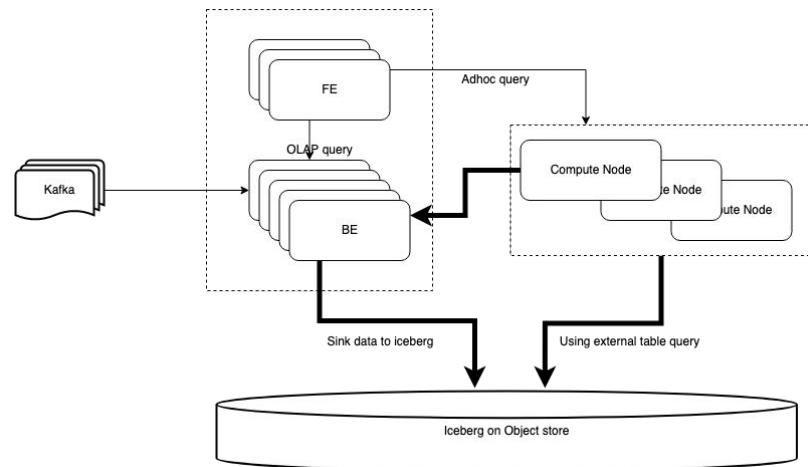


# Use Case: Game and User Behavior Analytics at Tencent IEG

- 400+ game data analysis and user behavior analysis
- Operation reports need to be real-time.
- Using ClickHouse for real-time analysis and Trino for Ad-hoc before, but they want to integrate them all.
- Using Iceberg + COS store, need better performance.
- Need elastic in ad-hoc query to deduce cost.

## StarRocks Solution:

- Using StarRocks Primary key to solve update problem.
- Using compute node on k8s to auto-scaling.
- Get much more performance in ad-hoc query.

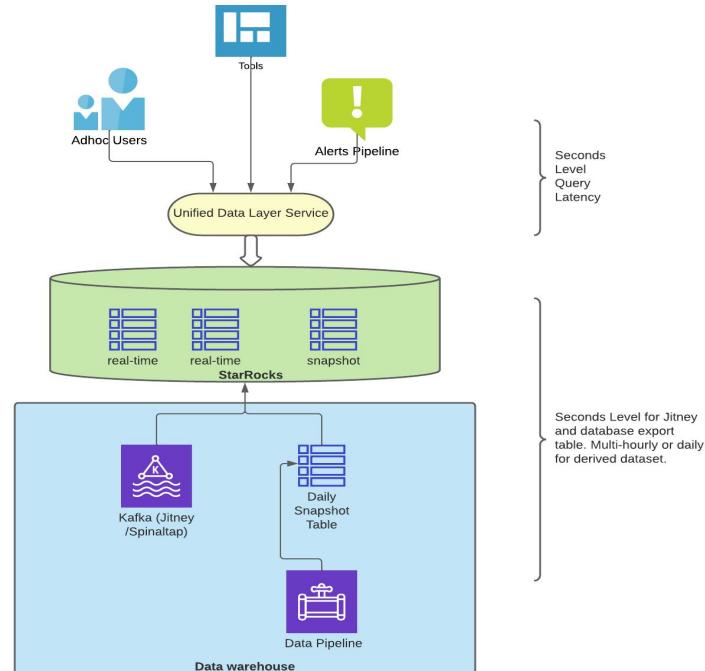


# Use Case: Trust Analytics at Airbnb

To enhance security, Airbnb needs a real-time fraud detection system (Trust Analytics) to identify various attacks and take actions ASAP. This system must support Ad-Hoc query and real-time update.

## StarRocks Solution:

- StarRocks hosts real-time updated datasets via Primary Key.
- Dataset import from Kafka has a sub-minute delay.
- StarRocks provides second-level query latency for complex joins.
- Alerting can be achieved by just running a SQL query regularly.



# Thank you.



● Community	<a href="https://starrocks.io">starrocks.io</a>
● Enterprise	<a href="https://celerdata.com">celerdata.com</a>
● Managed Service	<a href="https://cloud.celerdata.com">cloud.celerdata.com</a>

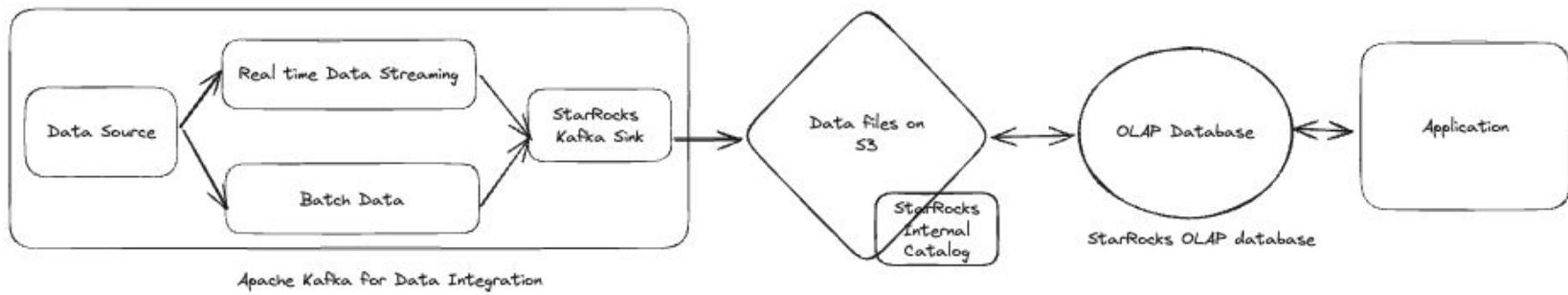
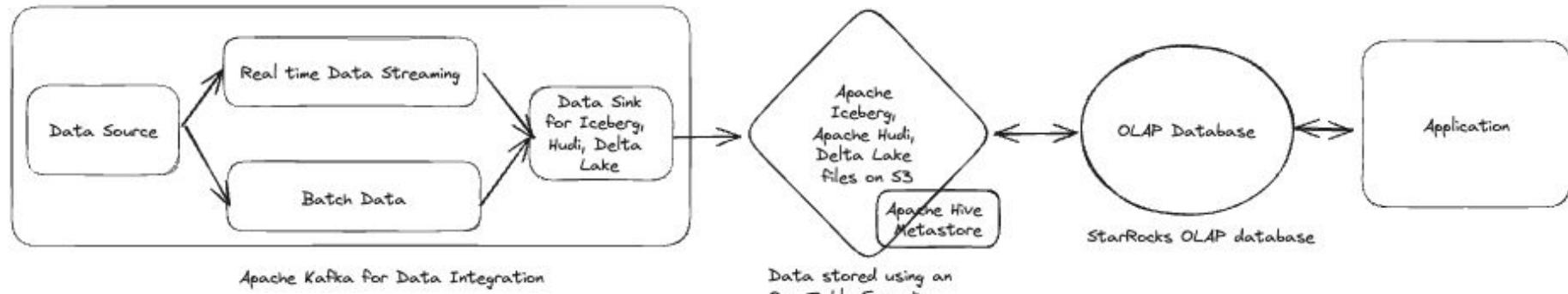
---

# Credits

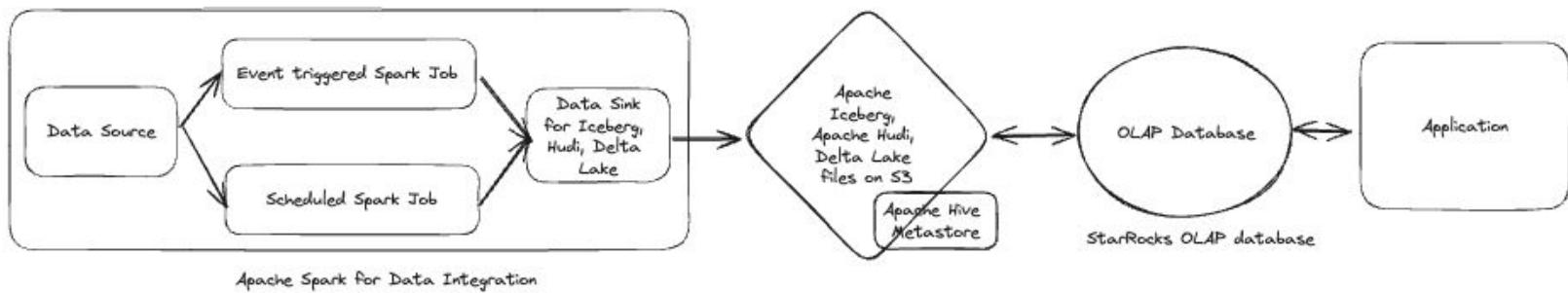
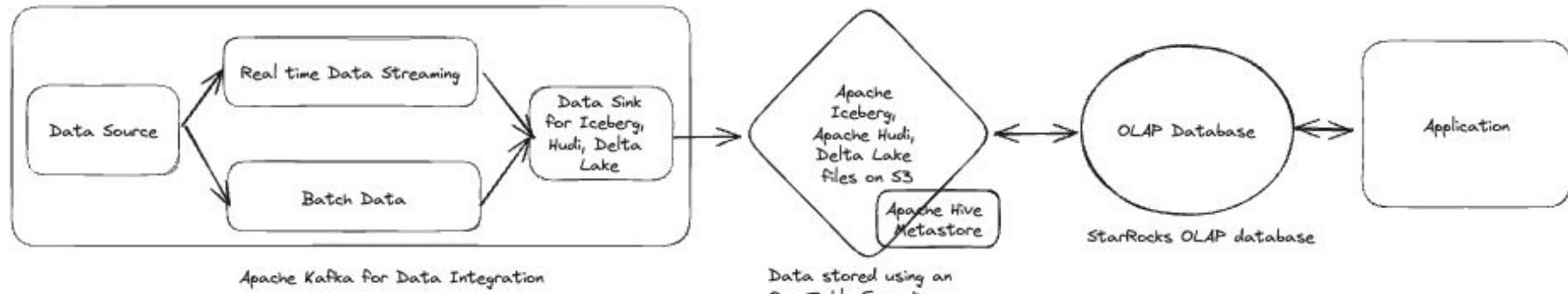
- This presentation is using images from Flaticon.com

# Architectural Patterns and Best Practices

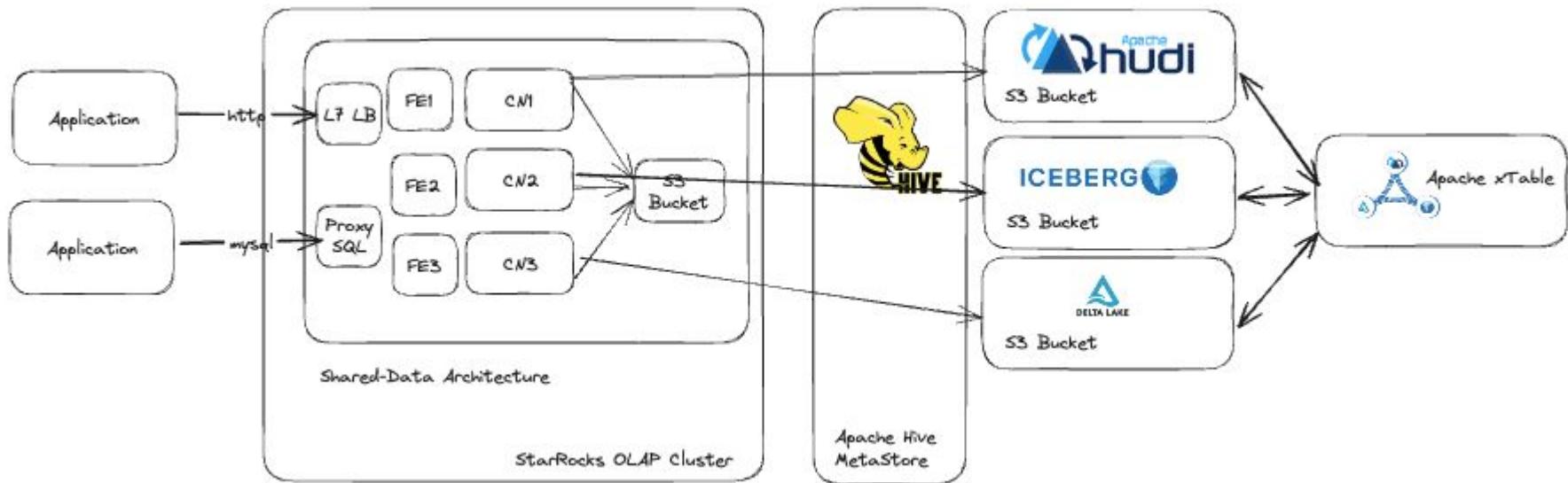
# Kappa and Lambda Architecture with StarRocks + Apache Kafka



# Kappa and Lambda Architecture with Open Lakehouse



# Open Data Lakehouse with Apache xTable



# Best Practices

