


CLICKHOUSE + APACHE ICEBERG:

READ & WRITE LAKEHOUSE TABLES

Unlocking the Future of Unified Data Analytics



28th Aug'25 | 8:30 PM IST

Presented by:  **OLake**



Shivji Jha

Staff Engineer,
Nutanix



Saurabh Ojha

MTS 2,
Nutanix

About Us

SHIVJI KUMAR JHA



Staff Engineer,
Data platform,
Nutanix

- **Areas of Interest**
 - Distributed Storage
 - Database & Streaming Internals
 - Application Architectures
- **Contributed to Clickhouse, MySQL, Apache Pulsar**
- **30+ talks at conferences & 10+ meetups**
 - Slides: <https://github.com/shiv4289/shiv-tech-talks#talks>
 - Recordings: <https://tinyurl.com/shiv-talks>
- **Text here:** [linkedin.com/in/shivjijha/](https://www.linkedin.com/in/shivjijha/)

SAURABH KUMAR OJHA



Software Engineer,
Data platform,
Nutanix

- **Areas of Interest**
 - Databases Internals
 - Streaming Internals
 - Linux kernel Features
- **Open-Source Enthusiast:** Contributions to Clickhouse server and Ecosystem, Nats, Transferia and more
- **Mostly active here:**
 - [linkedin.com/in/ojhasaurabh2099/](https://www.linkedin.com/in/ojhasaurabh2099/)



Agenda

Setting stage : Definitions

Let's play with a Demo!

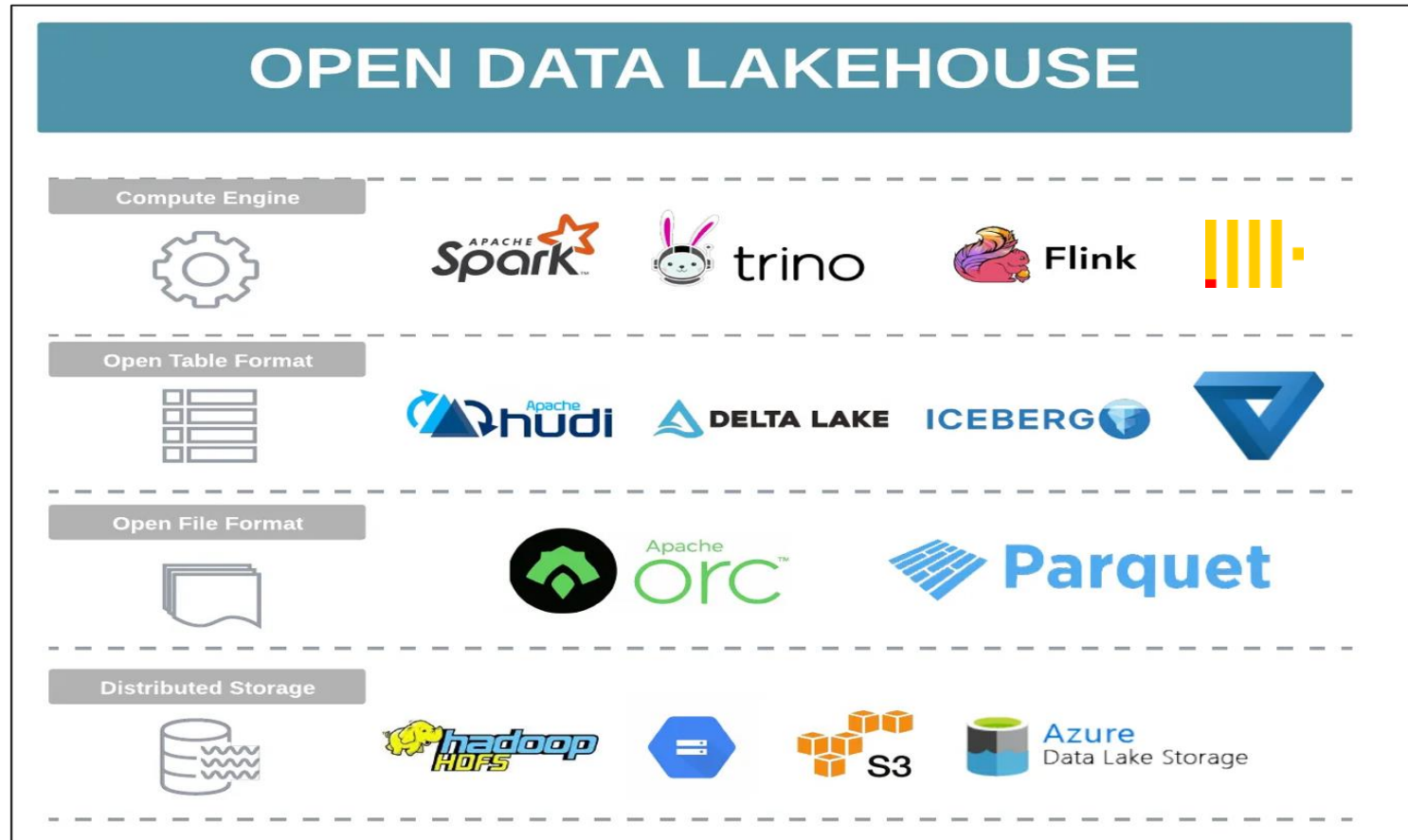
Architecture Patterns & Use cases

Lay of the land today: Clickhouse & Iceberg

Call for contributors

Wishlist for Future

The Lakehouse stack



<https://alirezasadeghil.medium.com/the-history-and-evolution-of-open-table-formats-0f1b9ea10e1e>

What is Apache Iceberg?

- Open Spec for table format, not a database engine
- SDKs with reference implementation, not for production
- Use Spark, Trino, DuckDB, Clickhouse, AWS S3 Tables etc..
- **Supports**
 - ACID transactions.
 - Time travel on data
 - Schema Evolution
 - Partition Evolution



Why Iceberg?

Interoperability across engines (Spark, Flink, Trino, Clickhouse)

Strong ecosystem: adoption & community

Own your data, not locked in proprietary format

Decoupling of **storage** and **compute**

Even stream-native platforms like Kafka (TableFlow), Pulsar(Ursa), and RisingWave are adopting Iceberg

A lot of mindshare at the moment!

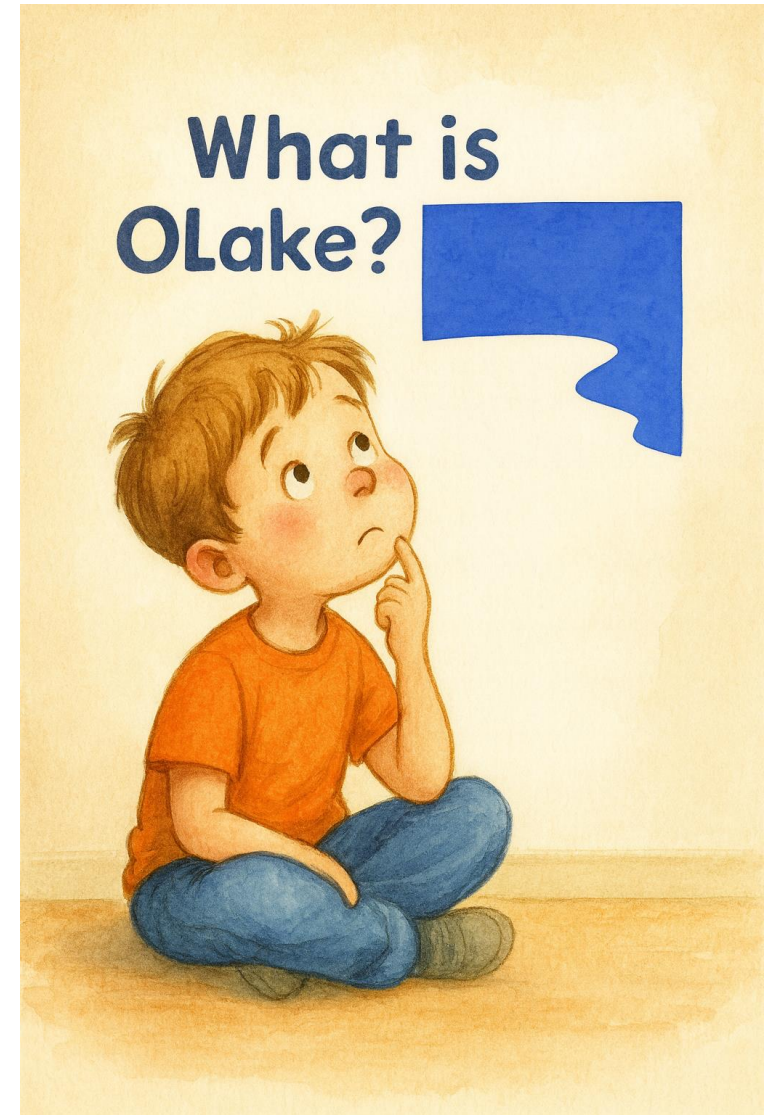
What is Clickhouse?

- Open source* OLAP database known for speed
- Very easy to get started
 - *one binary, regular SQL/Tables like RDBMS*
 - *Very popular, open, well adopted!*
- Storage Compute segregation(?)
 - *Possible with Open source, not trivial!*
- But a lot of integrations (engines, format Readers etc)
- Supports parquet, iceberg and now even writes (experimental)



What is OLake (CDC into iceberg)?

- Open-source, high-speed **Change Data Capture** tool
- Syncs data from Postgres, MySQL, MongoDB, etc.. → Iceberg
- **Ultra-fast throughput**
- Lightweight: no Spark or Flink required
- Supports both **full load + CDC** with schema evolution
- Simple to run via Docker & easy UI/CLI



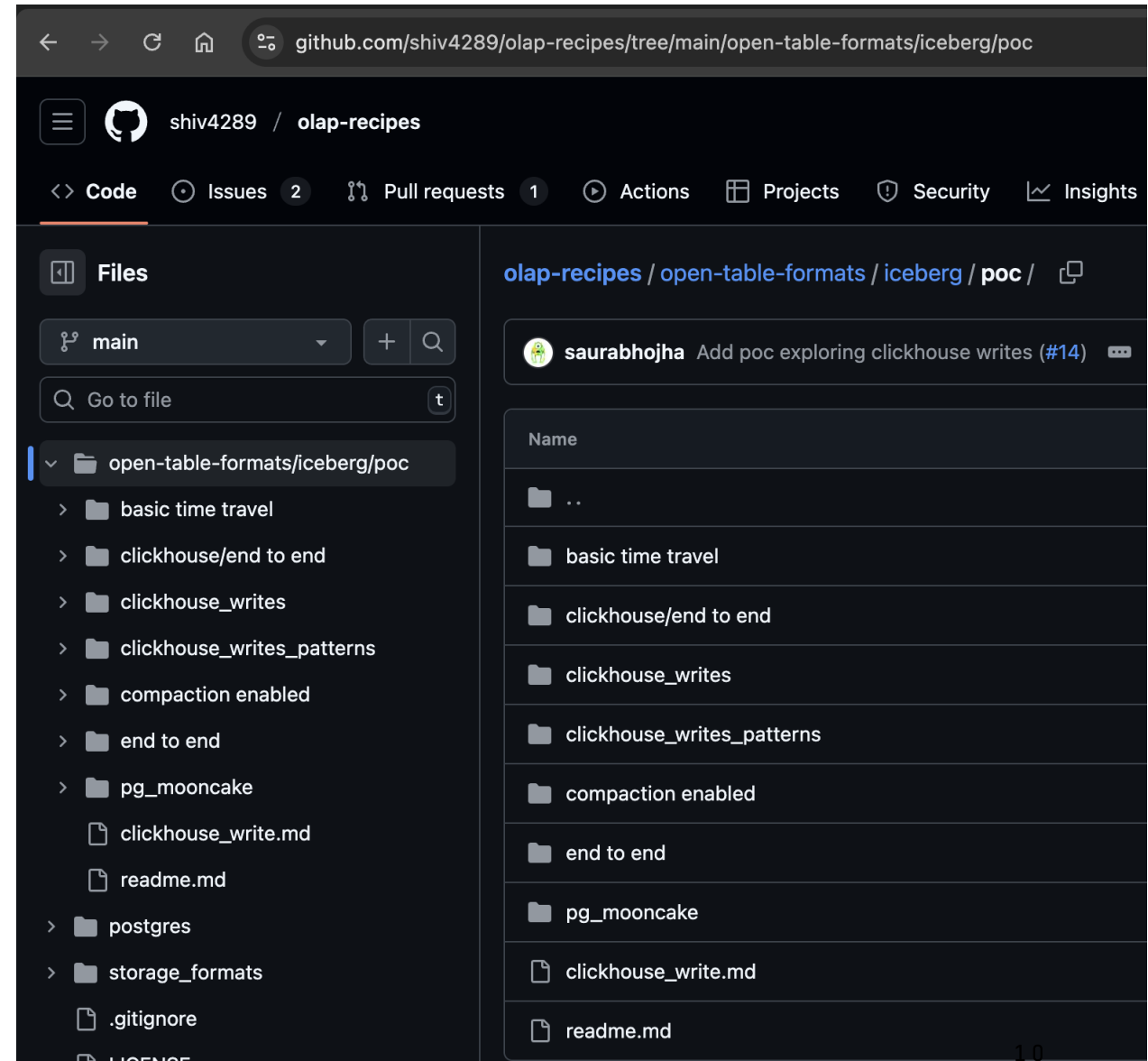
Demo Time!!



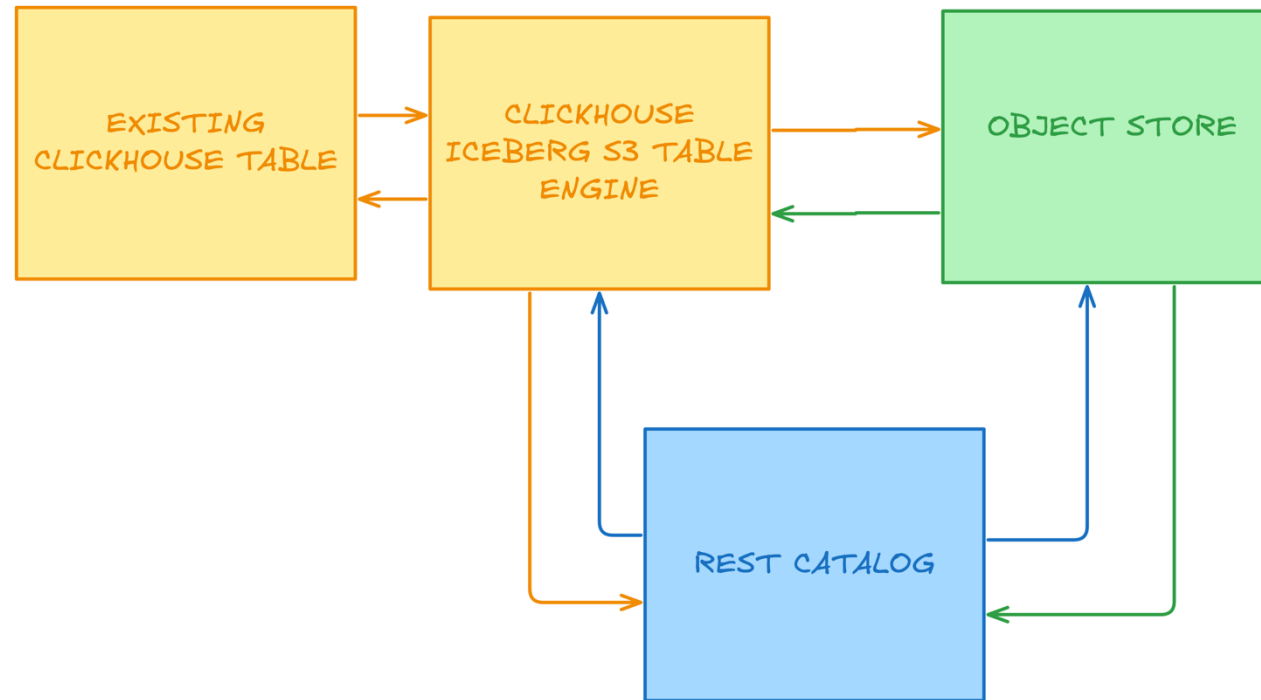
Want to get your hands dirty?

Prerequisites:

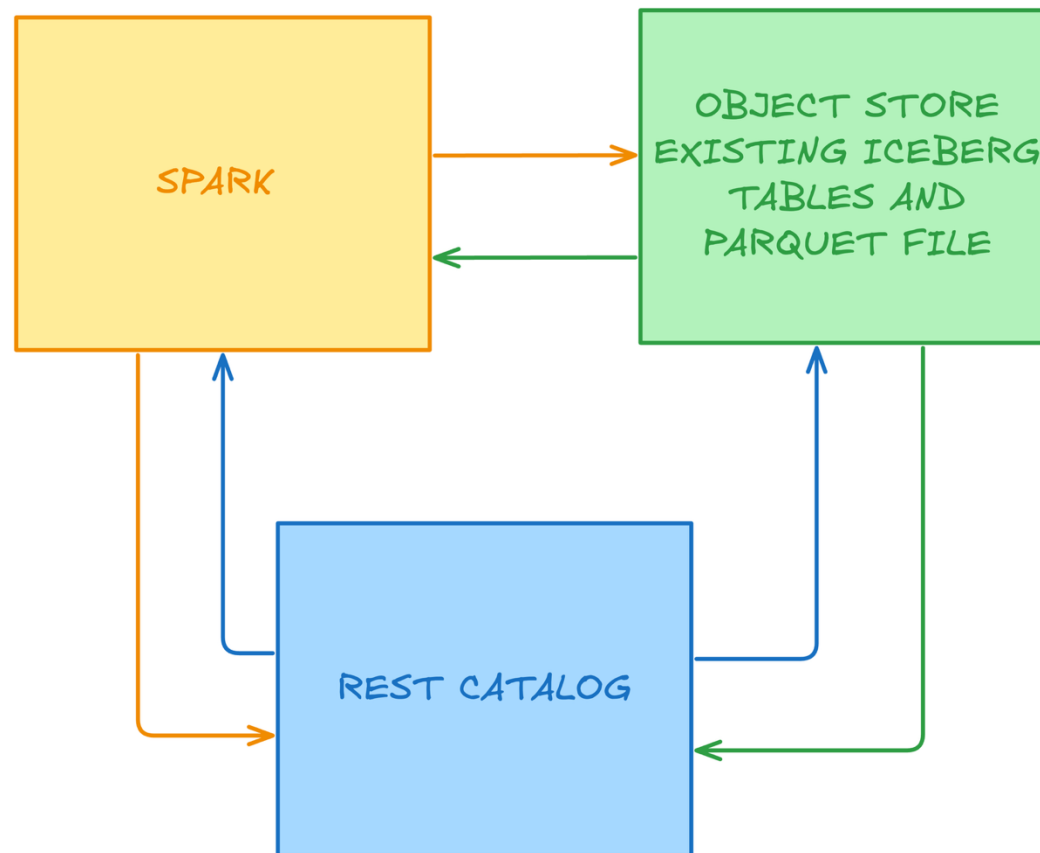
- Have a docker or equivalent setup.
- Linux / Mac OS
- Try the demo yourself !! [Find it here](#)



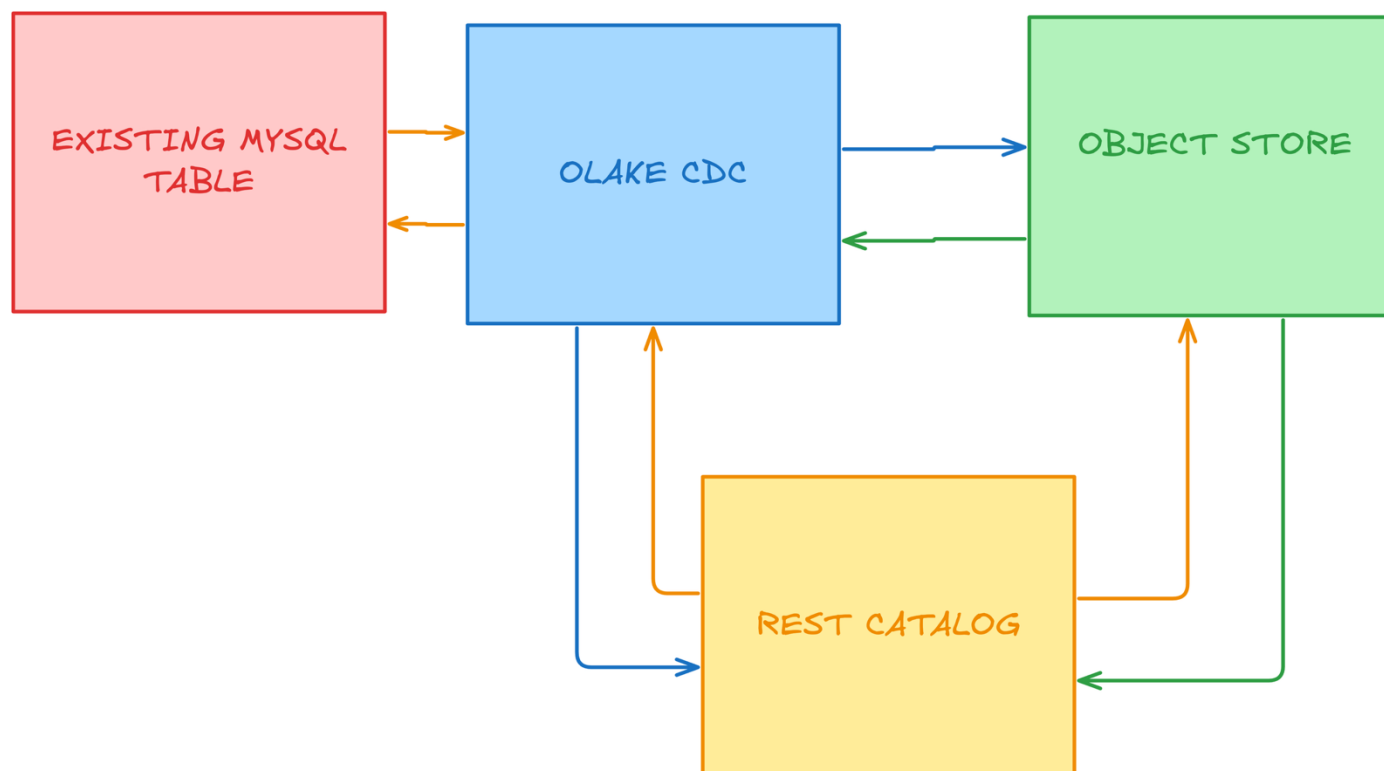
FLOW 1: DATA FROM CLICKHOUSE TO ICEBERG TABLE



FLOW 2: WRITING TO SAME ICEBERG TABLE USING SPARK



FLOW 3: WRITING TO THE SAME ICEBERG TABLE USING OLAKE CDC FROM MYSQL








Fresh off the oven (experimental)


Did a bunch of plumbing to get here!


We've already merged these into the Clickhouse source code, so you can avoid running into the same challenges!!

 **Fix incorrect S3 metadata path resolution when config is populated in Iceberg REST catalog /v1/config response**
✓ can be tested pr-bugfix pr-synced-to-cloud submodule changed
#80562 by saurabhajha was merged on Jul 18

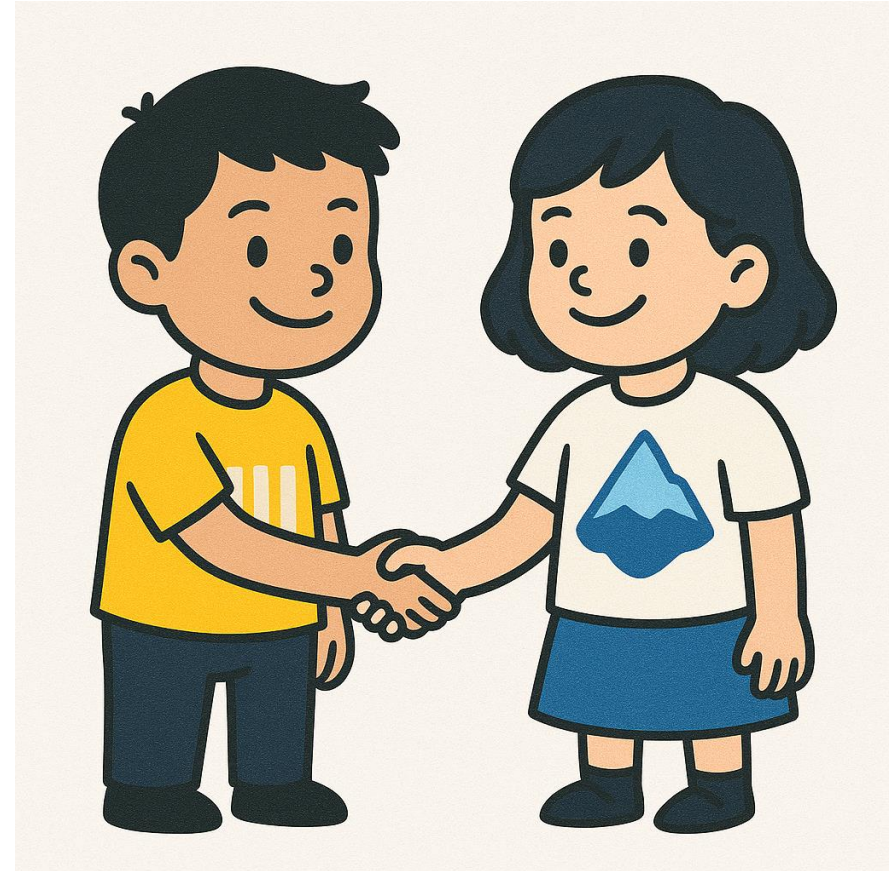
 **Add Nessie catalog integration tests for ClickHouse Iceberg support** ✓ can be tested pr-not-for-changelog
#85128 opened 3 weeks ago by somratdutta

 **Add Lakekeeper catalog support in docs** ✗
#4177 by somratdutta was merged 2 days ago • Approved 2 of 3 tasks

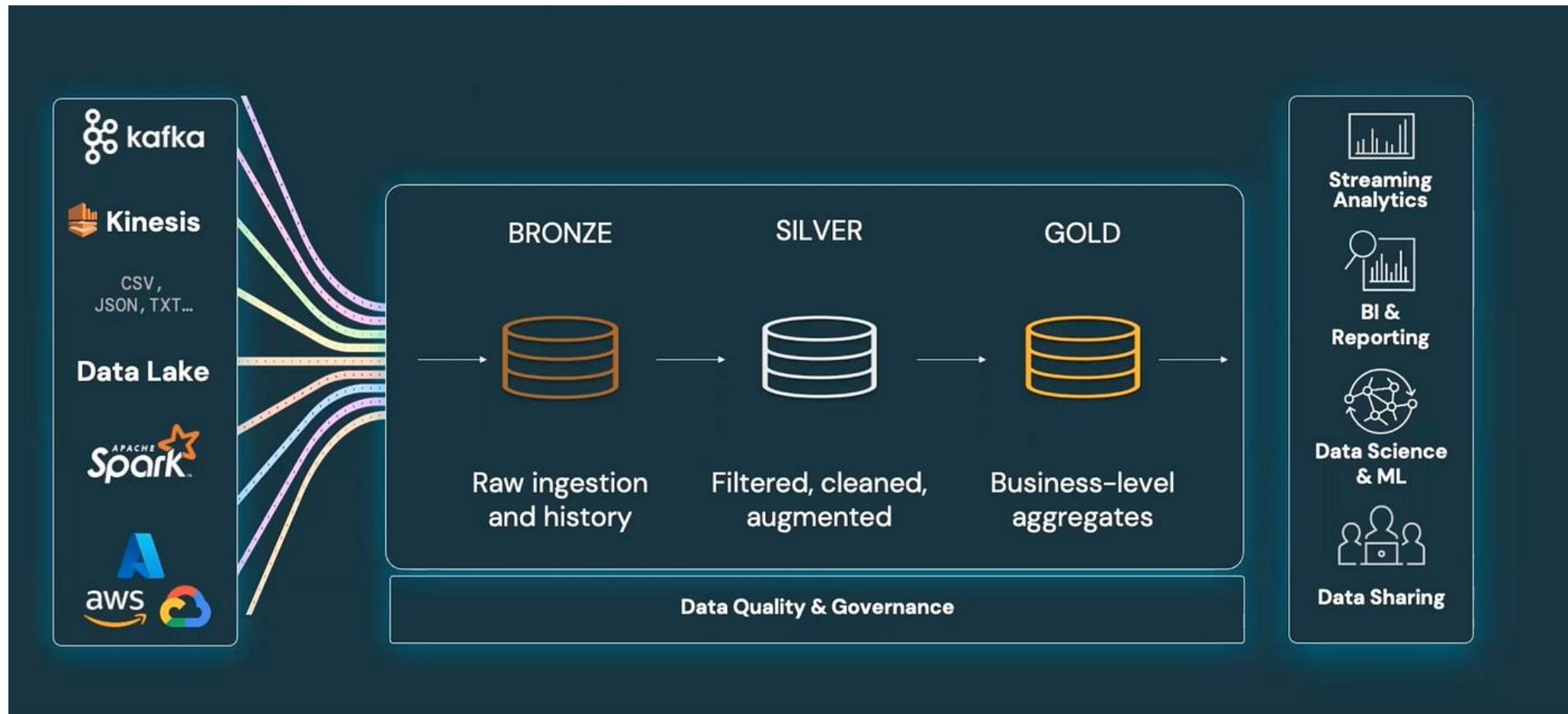
 **Add Nessie catalog support in docs** ✗
#4180 by somratdutta was merged 2 days ago • Review required 2 of 3 tasks

 **Add REST catalog support in docs (#1)** ✗ trademark-addendum-signed
#4031 by somratdutta was merged on Jul 21 • Approved 2 of 3 tasks

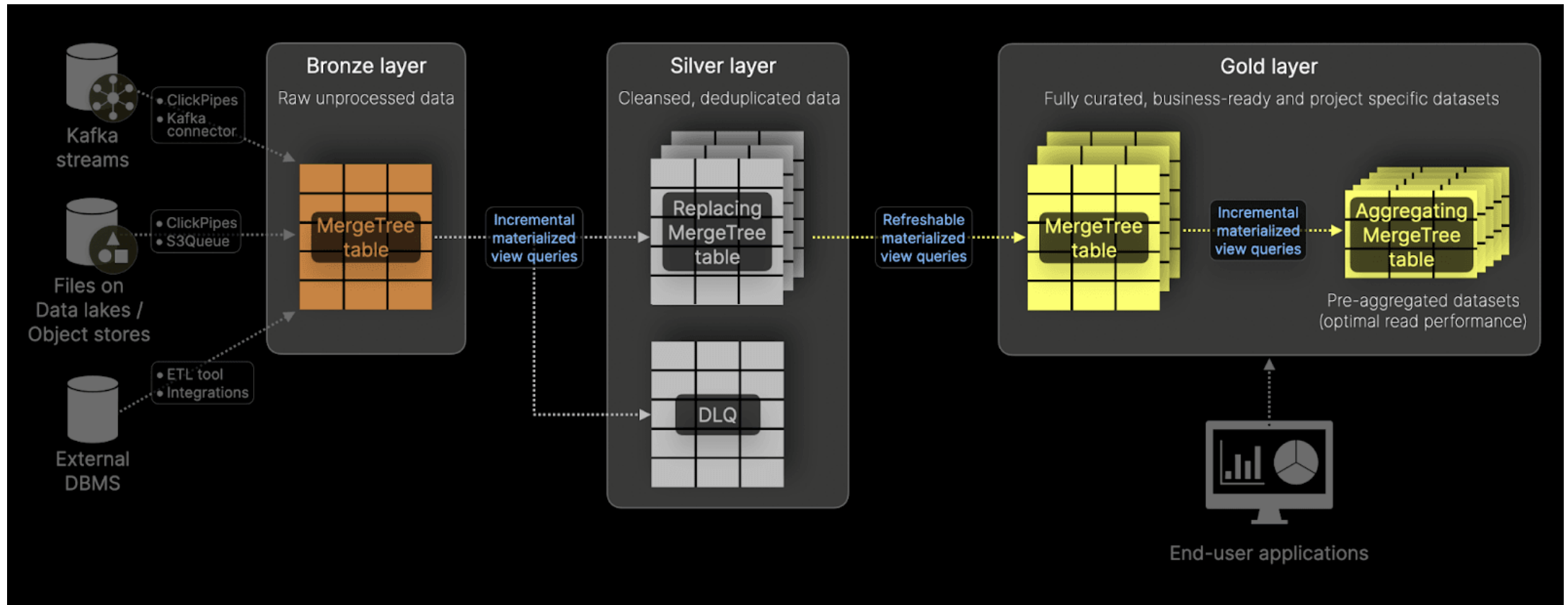
Design Patterns



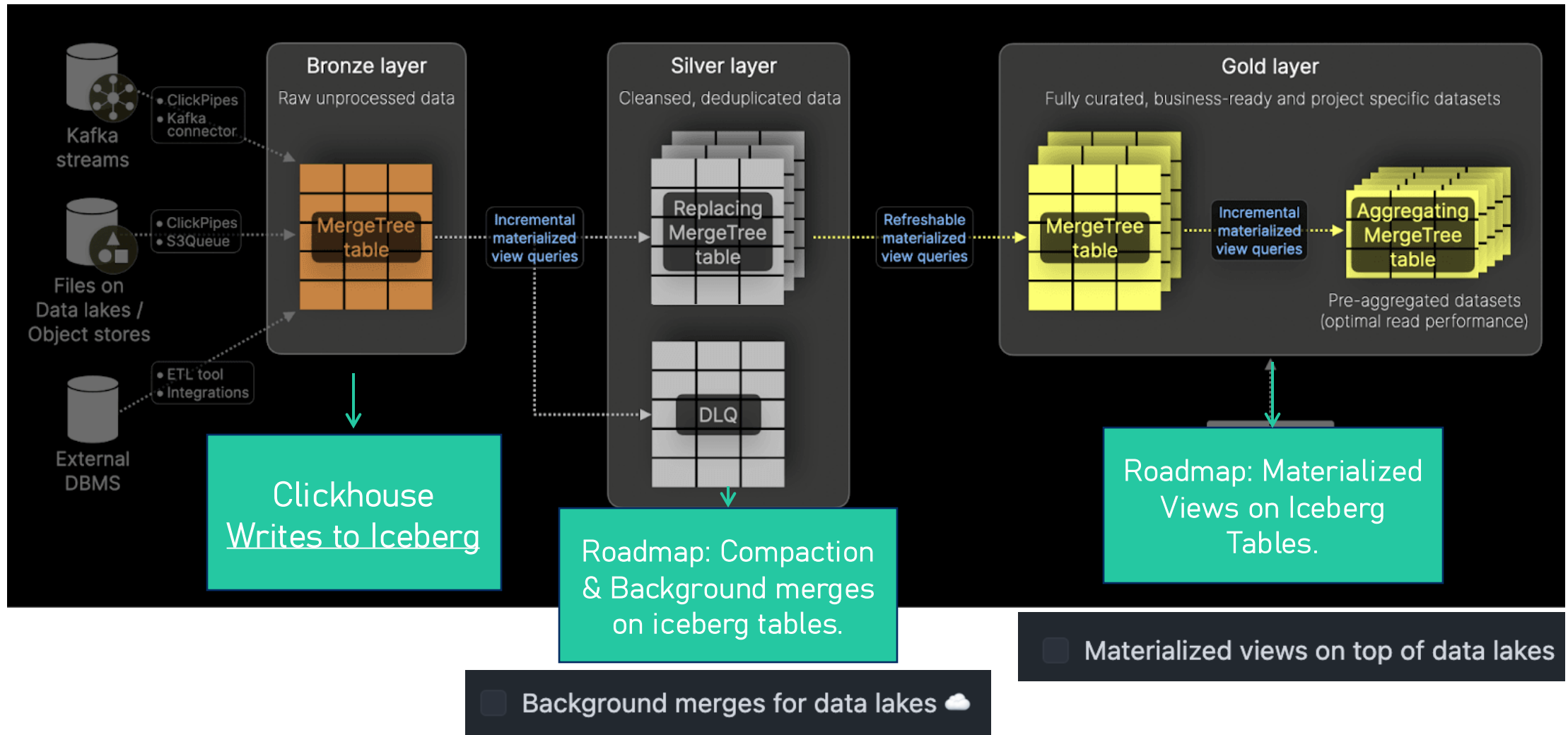
Medallion Architecture



Medallion Architecture ft. Clickhouse TODAY



Medallion Architecture ft Clickhouse / Iceberg



Use cases



Realtime (Fraud detection- eg: credit card fraud)



BI dashboards (quarterly sales trends)

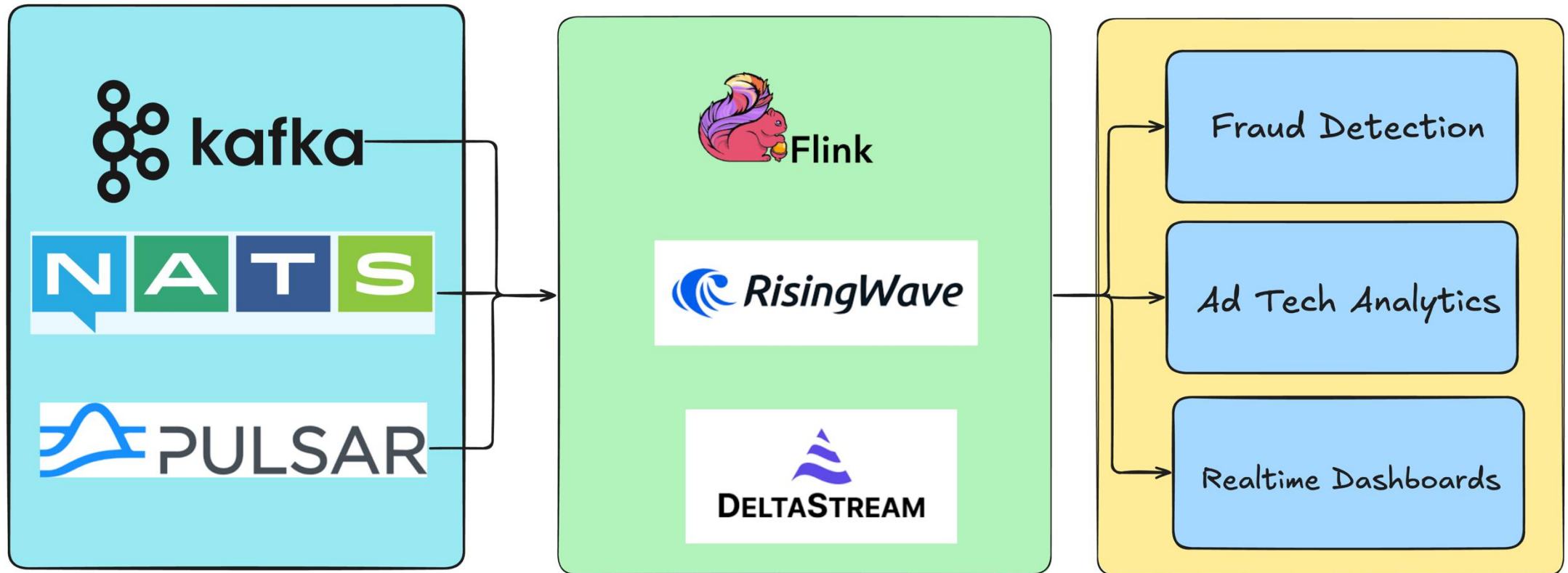


ML training pipelines

Architecture blueprints



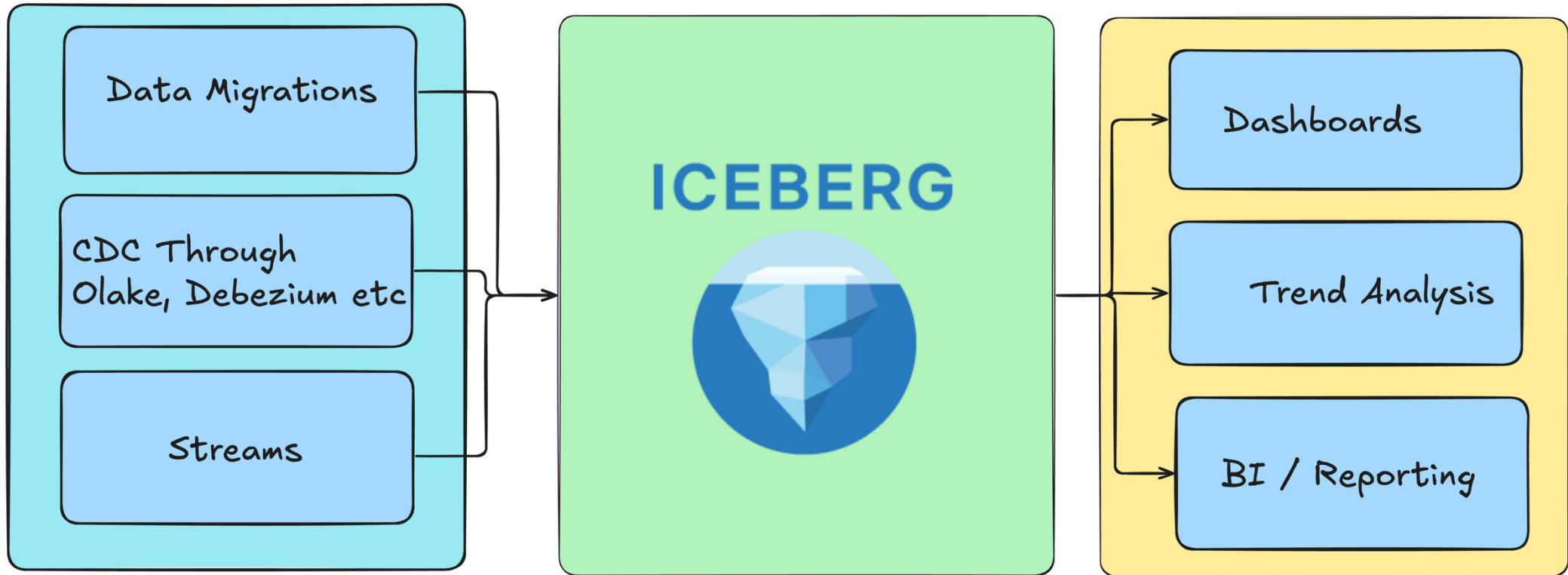
Realtime (Fraud detection-
eg: credit card fraud,
Clickstream Analytics etc..)



Architecture blueprints



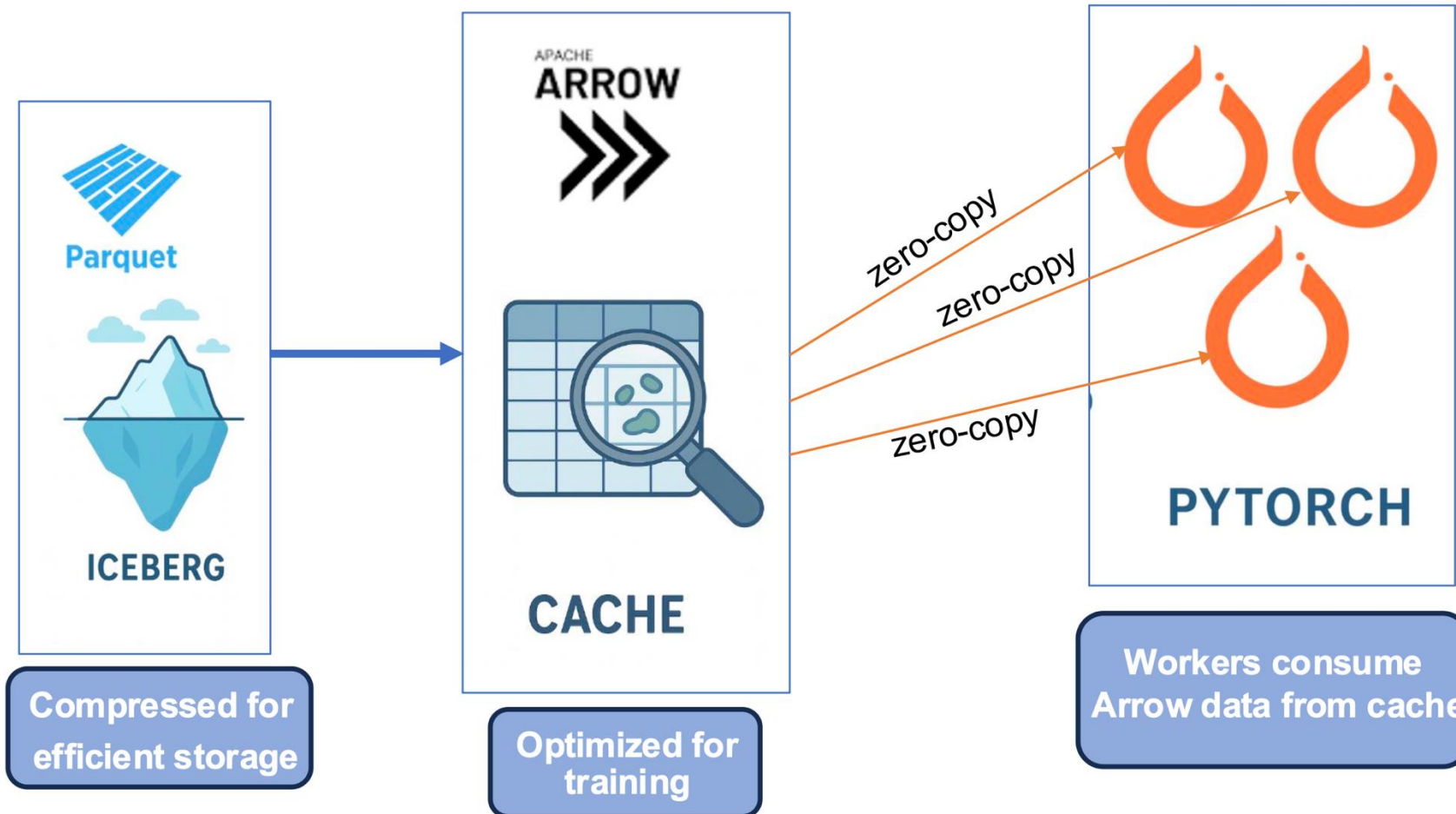
BI dashboards (quarterly sales trends)



Architecture blueprints



ML training pipelines



Callout for contributors



Experimental (preview), fresh off the oven!



By default, fields are nullable. No **required** fields yet



Doesn't support partitioning while creating tables.



Can't attach to an already existing table -issue [here](#).



No deletes on Iceberg tables – only INSERTs



Can't read compressed metadata files. Fixed in master!

Clickhouse Roadmap

[Background Merges](#) – Compactions & clustering for better performance.

Introducing external materialized views for better query efficiency

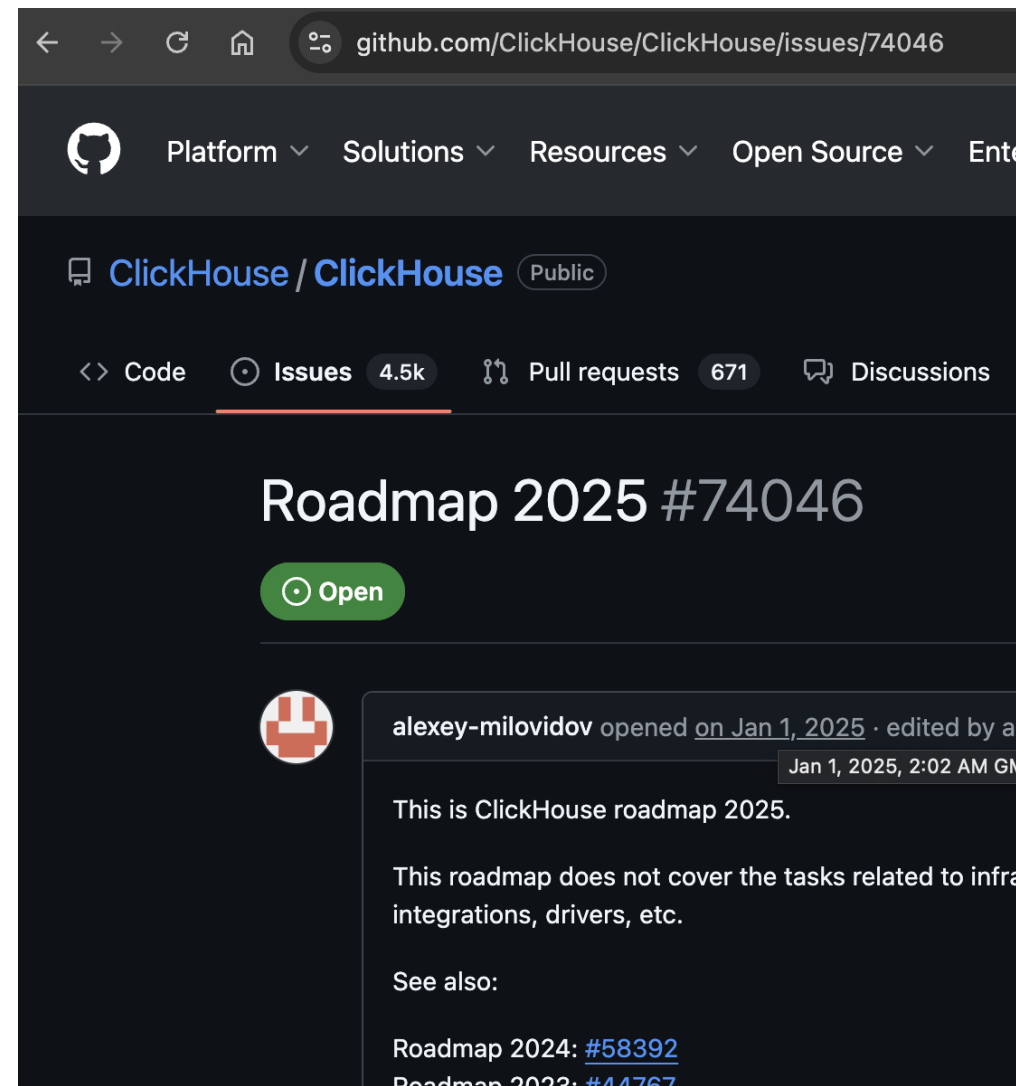
Support for Equality Deletes.

Schema updates directly from Clickhouse.

A [brand new parquet reader](#) from scratch!!

Iceberg CDC Connector in ClickPipes (Iceberg to clickhouse native table)

<https://github.com/ClickHouse/ClickHouse/issues/74046>



Takeaways

Iceberg is taking off! Who wouldn't want to own data, right?

Proof of the pudding: A lot of vendors are integrating with Lakehouses

We use Clickhouse & are looking closely at Clickhouse integrations.

No better time to be a builder - Freedom of choice for query engines.

No better time to be an open-source contributor - Checkout roadmap.



Questions?

Get in Touch:

[linkedin.com/in/shivjijha/](https://www.linkedin.com/in/shivjijha/)
[linkedin.com/in/ojhasaurabh2099/](https://www.linkedin.com/in/ojhasaurabh2099/)