

HANOI UNIVERSITY OF
SCIENCE AND TECHNOLOGY

scalable Lakehouse for Financial Analy

Big Data Storage and Processing Capstone Project

Trinh Hoang Anh
Bui Khac Chien
Bach Nhat Minh

Tran Quang Minh
Vo Ta Quang Nhat

Supervised by: Dr. Tran Viet Trung

Contents

1. Introduction & Problem Statement
2. The Architectural Blueprint
3. Infrastructure & Deployment
4. Notable Lessons Learned
5. Conclusion & Key Achievements

The Challenge: Taming Financial Data

Financial markets generate data at extreme volume, velocity, and variety, creating significant engineering challenges that demand a unified architecture.

Heterogeneous Data Sources

Integrating real-time streams (e.g., stock quotes) with historical datasets from multiple sources, GDELT global

Dual Analytical Demands

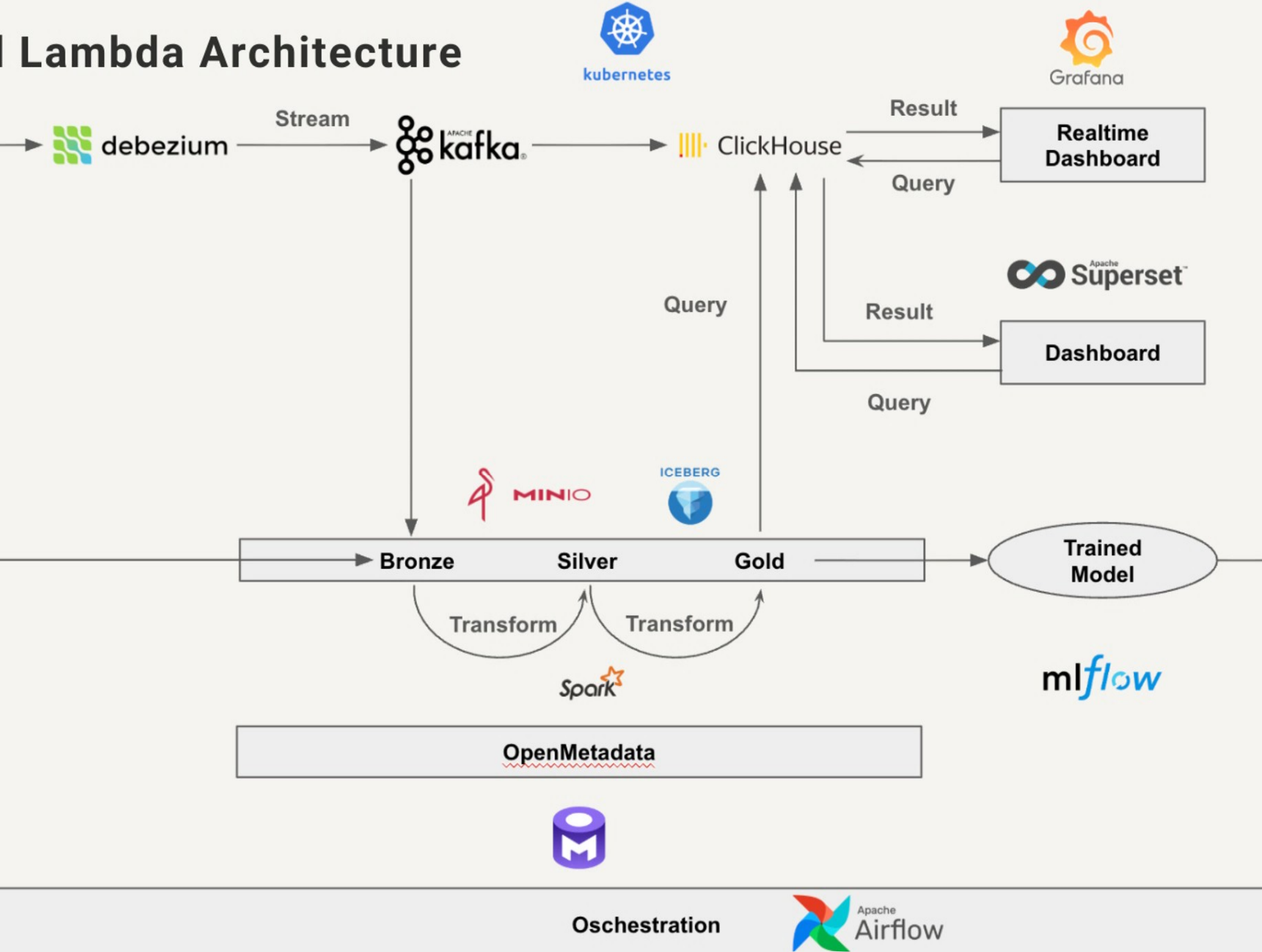
Supporting low-latency dashboards for real-time monitoring **and** deep, historical analysis for ML models.

Data Integrity & Governance

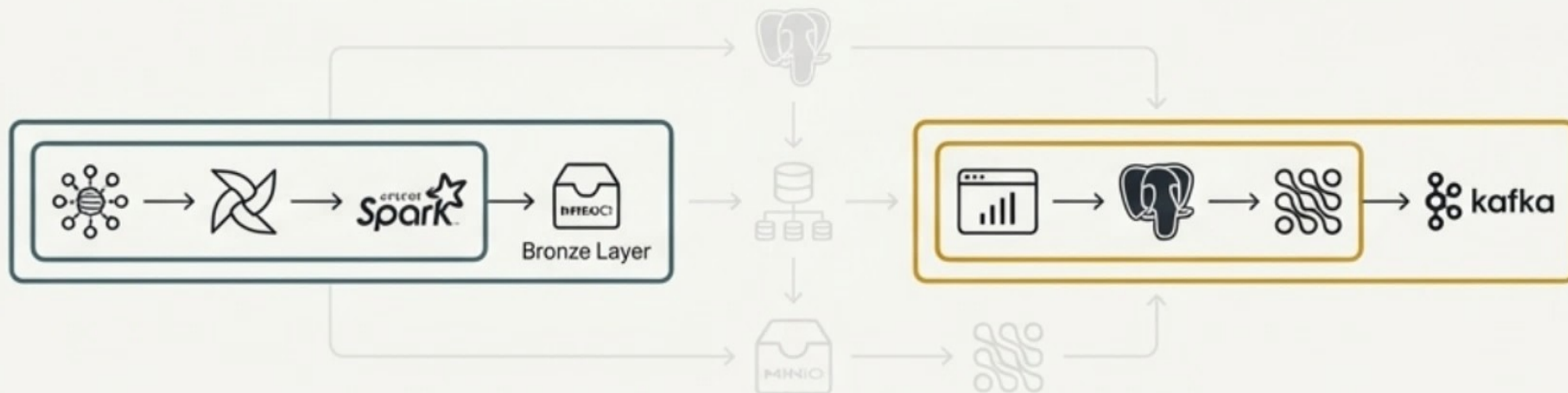
Ensuring consistency, quality, and traceability across a massive, evolving dataset from raw ingestion to aggregated insights.

Mission: To design and implement a scalable Big Data Lakehouse capable of ingesting, processing, and analyzing financial data in both real-time and batch modes.

Lambda Architecture



Question: A Hybrid, Purpose-Built Strategy



Stooq (High Throughput)

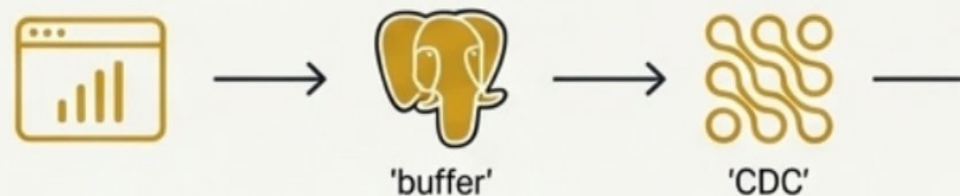
Stooq.



'Direct-to-Lake' approach bypasses relational databases to increase throughput and avoid I/O bottlenecks for bulk historical data.

Speed Layer (Low Latency & Integrity)

Source: Finnhub API.



Principle: Uses PostgreSQL as a transactional buffer for deduplication. Debezium's Change Data Capture on the buffer creates a real-time, event-driven stream without impact on the source.

Processing: The Medallion Lakehouse on Spark & Iceberg

v)

Raw data from sources (GDELT, Stooq) stored as-is. Partitioned by date for efficient writes. The raw data is available for reprocessing.

Curated & Enriched)

Validated, structured, and deduplicated data. Key transformations include schema enforcement, time zone conversion (UTC to EST), and linking news events to specific stock tickers.

Aggregated)

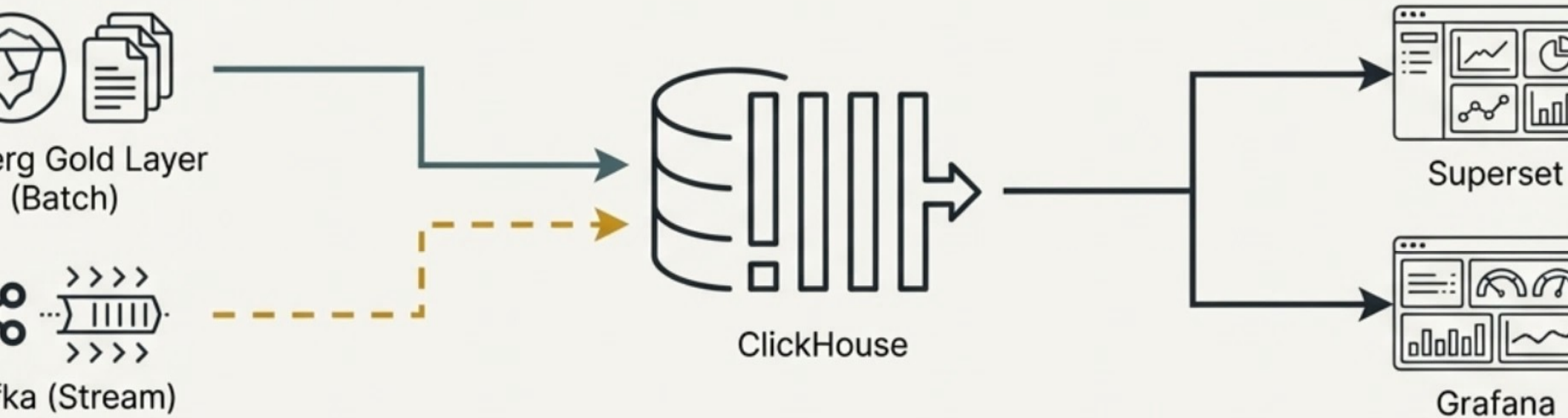
Pre-computed aggregates optimized for analytics and ML. Examples: daily sentiment scores, news volume, and sentiment by ticker.

Storage - Apache Iceberg

On top of S3, Iceberg provides ACID transactions, schema evolution, and time-travel directly on the data lake, combining the reliability with data lake scalability.

Serving: Unifying Batch and Stream with ClickHouse

se provides a single, high-performance OLAP access point for both deep hi
low-latency real-time streams.



Strategies:

Data: Utilizes the Iceberg Table Engine for **zero-copy queries** directly on Parquet files in M
gated data is also loaded into 'MergeTree' tables for maximum dashboard performance.

on Data: Ingests high-throughput streams from Kafka into optimized 'MergeTree' tables des
eries queries.

Sub-second query performance for interactive dashboards and immediate visibility of mar

Stack: A Containerized & Cloud-Native Foundation

The system is built on a containerized, cloud-agnostic infrastructure, enabling flexible and scalable deployments.

Container & Orchestration



Processing & Storage



Machine Learning



m

Flow & Metadata



Streaming & Database



Core Principles: Infrastructure as Code (IaC) • Scalability • Fault Tolerance

able Lesson #1: Hybrid Ingestion is a Necessity Choice

Problem

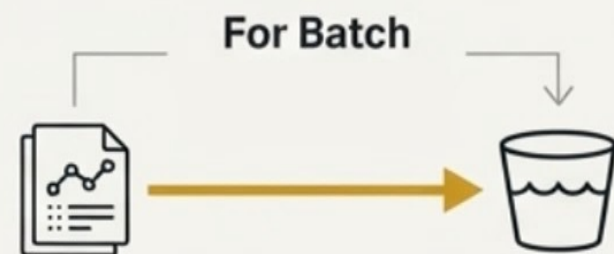
Clash Grotesk Medium

A "one-size-fits-all" ingestion approach created severe bottlenecks.

Routing massive batch files (GDELT) through PostgreSQL caused I/O overhead, while writing API streams directly to the lake risked data duplication.

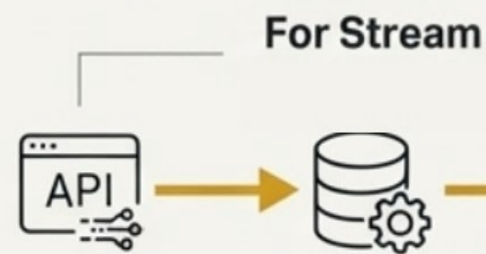
Solution

Clash Grotesk Medium



Direct-to-Lake

Spark writes directly to MinIO, bypassing transactional overhead for maximum throughput.



Transactional Buffering

API -> Postgres -> CDC, using database for its strengths in integrity and deduplication during streaming.

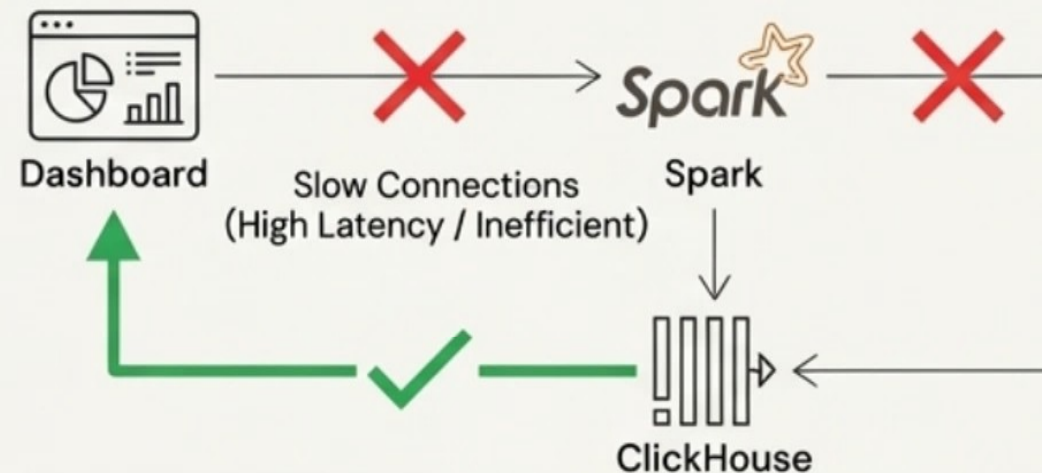
Takeaway

Use a database in the middle unless you have a transactional reason. Use raw object storage for bulk throughput.

Table Lesson #2: Separate Processing from Serving Low-Latency BI

Dashboards were sluggish and timed out.
Connecting Superset directly to Spark was
slow (high latency), and using PostgreSQL for
large-scale aggregations was inefficient (row-
oriented).

Solution



Dedicated OLAP Serving Layer

Implemented **ClickHouse** as a dedicated **OLAP serving** layer. Columnar storage and vectorized query execution are perfect for the fast aggregation queries required by BI dashboards, providing sub-second responses.

Takeaway

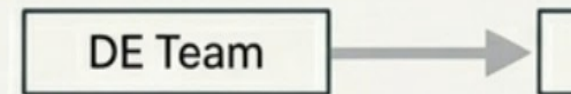
Use high throughput (ETL processing). ClickHouse is for latency (interactive serving). Use the right engine for the job.
Satoshi Regular

Table Lesson #3: Data Schemas as Contracts Engineering Work

The visualization team was blocked, waiting for the data engineering pipeline to be completed. This forced a slow, sequential "waterfall" process.

Solution

Traditional Waterfall



Contract-First Development



Adopted **Contract-First Development**. The teams agreed on a "Gold" layer data schema (column names, types, formats). The visualization team built dashboards using mock data matching the contract, while the DE team built the pipeline to produce data that required zero changes to the dashboard logic.

Takeaway

A shared data contract decouples teams, is the foundation for agile development in data projects, and reduces integration risk.

Lesson #4: Robust Deployment Requires Infrastructure

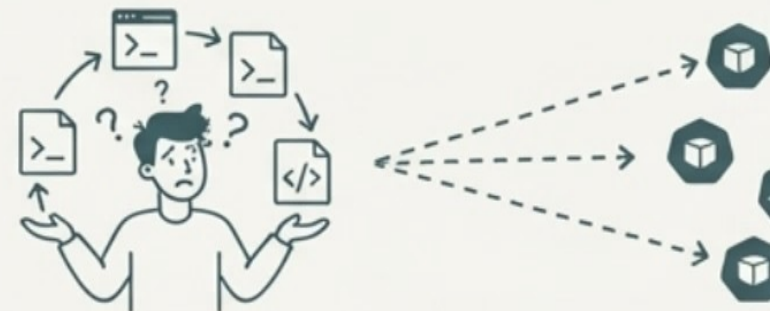
Managing 15+ Kubernetes resources with manual `kubectl apply` commands and shell scripts was brittle, error-prone, and lacked versioning.

Takeaway

Infrastructure as versioned code (IaC), configuration drift and making complex systems manageable, resilient, and repeatable.

Solution

Before



After



Migrated all Kubernetes YAML manifests into a single, unified **Helm chart**. This templated the entire application, centralizing configuration in a single `values.yaml` file and enabling versioned deployments and rollbacks with a single command.

Conclusion: An End-to-End Financial Data Platform

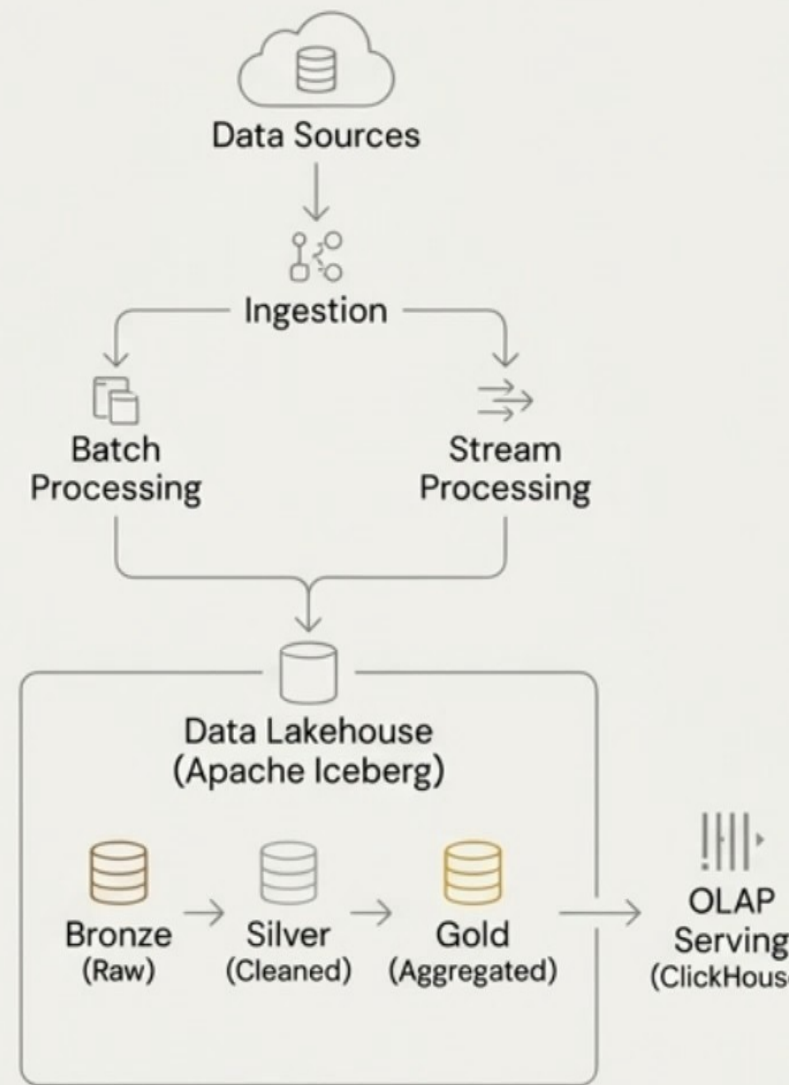
Successfully designed and implemented a scalable, modern Big Data Lakehouse that addresses the challenges of complex financial analysis.

Unified Architecture: Integrated batch and stream processing via a Lambda pattern.

Quality & Reliability: Ensured curated, reliable data with the Medallion model on Apache Iceberg.

Performance Analytics: Delivered low-latency dashboards through a dedicated OLAP serving layer (House).

Secure & Reproducible System: Built on a modern, containerized infrastructure managed with Infrastructure as Code (Helm).



Simplified Architecture Overview