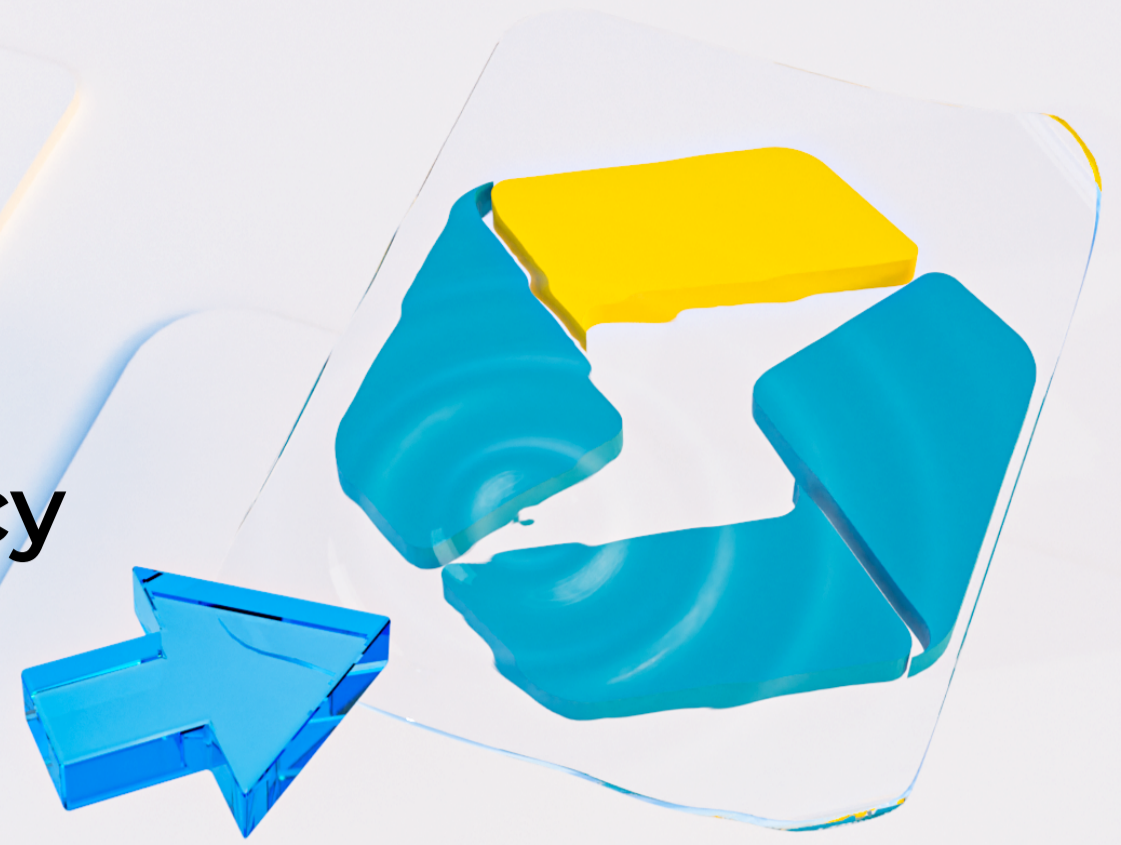




How WeChat's Data Lakehouse Architecture Enhances Efficiency for Trillions of Daily Records



Trillions

daily records

2h

shorter development cycle for offline tasks

65%

storage cost

About WeChat

WeChat is the world's largest standalone mobile app, serving over 1.3 billion monthly active users as a platform for instant messaging, social media, and mobile payments. To support its unprecedented and rapidly expanding user base, WeChat's technological backend has had to evolve quickly as well, transitioning from a Hadoop + data warehouse architecture to a modern open data Lakehouse architecture.

Challenges

With over a billion users, it comes as no surprise that WeChat manages extremely large data volumes. In some cases, single tables are growing by trillions of records daily and queries regularly scan over 1 billion records.

WeChat's business scenarios demand rapid end-to-end response times, with a query latency P90 target of under 5 seconds, and data freshness requirements that vary from seconds to minutes. This complexity is elevated by the need to process often more than 50 dimensions and 100 metrics at a time.

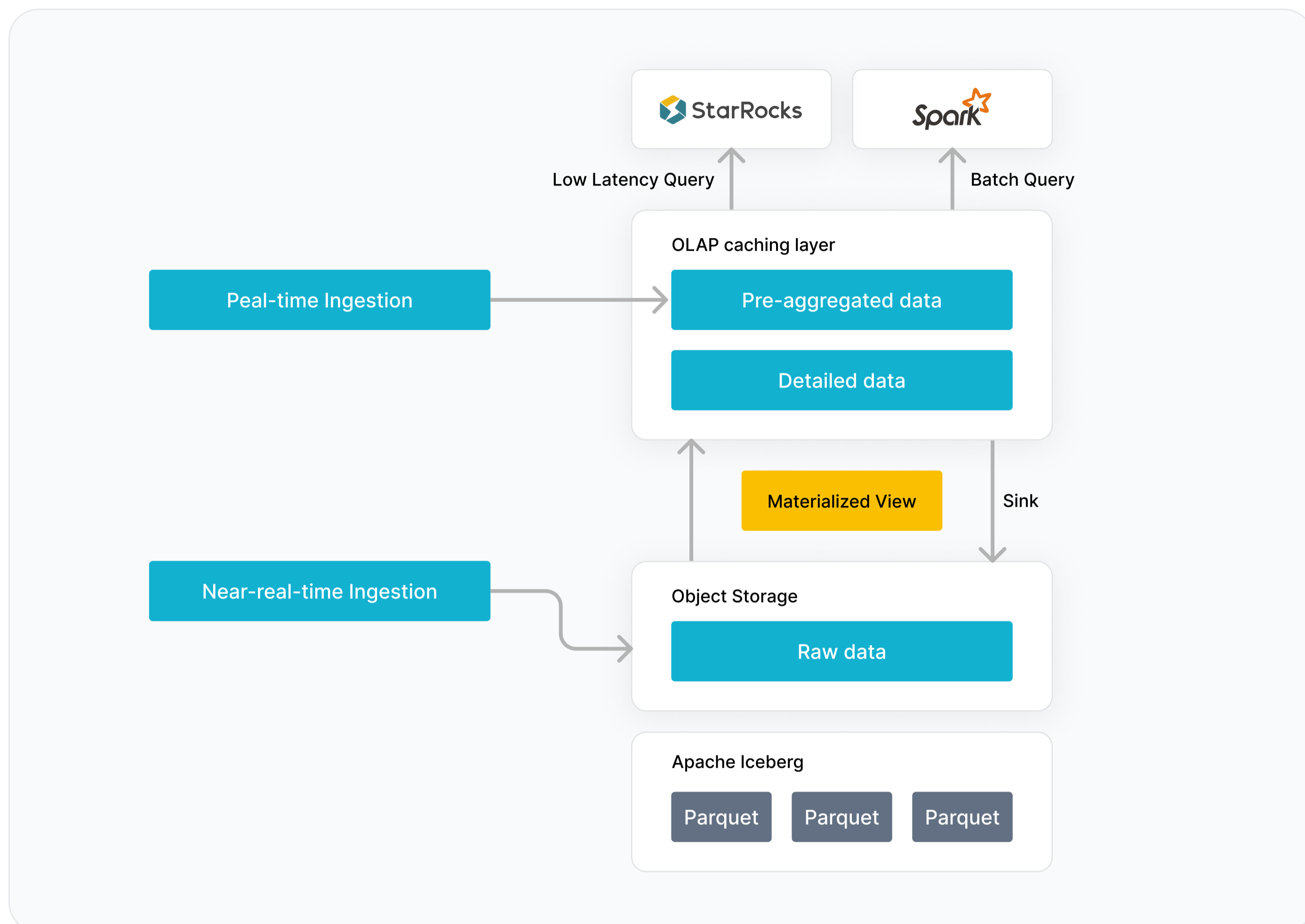
WeChat's legacy data architecture involved a Hadoop-based data lake system along with a variety of data warehouses. This resulted in significant operational overhead and data governance challenges including:

- Juggling multiple systems from separated real-time and batch analytics pipelines
- Maintaining data ingestion pipelines for data warehouses
- Governance challenges from managing multiple copies of the same data
- Managing incompatible APIs of different systems
- Challenges in standardizing data analysis processes

To address these issues, WeChat pursued a unified approach to their analytics which necessitated a redesign of their data architecture.

Solution

WeChat's revamped Data Lakehouse architecture now features StarRocks as the low latency query engine with Apache Spark as its batch processing engine. Data is stored as Parquet files on Tencent Cloud's COS (cloud object storage) with Apache Iceberg as the data lake table format.



This new architecture supports both real-time and near-real-time data ingestion:

- For real-time ingestion: data is first ingested into a warehouse in real-time, then cold data is sunk into the data lake and queries can automatically union cold and hot data.
- For near-real-time ingestion: raw data is directly ingested into Apache Iceberg, and then cleaned and transformed using [StarRocks' materialized view](#).

Result

WeChat's StarRocks-based data lakehouse solution is now in production across multiple business scenarios within the company including livestreaming, WeChat Keyboard, WeChat Reading, and Public Accounts.

By unifying all workloads with one system, WeChat has experienced significant operational benefits from this improved efficiency. Their

live streaming business is one example: the new lakehouse architecture halved the number of tasks data engineers are required to manage, reduced storage costs by over 65%, and shortened the development cycle of offline tasks by two hours.

Not only is their architecture now simplified, data freshness and query latency improved as well, with batch ingestion being eliminated and their near-real-time data pipeline and query latencies being brought down to the mostly sub-second level.

Table 1. WeChat Data Lakehouse Performance Numbers

Concurrency	Near-Real-Time Pipeline	Real-Time Ingestion Pipeline
Data Freshness	5 mins - 10 mins	Seconds - 2 mins
Query Latency: Detailed Data	Seconds	Sub-second
Query Latency: Pre-Processed Data	Sub-second	Sub-second

What's Next For WeChat

Looking ahead, WeChat's goal is to continually explore and refine their existing data lakehouse architecture to further integrate it across more critical operations including:

- Enabling users to interact with the system using SQL without needing to understand the underlying architecture.
- Unifying data access, querying, and the storage system.
- Unifying SQL interaction standards across the platform.



[Join the StarRocks Slack Community](#) 

Share notes, ask questions, and get feedback from thousands of your peers working at world-class companies.