# From 24h to 15min

*How Robinhood Built a Streaming Lake House to Bring Data Freshness from 24h to <15min*

ORGANIZED BY ◈ databricks

**Balaji Varadarajan**

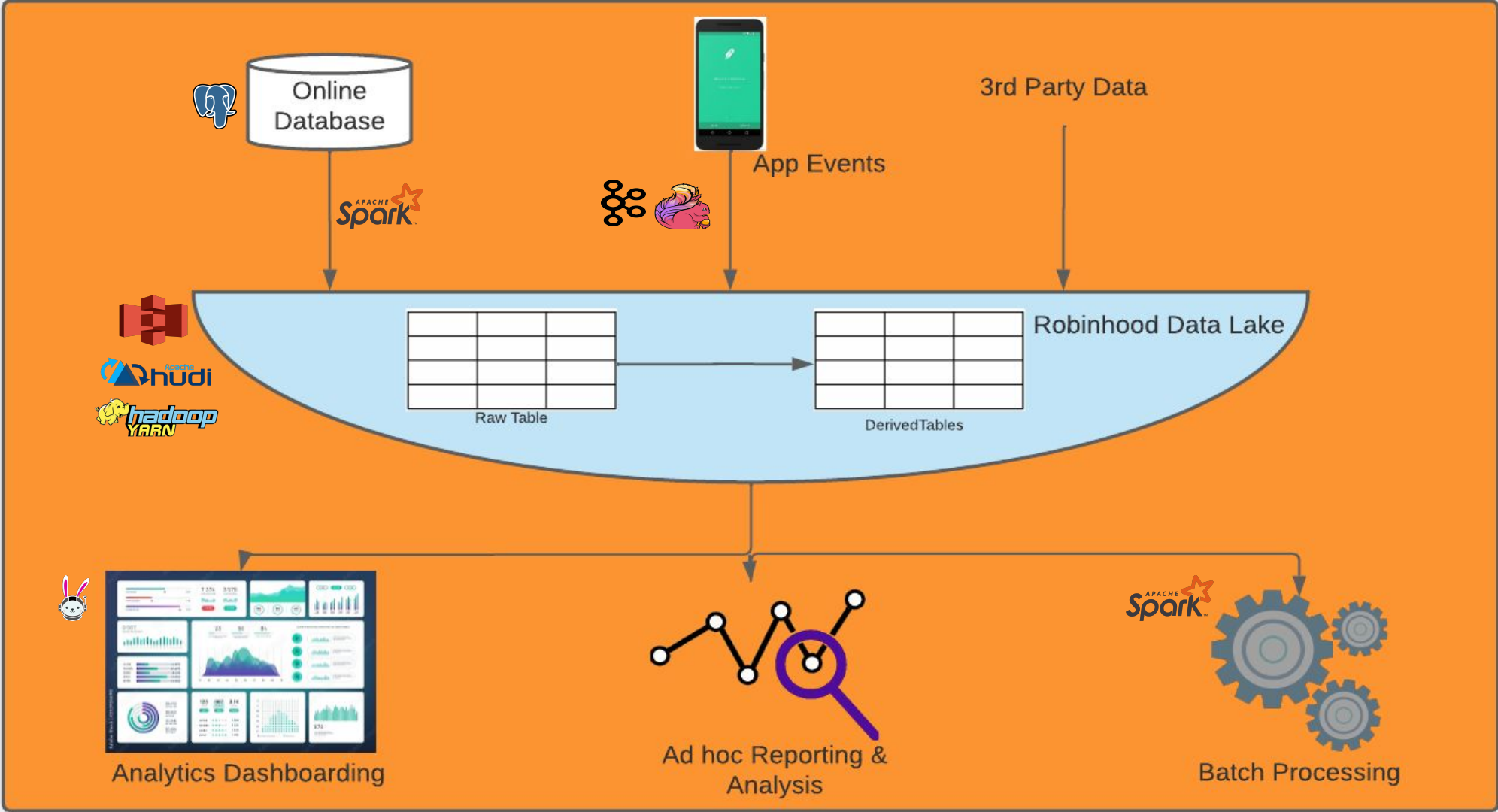Senior Staff Engineer, Robinhood Markets

Apache Hudi PMC

**Vikrant Goel**

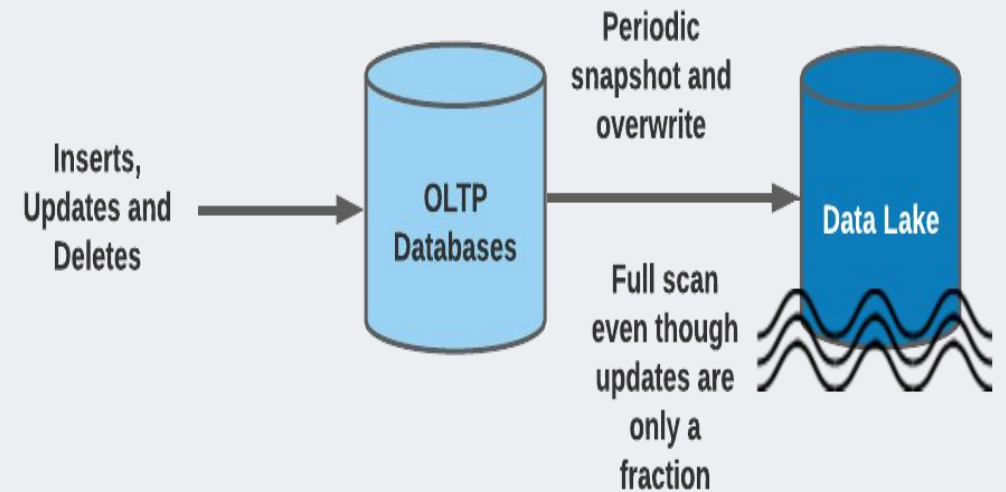Engineering Manager, Robinhood Markets

1

# Agenda

- Legacy Data Lake and Ingestion Framework

- Deep Dive - Change Data Capture (CDC)
  - Design
  - Lessons Learned

- Deep Dive - Data Lakehouse Ingestion
  - Apache Hudi
  - End to End Setup
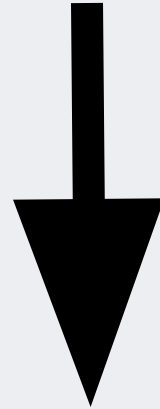
# Data Lake

**DATA+AI**
SUMMIT 2022

# Daily Snapshots

- Daily snapshotting of tables in RDBMS (RDS)

- High Read & Write amplifications

- Dedicated Replicas to isolate snapshot queries

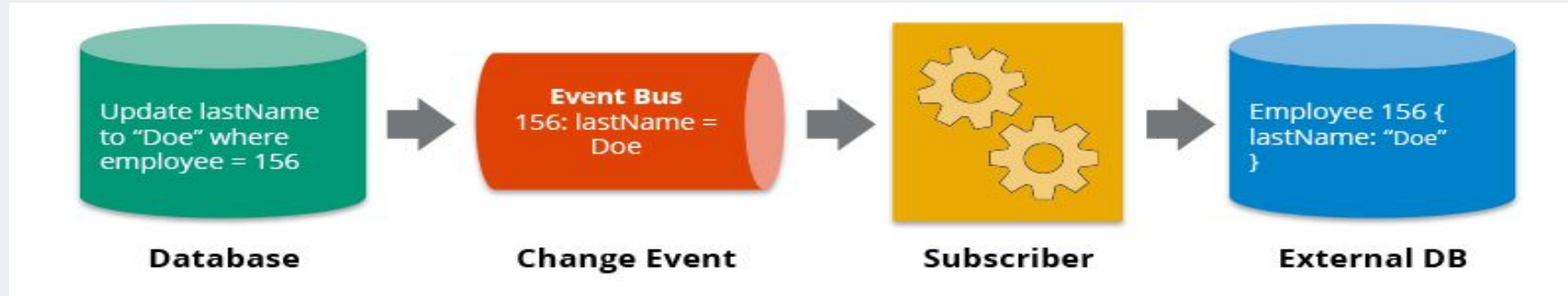- Bottlenecked by Replica I/O

- 24+ hours data latency

# Need Faster Intraday Ingestion Pipeline

↓

# *Unlock Data Lake for business critical applications*
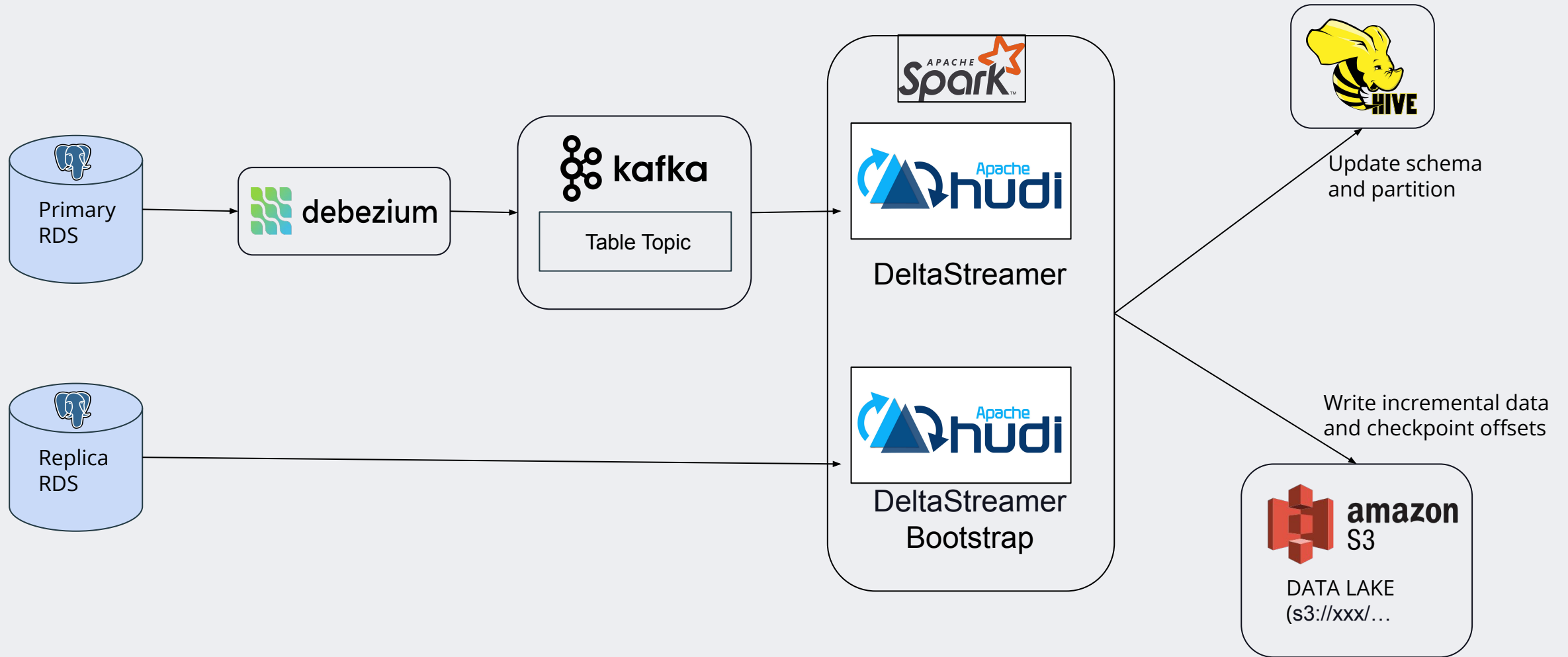
# Change Data Capture



- Each CRUD operation streamed from DB to Subscriber

- Merge changes to lake house

- Efficient & Fast -> *Capture and Apply only deltas*

# Deep Dive
# CDC using Debezium

# Debezium

- Open source & distributed Kafka-Connect Service for change data capture

- Support CDC from diverse RDBMS (Postgres, MySQL, MongoDB, etc.)

- Pluggable Sinks through Kafka

# High Level Architecture



Primary RDS → debezium → kafka (Table Topic) → Apache hudi DeltaStreamer (Spark)

Replica RDS → Apache hudi DeltaStreamer Bootstrap (Spark)

HIVE — Update schema and partition

amazon S3 — DATA LAKE (s3://xxx/... — Write incremental data and checkpoint offsets

DATA+AI
SUMMIT 2022

# Debezium - Zooming In



2. Debezium consumes WALs using Postgres Logical Replication

4. Debezium writes avro serialized updates into table level Kafka topics

Primary DB (RDS) WriteAheadLogs (WALs)

LOG   LOG   LOG

Table_1 Topic

Table_2 Topic

Table_n Topic

3. Debezium updates and validates avro schemas using Kafka Schema Registry

AVRO Schema Registry

1. Enable logical-replication. All updates to the Postgres RDS database are logged into binary files called WriteAheadLogs (WALs)

# Why did we choose Debezium over alternatives?

| | Debezium | AWS Database Migration Service (DMS) |
|---|---|---|
| Operational Overhead | High | Low |
| Cost | Free, with engineering time cost | Relatively expensive, with negligible engineering time cost |
| **Speed** | **High** | **Not enough** |
| Customizations | Yes | No |
| Community Support | Debezium has a very active and helpful Gitter community. | Limited to AWS support. |

# Debezium: Lessons Learned



Postgres Primary:
- Publishes WALs
- Record LogSequenceNumber (LSN) for each consumer

## 1. Postgres Primary Dependency

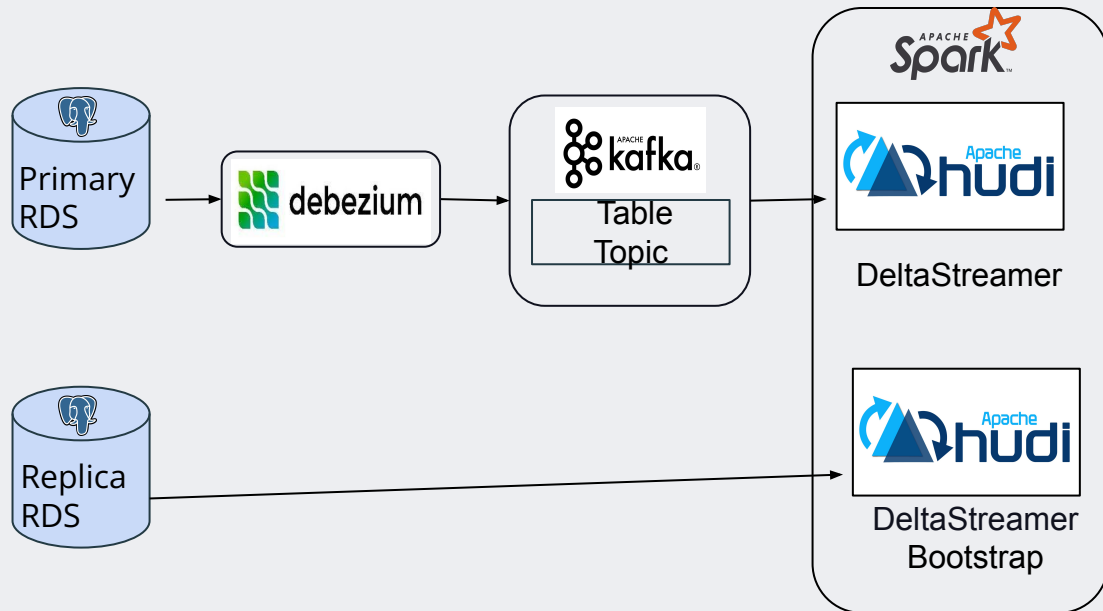ONLY the Postgres Primary publishes WriteAheadLogs (WALs).

### Disk Space:
- If a consumer dies, Postgres will keep accumulating WALs to ensure Data Consistency
- Can eat up all the disk space
- Need proper monitoring and alerting

### CPU:
- Each logical replication consumer uses a small amount of CPU
- Postgres10+ uses pgoutput (built-in)        : Lightweight
  Postgres9 uses wal2Json (3rd party)        : Heavier
- Need upgrades to Postgres10+

All trademarks, logos and brands are the property of their respective owners.

# Debezium: Lessons Learned



## 2. Initial Table Snapshot (Bootstrapping)

**Need for bootstrapping:**
- Each table to replicate requires initial snapshot, on top of which ongoing logical updates are applied

**Problem with Debezium:**
- Debezium processes data row at row level
- Large tables are slow
- Too much pressure on Kafka Infrastructure and Postgres primary
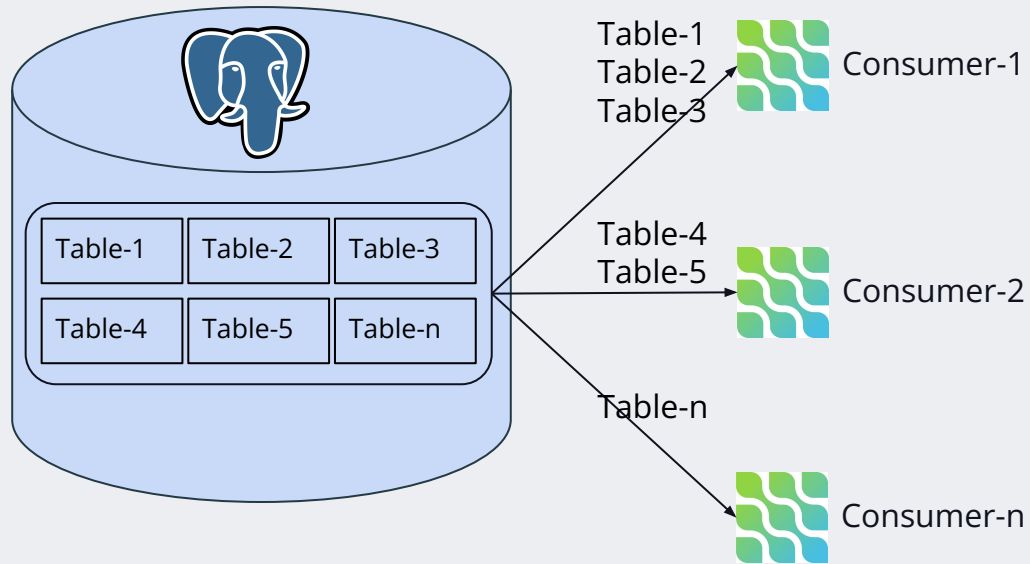
**Solution using Hudi Deltastreamer:**
- Custom bootstrapping framework using partitioned and distributed spark reads
- Can use read-replicas instead of the master

# Debezium: Lessons Learned

## 3. AVRO vs JSON

| | AVRO | JSON | JSON + Schema |
|---|---|---|---|
| Throughput (Benchmarked using db.r5.24xlarge Postgres RDS instance) | Up to **40K mps** | Up to **11K mps.** JSON records are larger than AVRO. | Up to **3K mps.** Schema adds considerable size to JSON records. |
| Data Types | - Supports considerably high number of primitive and complex data types out of the box. - Great for type safety. | Values must be one of these 6 data types: - String - Number - JSON object - Array - Boolean - Null | Same as JSON |
| Schema Registry | Required by clients to deserialize the data. | Optional | Optional |

# Debezium: Lessons Learned



## 4. Multiple logical replication streams for horizontal scaling

- Multiple large tables can overwhelm a single Debezium connector

- Split the tables across multiple Debezium connectors
  **Total throughput = throughput_per_connector * num_connectors**

- Each connector does have small CPU cost

# Debezium: Lessons Learned



## 5. Schema evolution and value of **Freezing Schemas**

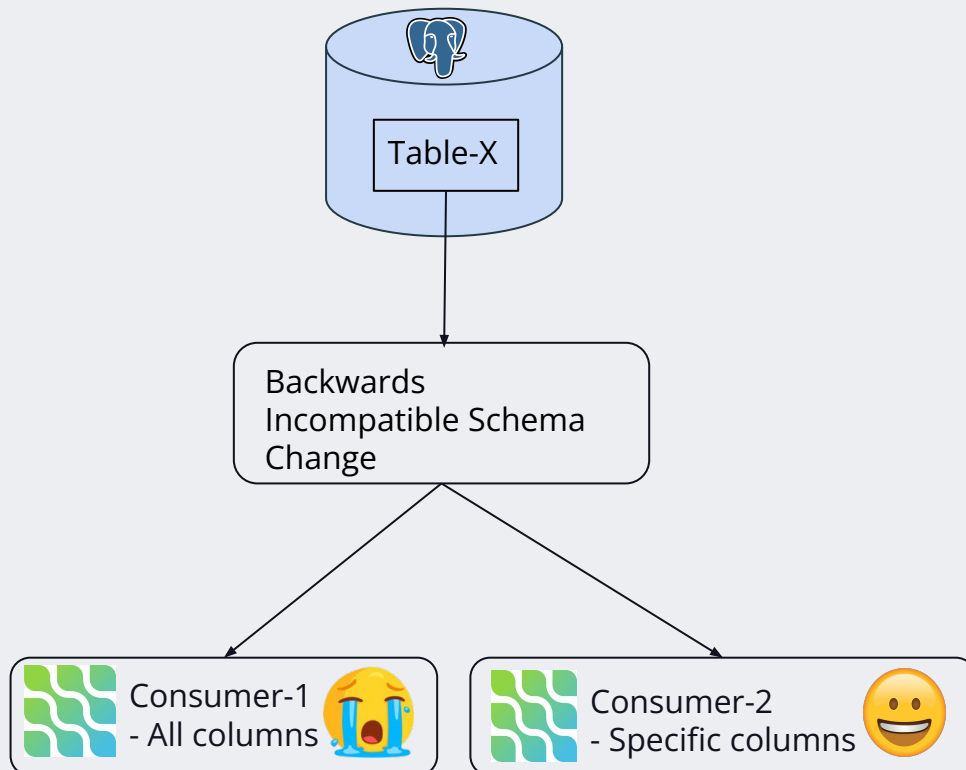**Failed assumption: Schema changes are infrequent and always backwards compatible.**
- Examples:
    1. Adding non-nullable columns (Most Common 99/100)
    2. Deleting columns
    3. Changing data types
    4. Can happen anytime during the day #always_on_call

**How to handle the non backwards compatible changes?**
- Re-bootstrap the table

**Alternatives? Freeze the schema**
- Debezium allows to specify the list of columns per table.
- Pros:
    - #not_always_on_call
    - Batch the changes for management window
- Cons:
    - Schema is temporarily out of sync
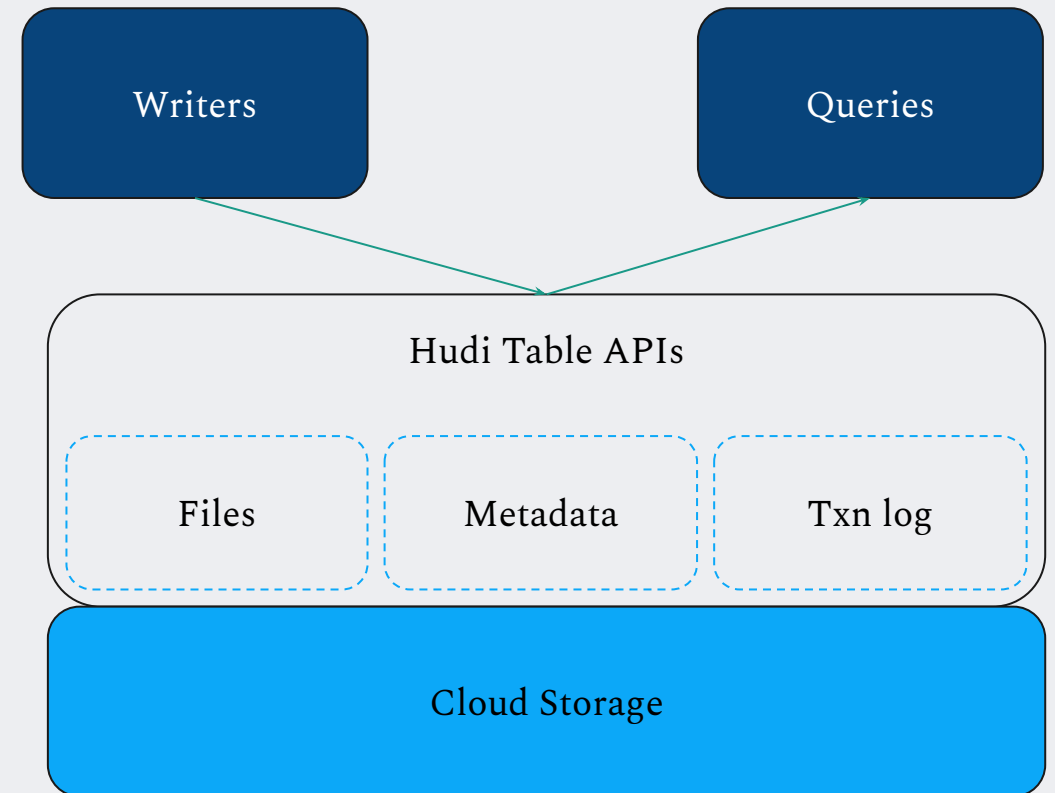
# Deep Dive
# Lakehouse

# Lakehouse - Requirements

- Transaction support

- Scalable Storage and compute

- Openness

- Direct access to files
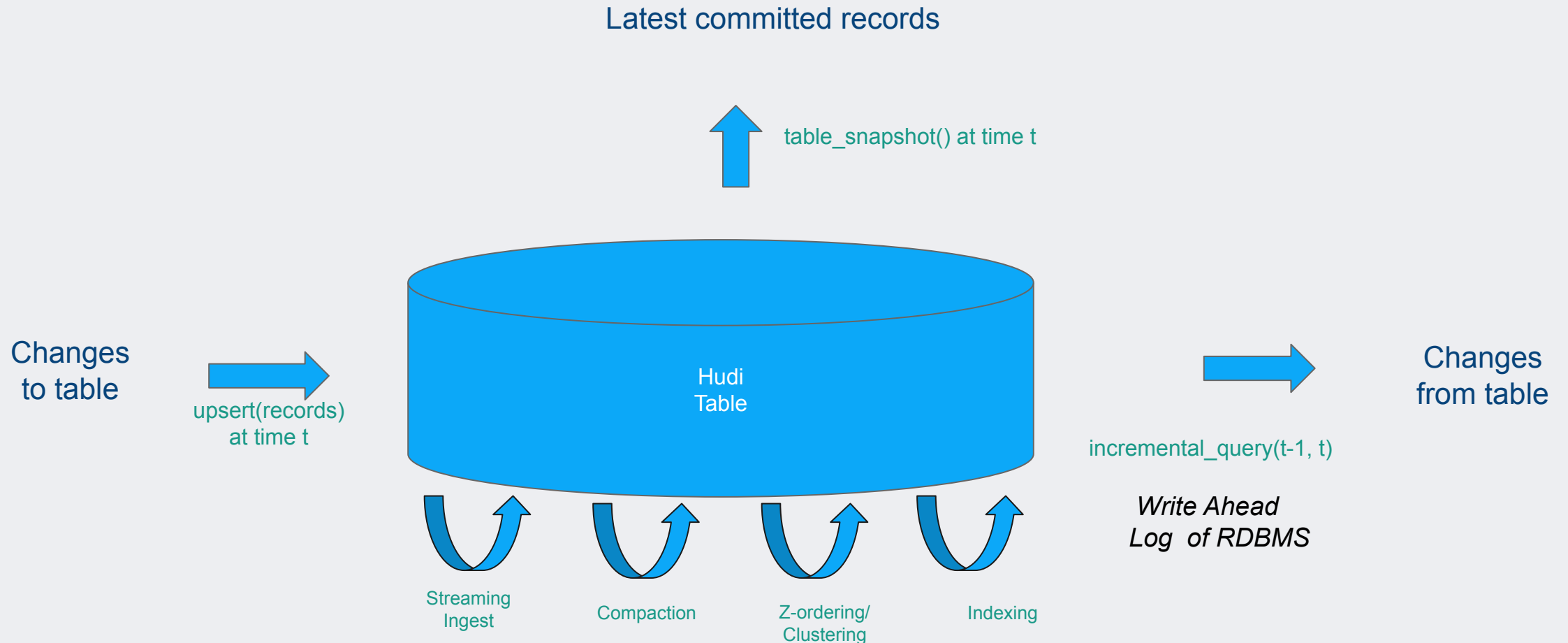
- End-to-end streaming

- Diverse use cases

# Apache Hudi - Introduction

- Transactional Lakehouse pioneered by Hudi
- Serverless, transactional layer over lakes.
- Multi-engine, Decoupled storage from engine/compute
- Upserts, Change capture on lakes
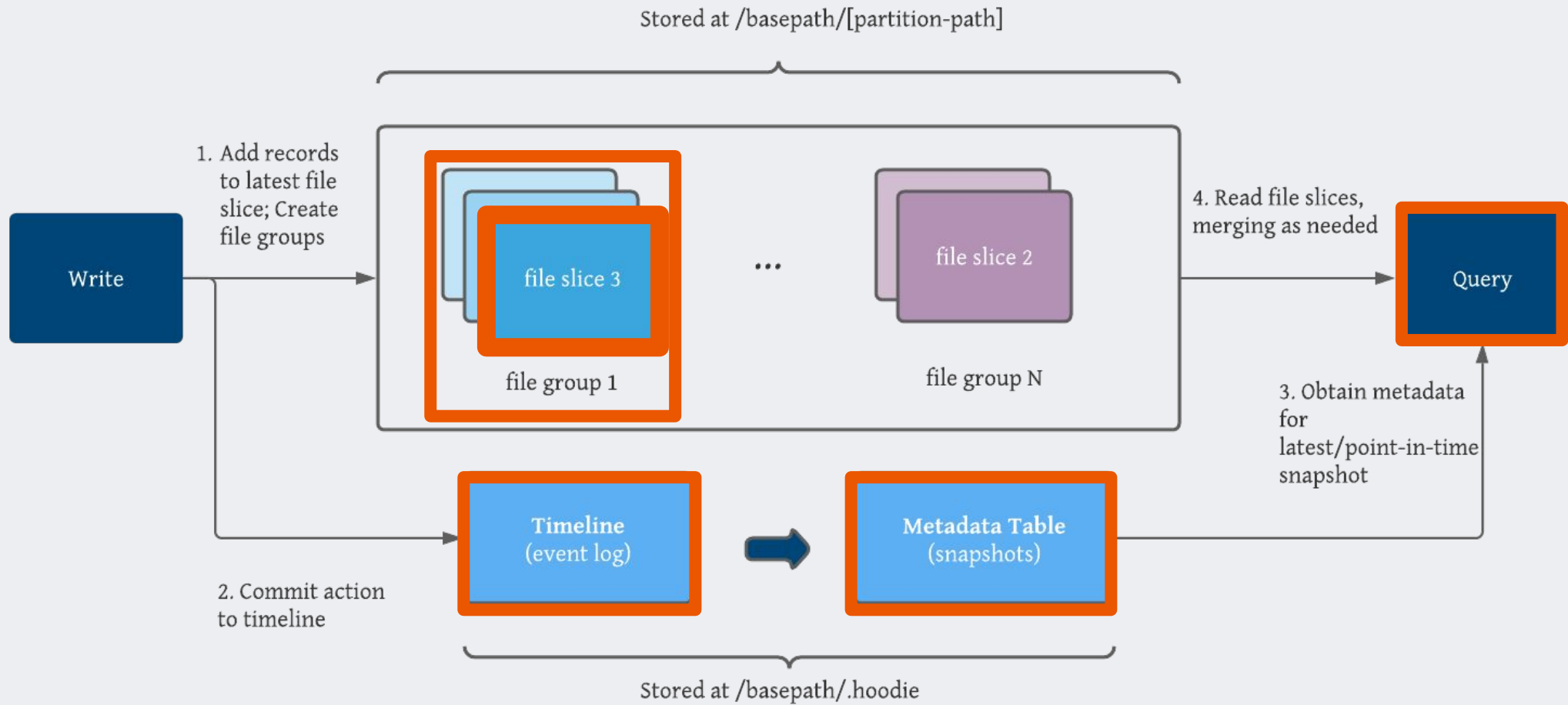- Introduced Copy-On-Write and Merge-on-Read
- Ideas now heavily borrowed outside

Writers

Queries

Hudi Table APIs

Files

Metadata

Txn log

Cloud Storage

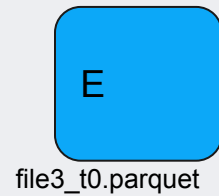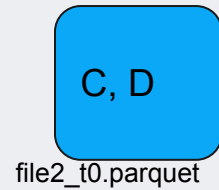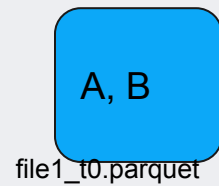https://eng.uber.com/hoodie/ Mar 2017

# Apache Hudi - Upserts & Incrementals
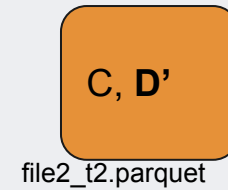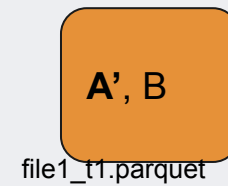
# Apache Hudi - Storage Layout

# Apache Hudi - Copy-On-Write Table

| | Insert: A, B, C, D, E | Update: A => A', D => D' | Update: A' => A", E => E', Insert: F |
|---|---|---|---|
| | commit time=0 | commit time=1 | commit time=2 |

| | A, B | A', B | A", B |
|---|---|---|---|
| | file1_t0.parquet | file1_t1.parquet | file1_t2.parquet |

| | C, D | C, D' | |
|---|---|---|---|
| | file2_t0.parquet | file2_t2.parquet | |

| | E | | E',F |
|---|---|---|---|
| | file3_t0.parquet | | file3_t2.parquet |

| | | | |
|---|---|---|---|
| Snapshot Query | A,B,C,D,E | A',B,C,D',E | A",B,C,D',E',F |
| Incremental Query | A,B,C,D,E | A',D' | A",E',F |

**DATA+AI**
SUMMIT 2022

All trademarks, logos and brands are the property of their respective owners.

# Apache Hudi - Merge-On-Read Table

| Insert: A, B, C, D, E | Update: A => A', D => D' | Update: A'=>A", E=>E',Insert: F | Compaction |
|---|---|---|---|
| commit time=0 | commit time=1 | commit time=2 | commit time=3 |

| A, B | A' | A" | A", B |
|---|---|---|---|
| file1_t0.parquet | .file1_t1.log | .file1_t2.log | file1_t3.parquet |

| C, D | D' | | C, D' |
|---|---|---|---|
| file2_t0.parquet | .file2_t2.log | | file2_t3.parquet |

| E | | E',F | E',F |
|---|---|---|---|
| file3_t0.parquet | | .file3_t2.log | file3_t3.parquet |

| | commit time=0 | commit time=1 | commit time=2 | commit time=3 |
|---|---|---|---|---|
| Snapshot Query | A,B,C,D,E | A',B,C,D',E | A",B,C,D',E',F | A",B,C,D',E',F |
| Incremental Query | A,B,C,D,E | A',D' | A",E',F | A",E',F |
| Read Optimized Query | A,B,C,D,E | A,B,C,D,E | A,B,C,D,E | A",B,C,D',E',F |

# The Community

## Pre-installed on 5 cloud providers

aws · Google Cloud · Alibaba Cloud · Tencent Cloud · HUAWEI

## Diverse PMC/Committers

Uber · Apple · Snowflake · aws · lyft · Robinhood · zendesk · ByteDance · Alibaba · Tencent Cloud · Ant Financial 蚂蚁金服 · DoubleVerify · Onehouse

## Rich community of participants

ebay · Walmart · Udemy · Disney+ · Microsoft · Twitter · Baidu 百度 · Roku · zoom · shopify · amazon · Morgan Stanley · SK telecom · SAMSUNG · zomato · HSBC · Analytics Vidhya · Alipay · GE · H3C 新华三集团 · DATA+AI SUMMIT 2022

| 800B+ Records/Day (from even just 1 customer!) | 2000+ Slack Members | 1M DLs/month (400% YoY) |
|---|---|---|
| 1000+ GH Engagers | 225+ Contributors | 20+ Committers |

Apr '20, May '20, Jun '20, Jul '20, Aug '20, Sep '20, Oct '20, Nov '20, Dec '20, Jan '21, Feb '21, Mar '21, Apr '21, May '21, Jun '21, Jul '21, Aug '21, Sep '21, Oct '21, Nov '21, Dec '21, Jan '22
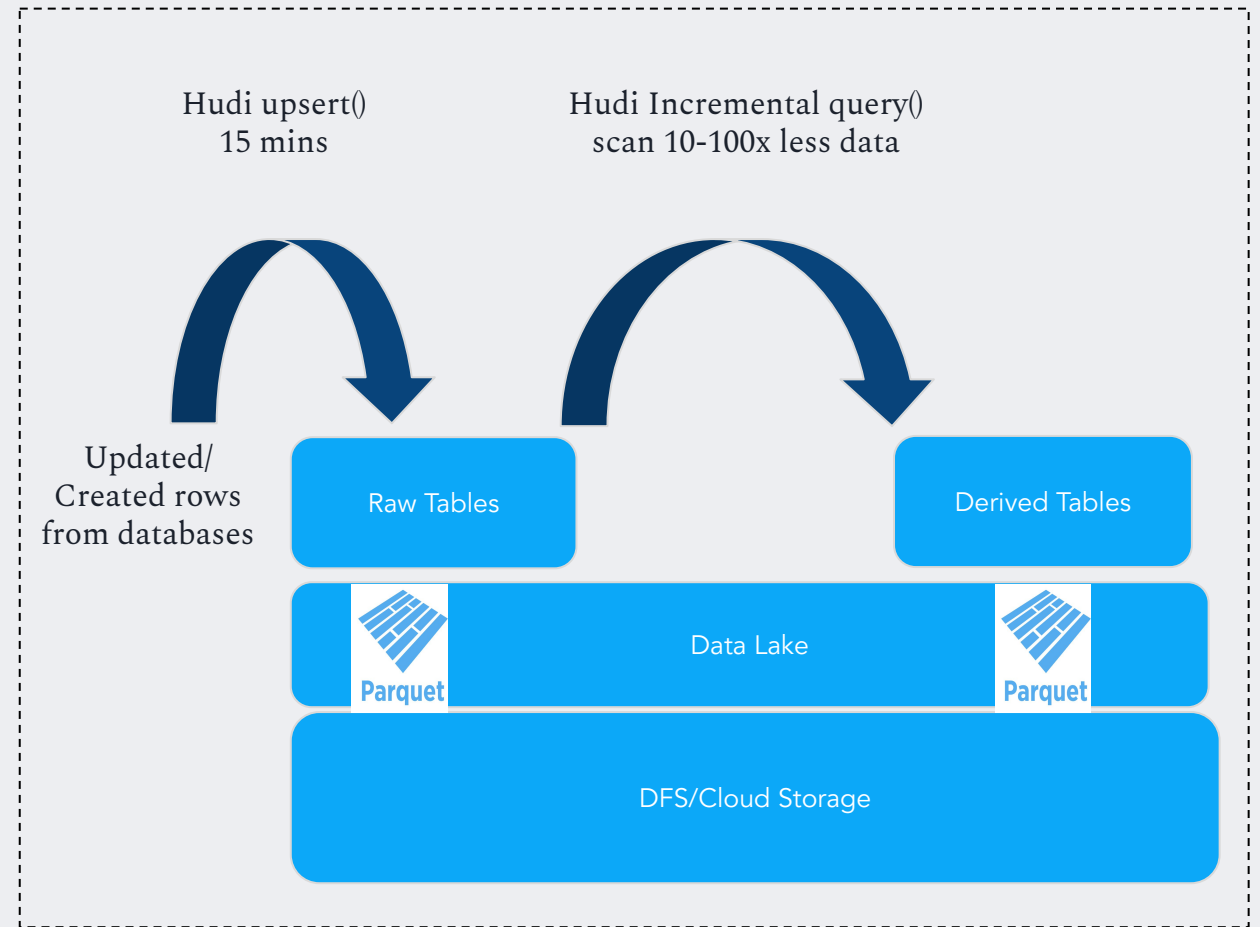
All trademarks, logos and brands are the property of their respective owners.
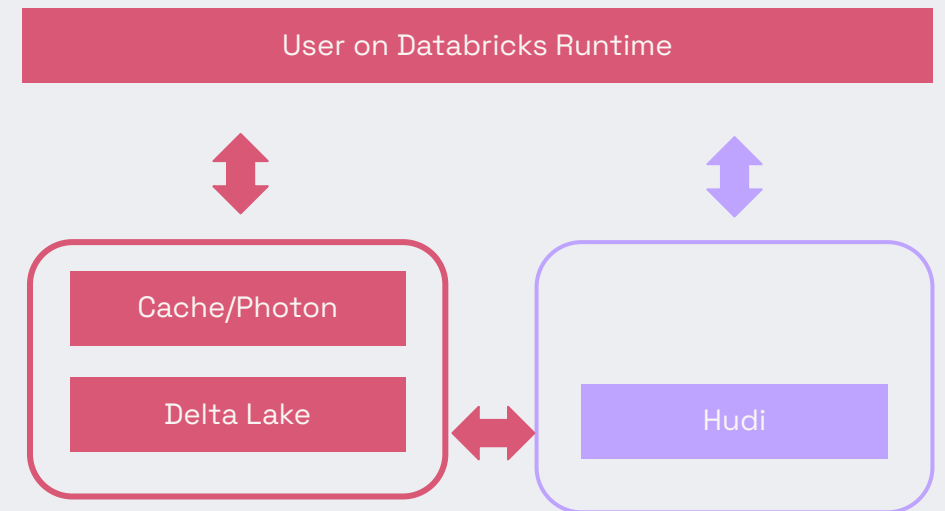
# Apache Hudi - Community

# Apache Hudi - Relevant Features

- Database abstraction for cloud storage/hdfs

- Near real-time ingestion

- Incremental, Efficient ETL downstream

- ACID Guarantees

Hudi upsert()
15 mins

Hudi Incremental query()
scan 10-100x less data

Updated/
Created rows
from databases

Raw Tables

Derived Tables

Parquet

Data Lake

Parquet

DFS/Cloud Storage
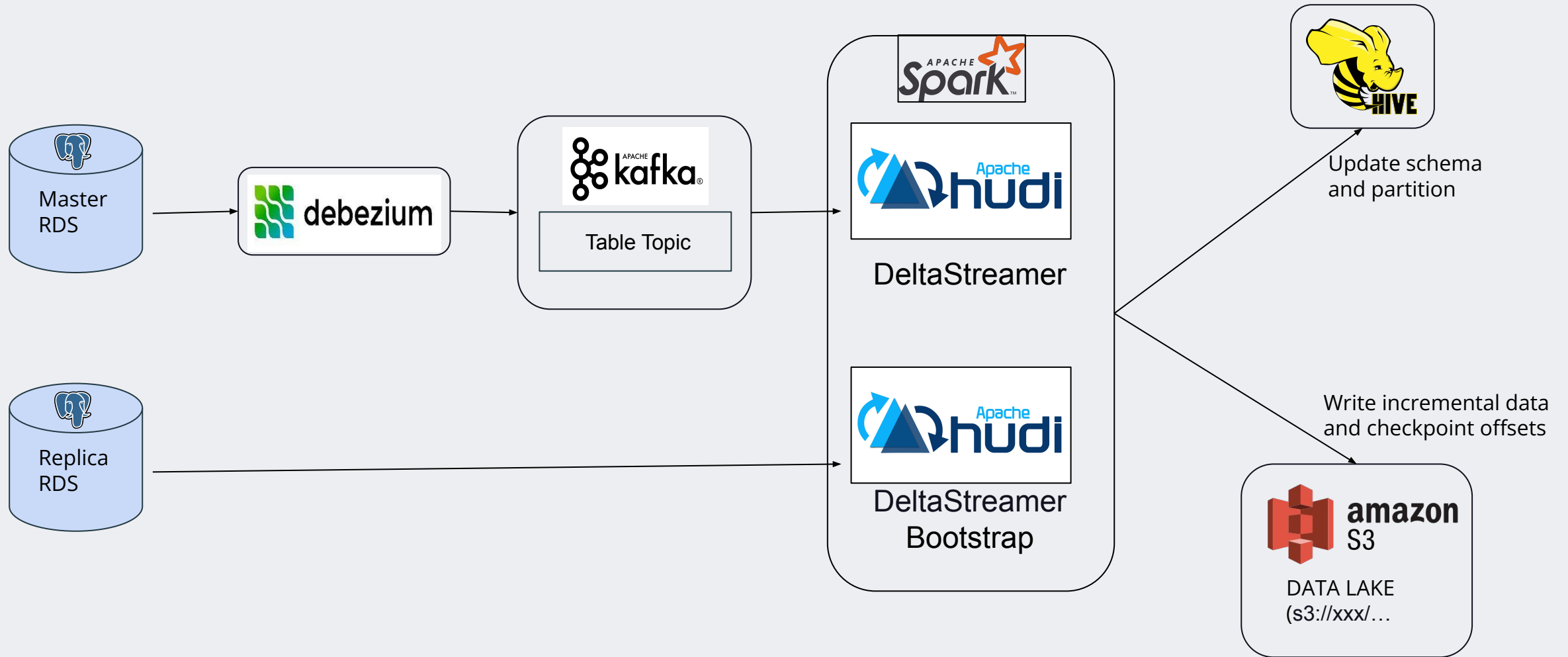
26

**DATA+AI**
SUMMIT 2022

# Apache Hudi + Databricks = Best of Both Worlds!

- We believe
    - Hudi offers the most complete lakehouse storage platform.
    - Databricks offers great Spark experience, but cache/photon don't work natively on Hudi.

- Hudi supports pluggable metadata syncing - BigQuery, Hive, Glue, DataHub, …
- WIP - A bridge between Hudi and Delta Lake.
    - Query Hudi table snapshot as a Delta Lake table
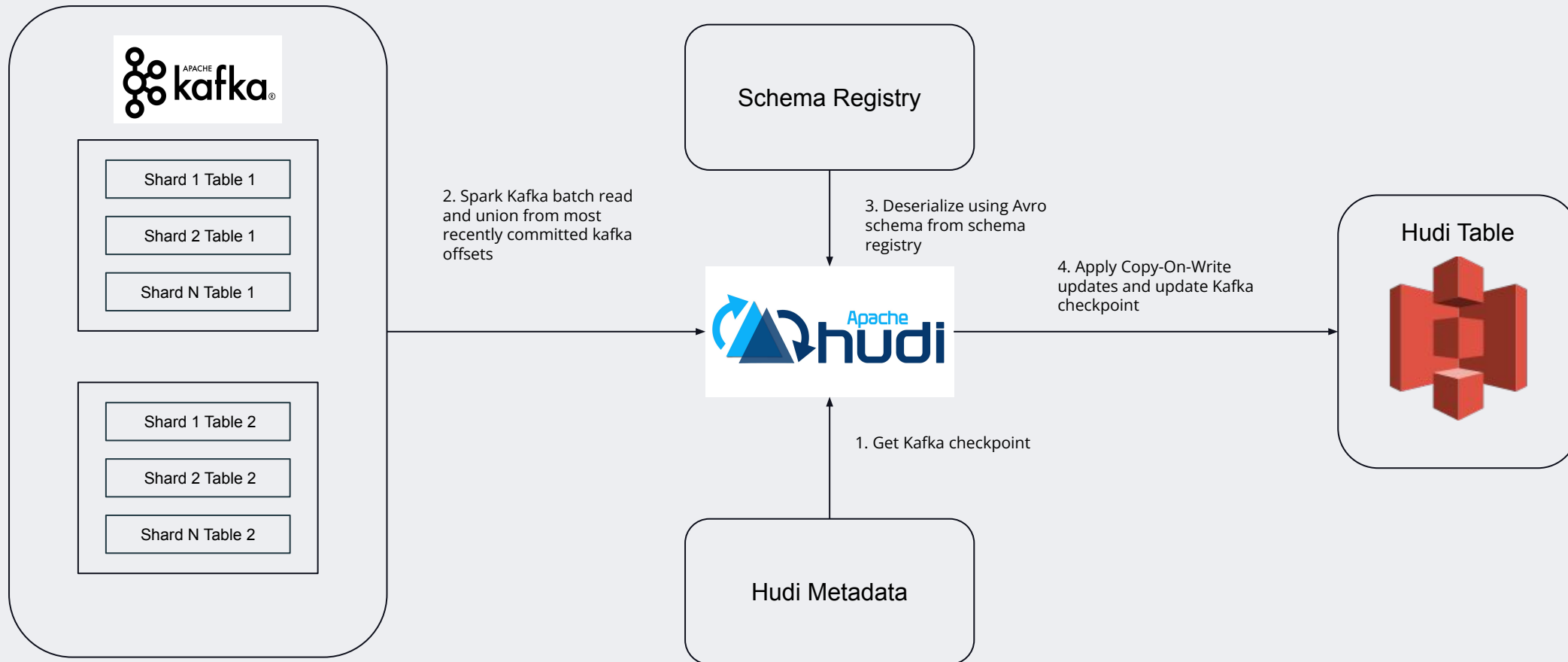
**User on Databricks Runtime**
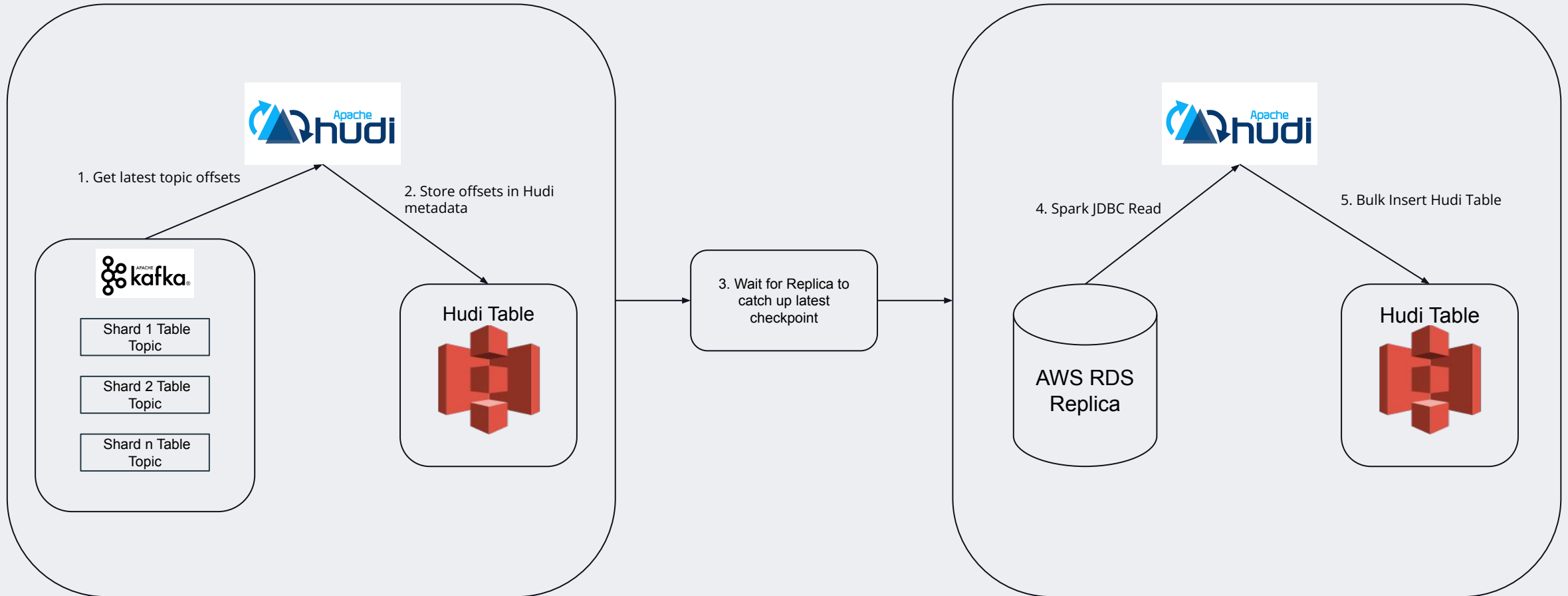
Cache/Photon

Delta Lake

Hudi

# Deep Dive
# Ingestion

# Recap - High Level Architecture



All trademarks, logos and brands are the property of their respective owners.

**DATA+AI**
SUMMIT 2022

# Data Lake Ingestion - CDC Path



APACHE kafka

Shard 1 Table 1

Shard 2 Table 1

Shard N Table 1

Shard 1 Table 2

Shard 2 Table 2

Shard N Table 2

Schema Registry

2. Spark Kafka batch read and union from most recently committed kafka offsets

3. Deserialize using Avro schema from schema registry

Apache hudi

4. Apply Copy-On-Write updates and update Kafka checkpoint

Hudi Table

1. Get Kafka checkpoint

Hudi Metadata

# Data Lake Ingestion - Bootstrap Path



1. Get latest topic offsets

2. Store offsets in Hudi metadata

Shard 1 Table Topic

Shard 2 Table Topic

Shard n Table Topic

Hudi Table

3. Wait for Replica to catch up latest checkpoint

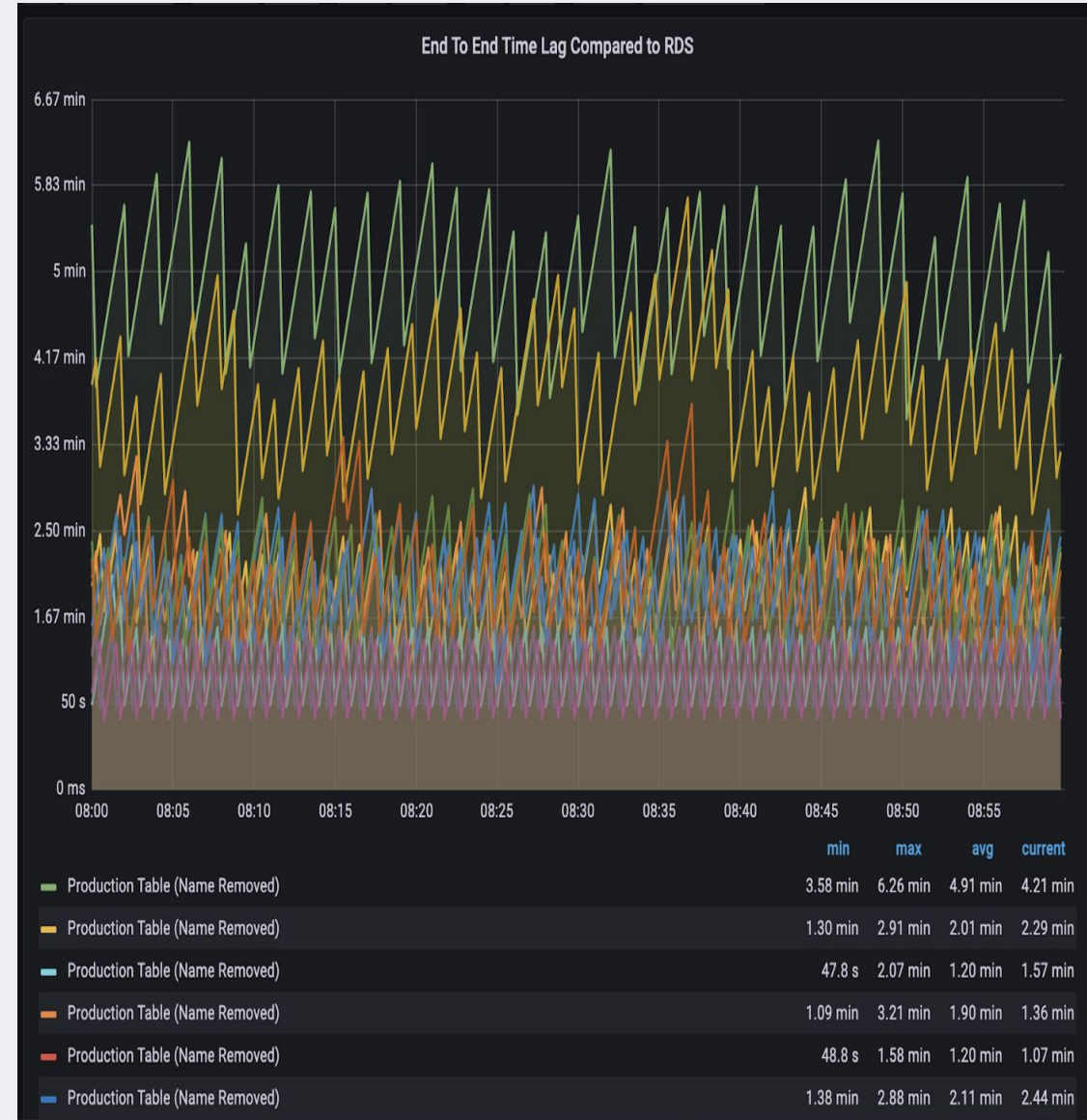4. Spark JDBC Read

5. Bulk Insert Hudi Table

AWS RDS Replica

Hudi Table
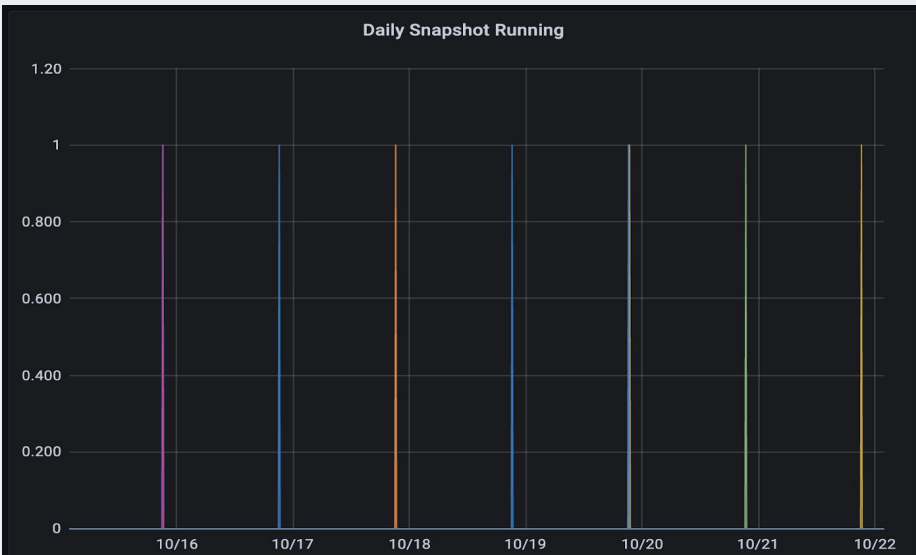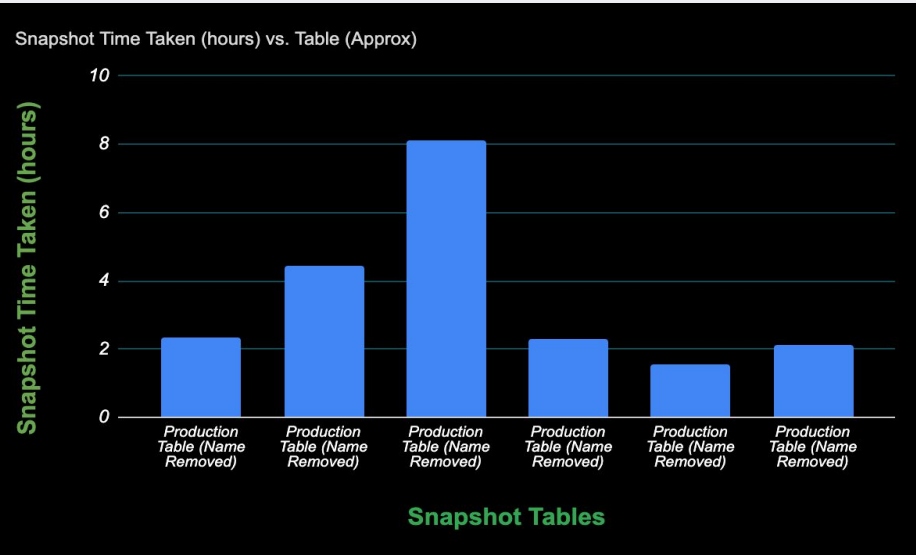
DATA+AI
SUMMIT 2022

# Running Ingestion at Scale

- 4000+ tables

- Automation, Self Healing

- Tiered SLAs - Provisioning and Isolation

- Pre-Commits and Validation for Quality Checks

- Monitoring & Alerting

**DATA+AI**
SUMMIT 2022

# Improved Freshness

# DATA+AI
## SUMMIT 2022

# Thank you!

We'll hangout at back of room for any QA.

**Balaji Varadarajan**

Senior Staff Engineer, Robinhood Markets

Apache Hudi PMC

**Vikrant Goel**

Engineering Manager, Robinhood Markets