# Open Data Foundations across Hudi, Iceberg, and Delta Lake

Data Council

# Speaker Bio

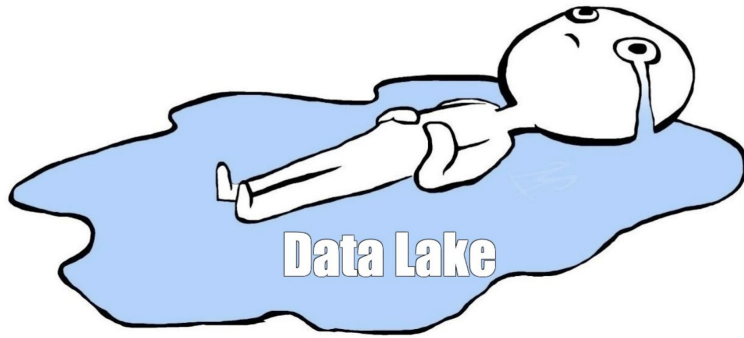## Kyle Weller
Head of Product @ ONEHOUSE

https://www.linkedin.com/in/lakehouse/

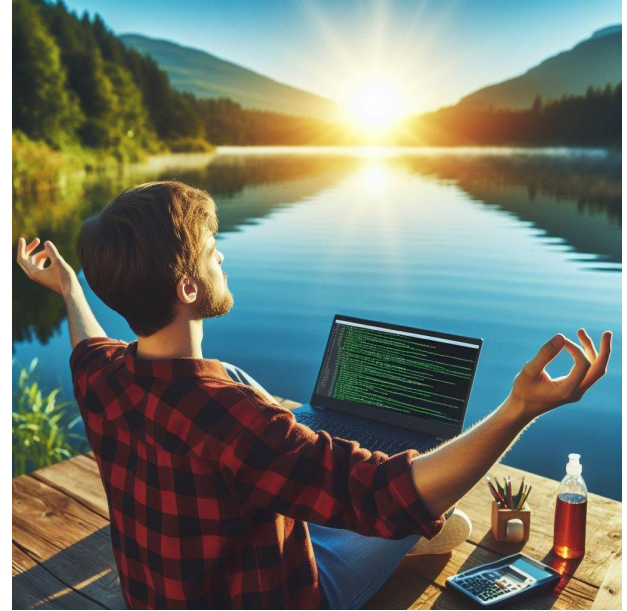**10+ years of building data platforms and data products**

- Currently building Onehouse.ai
- Product lead for Azure Databricks 0 -> 9-fig ARR
- Built Azure Machine Learning services inside SQL Server
- Ran data and growth strategy for Cortana (MSFT AI assistant)
- Worked on PB scale data lake platform for Bing Search
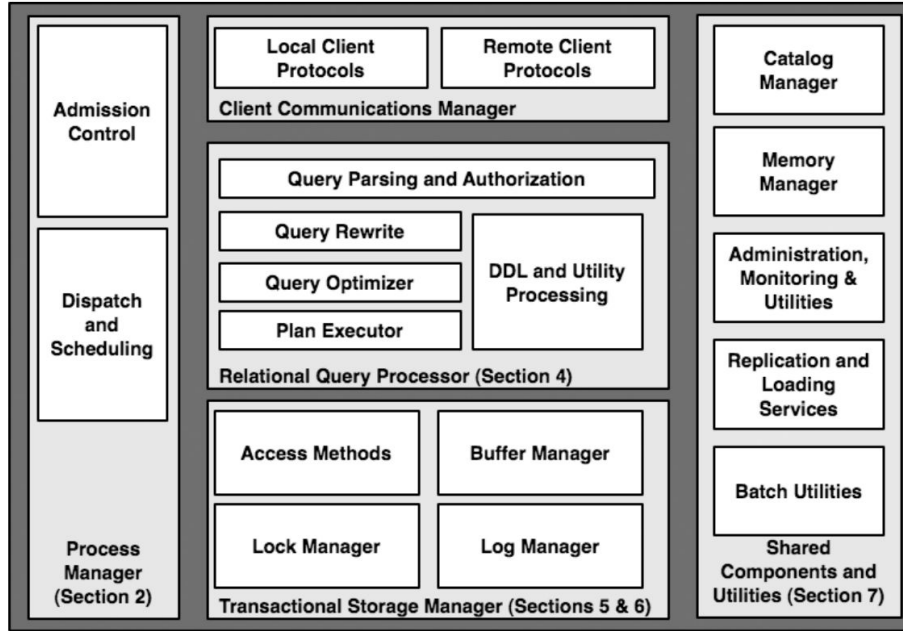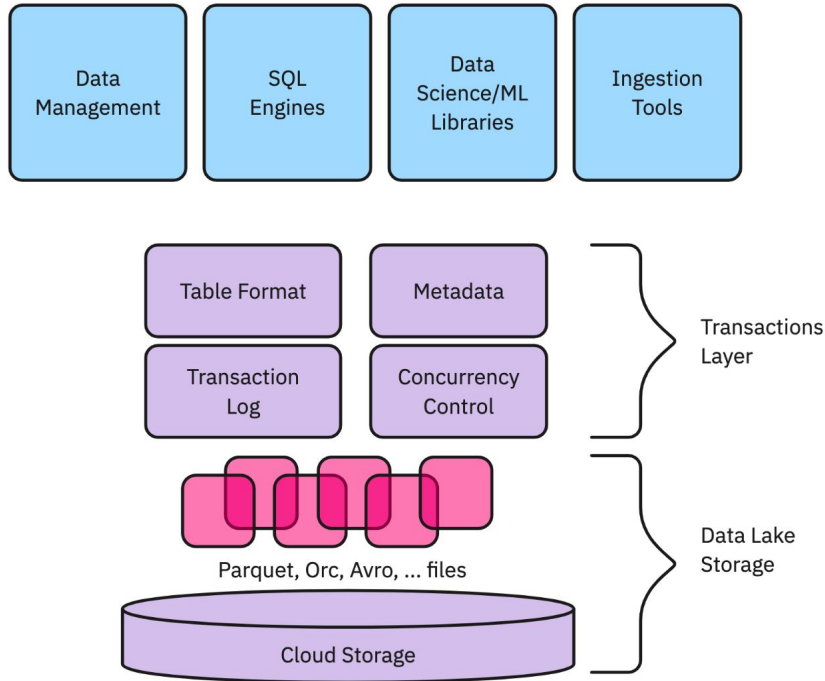- Designed v1 TB scale data lake for MSFT Office

# Data Lakes…



or

# Architecture of a Database System

# S3 Data Lake Storage

# Data Lakehouse - Unbundling of the DBMS

Data Management | SQL Engines | Data Science/ML Libraries | Ingestion Tools

Table Format | Metadata
Transaction Log | Concurrency Control

Transactions Layer

Parquet, Orc, Avro, … files

Cloud Storage

Data Lake Storage

Apache Hudi | DELTA LAKE | ICEBERG

**A lakehouse has the following key features:**

- **Transaction support:** In an enterprise lakehouse many data pipelines will often be reading and writing data concurrently. Support for ACID transactions ensures consistency as multiple parties concurrently read or write data, typically using SQL.

- **Schema enforcement and governance:** The Lakehouse should have a way to support schema enforcement and evolution, supporting DW schema architectures such as star/snowflake-schemas. The system should be able to reason about data integrity, and it should have robust governance and auditing mechanisms.

- **BI support:** Lakehouses enable using BI tools directly on the source data. This reduces staleness and improves recency, reduces latency, and lowers the cost of having to operationalize two copies of the data in both a data lake and a warehouse.

- **Storage is decoupled from compute:** In practice this means storage and compute use separate clusters, thus these systems are able to scale to many more concurrent users and larger data sizes. Some modern data warehouses also have this property.

- **Openness:** The storage formats they use are open and standardized, such as Parquet, and they provide an API so a variety of tools and engines, including machine learning and Python/R libraries, can efficiently access the data **directly**.

- **Support for diverse data types ranging from unstructured to structured data**: The lakehouse can be used to store, refine, analyze, and access data types needed for many new data applications, including images, video, audio, semi-structured data, and text.

- **Support for diverse workloads:** including data science, machine learning, and SQL and analytics. Multiple tools might be needed to support all these workloads but they all rely on the same data repository.

- **End-to-end streaming:** Real-time reports are the norm in many enterprises. Support for streaming eliminates the need for separate systems dedicated to serving real-time data applications.

# Origin Stories

**2017** — open sourced **Uber**

```
1   + # Hudi
2   + Hudi (pronounced Hoodie) stands for `Hadoop Upserts anD Incrementals`. Hudi manages storage of large
      analytical datasets on [HDFS](http://hadoop.apache.org/docs/stable/hadoop-project-dist/hadoop-
      hdfs/HdfsDesign.html) and serve them out via two types of tables
3
4     * **Read Optimized Table** - Provides excellent query performance via purely columnar storage (e.g.
      [Parquet](https://parquet.apache.org/))
5     * **Near-Real time Table (WIP)** - Provides queries on real-time data, using a combination of columnar &
      row based storage (e.g Parquet + [Avro](http://avro.apache.org/docs/current/mr.html))
```

**2018** — open sourced **NETFLIX**

```
1   + ## Iceberg
2   +
3   + Iceberg is a new table format for storing large, slow-moving tabular
      data. It is designed to improve on the de-facto standard table layout
      built into Hive, Presto, and Spark.
4   +
```

**2019** — open sourced

```
    - Delta Lake Core is .... (copy text from delta docs)
3   + Delta Lake is a next-generation engine built on top of Apache Spark. Delta Lake
      provides ACID transactions, optimized layouts and indexes, and execution engine
      improvements for building data pipelines to support big data use cases: batch
      and streaming ingests, fast interactive queries, and machine learning.
      Specifically, Delta offers:
```

# 🌶️ My Hot Take - They are divergent!



- Technical vision and goals are divergent

- The community needs are specialized

- All three projects are on fast growth trajectories

- New table formats are gaining traction: Apache Paimon, YOHB?



Ali Ghodsi 🔗 · 1st
CEO & Co-Founder at Databricks, Adjunct Professor at UC Berkeley
4mo · 🌐

Actually think **Vinoth Chandar** put the truth out there really well:
"There's already 3 major projects - Delta Lake, Hudi & Iceberg with thousands of users for each project. From a data OSS community perspective, we're way past having a standard open table format, what's important is to now make progress and move the industry forward in interoperability.



stackoverflow    Products    🔍 Search...

🏠 Home          Thrift, Avro, Protocolbuffers - Are they all dead?
🅀 Questions      Asked 7 years, 3 months ago    Modified 5 years, 11 months ago    Viewed 40k times

# Technical Fundamentals

- Metadata abstractions on files in cloud object storage
- Tables with SQL semantics and schema evolution
- ACID transactions
- Updates and deletes (merge/upsert)
- Data layout optimizations for performance tuning



(a) Delta

(b) Iceberg

(c) Hudi

# - How it looks on cloud storage

- Fundamentals of table formats Hudi, Delta, Iceberg are not that different
- Each has a special metadata layer on top of parquet files



```
s3_bucket/my_table/
 |- .hoodie/
 |     |- hoodie.properties
 |     |- metadata/
 |- file_1.parquet
 |- file_2.parquet
 |- file_N.parquet
```

```
s3_bucket/my_table/
 |- _delta_log/
 |     |- 000000.json
 |- file_1.parquet
 |- file_2.parquet
 |- file_N.parquet
```

```
s3_bucket/my_table/
 |- metadata/
 |     |- v1.metadata.json
 |     |- snap-9fa1-2-16c3.avro
 |     |- 0d9a-98fa-77.avro
 |- file_1.parquet
 |- file_2.parquet
 |- file_N.parquet
```

Amazon S3    Google Cloud Storage    Azure Data Lake Storage    HDFS

# - Which Format Should I Choose?

**Choose ![Apache Hudi]  if:**

1. Mutable data - GDPR Deletes, Updates
2. CDC workloads
3. Low latency requirements
4. Large ETL pipelines - perf/cost w/ incremental ETL

**Choose ![Delta Lake]  if:**

1. Best Databricks experience
2. Needs fastest premium Spark with Photon
3. Wants an "easy-to-get-started" table format

**Choose ![Iceberg]  if:**

1. Trino or Athena writes
2. Snowflake writes
3. Not sensitive to performance
4. Partition evolution

# - Which Format Should I Choose?

**Choose** Apache **hudi** **if:**
1. Mutable data - GDPR Deletes, Updates
2. CDC workloads
3. Low latency requirements
4. Large ETL pipelines - perf/cost w/ incremental ETL

**Choose** △ **DELTA LAKE** **if:**
1.
2. Needs fastest perf on Spark with Photon
3. Wants an "easy-to-get started" table format

**Choose** **ICEBERG** **if:**
1. Trino or Athena writes
2. Snowflake writes
3. Insensitive to performance
4. Partition evolution

**What if you could work across all 3?**

# - Example benefits of mix-and-match

## Writing

**Choose** **hudi** **writing w/ EMR** (Spark)

1. Fastest writes for mutable workloads
2. Most flexible tuning parameters for ingestion

**Choose** **DELTA LAKE** **writing w/ Fabric:**

1. Easy-to-get-started out of the box
2. Makes data available to the entire Azure portfolio

**Choose** **ICEBERG** **writing w/ BigQuery**

1. Only table format supported for writes
2. Partition evolution

## Reading

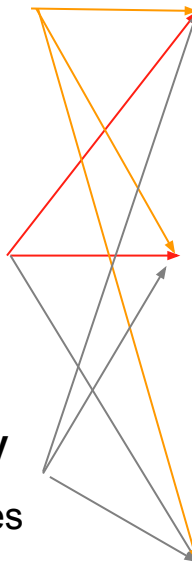**Choose** **DELTA LAKE** **reading w/ Databricks**

1. Get fastest queries with Photon acceleration
2. Great experience for Data Science

**Choose** **ICEBERG** **reading w/ Snowflake**

1. Only supported table format in Snowflake
2. Decouple data storage using external tables

**Choose** **hudi** **reading w/ DataProc** (Spark)

1. Fast record level indexes for point queries
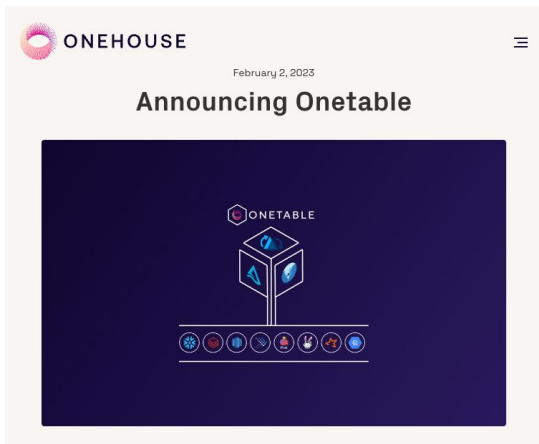2. Powerful secondary indexing capabilities for Spark

# Introducing:


Apache XTable™ (Incubating)

⭐Celebrate by adding a little star ⭐
https://github.com/apache/incubator-xtable



>600 GH Stars ⭐

>80 Forks

# Apache XTable™ - Timeline

**Onehouse announces OneTable**

**OSS Co-Launch with Microsoft, Google, Onehouse**

**Donation to ASF and incubation as Apache XTable**

| Feb 2023 | Nov 2023 | Mar 2024 |

# Apache XTable™ - How it Works

1: Choose your "source" format
2: Choose your "target" format(s)
3: XTable translates the metadata layers

Read your table as any of the formats

| Hudi | Delta | Iceberg |
| --- | --- | --- |

yaml

```
sourceFormat: HUDI
targetFormats:
  - DELTA
  - ICEBERG
datasets:
  -
    tableBasePath: s3://path/to/hudi-dataset/people  # replace this with gs://path/to/hudi-dataset/people if yo
    tableName: people
    partitionSpec: city:VALUE
```
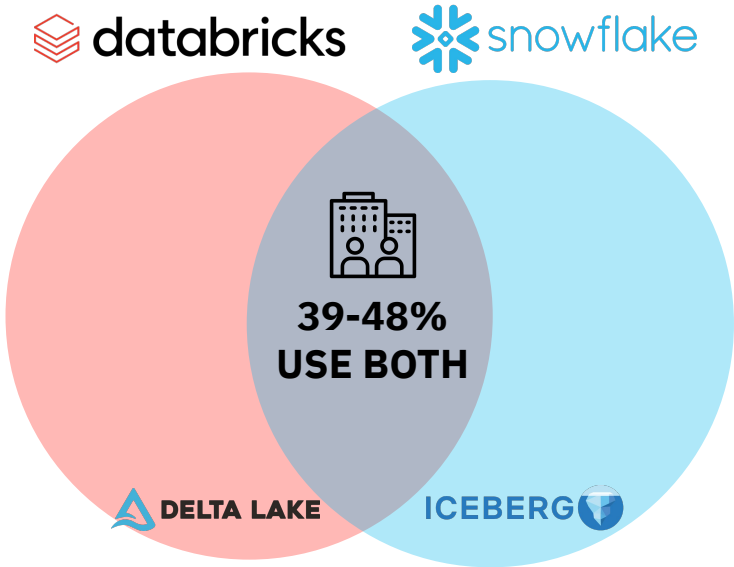
```
s3_bucket/my_table/
 |- .hoodie/
 |    |- hoodie.properties
 |    |- metadata/
 |- file_1.parquet
 |- file_2.parquet
 |- file_N.parquet
```
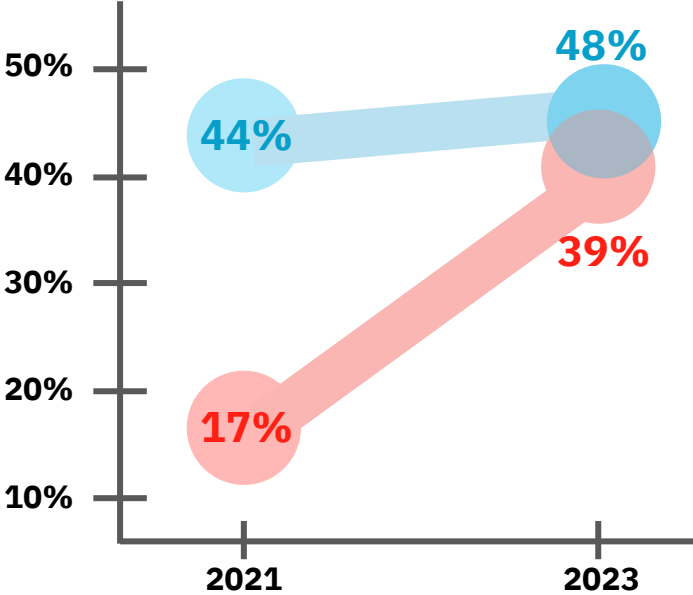
```
s3_bucket/my_table/
 |- .hoodie/
 |    |- hoodie.properties
 |    |- metadata/
 |- _delta_log/
 |    |- 000000.json
 |- metadata/
 |    |- v1.metadata.json
 |    |- snap-9fa1-2-16c3.avro
 |    |- 0d9a-98fa-77.avro
 |- file_1.parquet
 |- file_2.parquet
 |- file_N.parquet
```

```
# any of these work on the same table
spark.read.format("hudi")
spark.read.format("delta")
spark.read.format("iceberg")
```
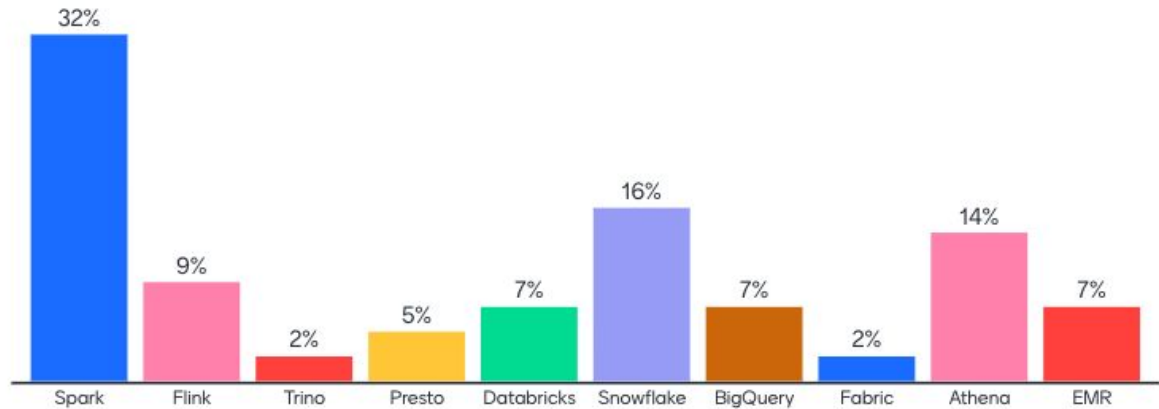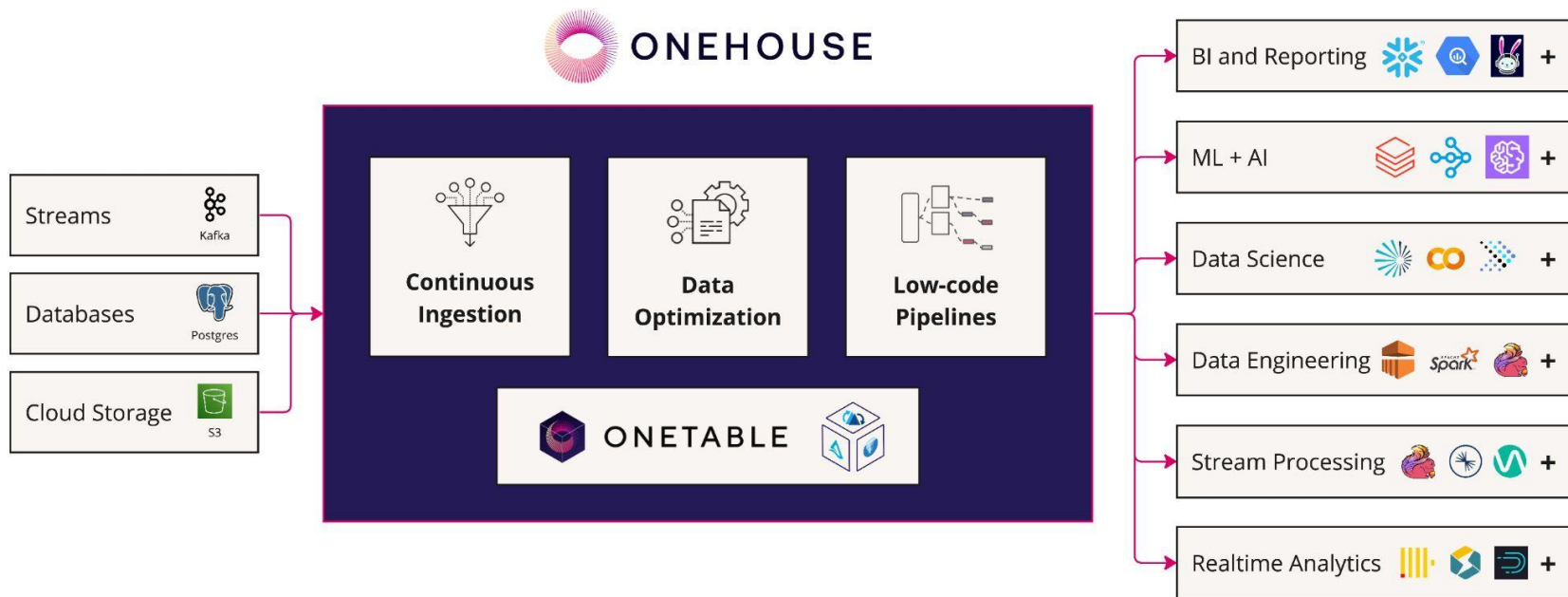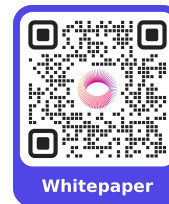
# A tale of two...



databricks     ❄ snowflake



**39-48%**
**USE BOTH**

DELTA LAKE     ICEBERG

## Overlap Growth 2021-2023



50%
44%          48%
40%
39%
30%
20%
17%
10%

2021          2023

Demo Time!

# Apache XTable™ - Vision for the Future

**Goals**
- Seamless and efficient interoperability
- Eliminate data silos
- Project sustainability and evolution

**Features**
- Real-time and transparent replication in any direction
- Accurate and lossless model
- Extensibility and flexibility

**Community**
- Neutral and inclusive: Vendors, Cloud providers, Users
- Graduate ASF Incubation

# Apache XTable™ - Join the community

## Initial Committers

- Tim Brown : *Onehouse*
- Vamshi Gudavarthi : *Onehouse*
- Ashvin Agrawal : *Microsoft*
- Jesus Camacho Rodriguez : *Microsoft*
- Anoop Johnson : *Google*
- Stamatis Zampetakis : *Cloudera*
- Hitesh Shah : *Adobe*
- Jean-Baptiste Onofré : *Dremio*
- Baljinder Singh : *Walmart*
- Vinish Reddy: *Onehouse*
- Vinoth Chandar: *Onehouse*

# Apache XTable™ - Roadmap

| Current Status | Roadmap (6-12 months) | Roadmap (long term) |
|---|---|---|

**Current Status**

- ➤ Supported formats: **Apache Hudi, Apache Iceberg, and Delta Lake**

- ➤ Tested with: Apache Spark, Trino, Microsoft Fabric, Databricks, BigQuery, Snowflake, Redshift, and more

- ➤ Features: on-demand incremental conversion, copy-on-write, catalog integration, change-history

**Roadmap (6-12 months)**

- ➤ **Merge-on-Read** (delete vectors)

- ➤ Apache Paimon (incubating)

- ➤ **Performance**, **efficiency**, and resiliency

- ➤ Deployment: as-a-service and in-memory

- ➤ Native engine integration

**Roadmap (long term)**

- ➤ Multi-writer (duplex)

- ➤ Synchronized commit timestamp

- ➤ Feature parity (superset)

- ➤ New technology stack

- ➤ Support new formats & versions
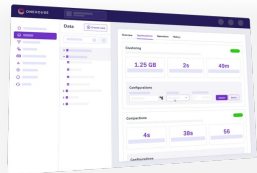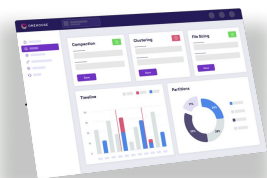
# Apache XTable™ - Let's Build Together

**Github:** https://github.com/apache/incubator-xtable

**Docs :** https://xtable.apache.org/docs/how-to

**Twitter :** https://twitter.com/apachextable

**LinkedIn :** https://www.linkedin.com/company/apache-xtable/

**Mailing List :** dev-subscribe@xtable.apache.org
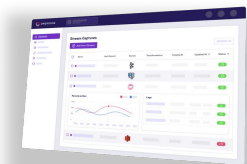
ONEHOUSE

# ONEHOUSE - 3 Ways to Engage with Onehouse

**Lakehouse Monitoring - ($0 Free)**
- No install, no permissions
- Monitoring and tuning insights for your Hudi, Iceberg, and Delta Lake pipelines

**Lakehouse Table Services**
- Keep your existing Hudi, Iceberg, Delta pipelines
- Onehouse will automate advanced Table Optimizations for 10x faster analytics

**Lakehouse Streaming Ingest/ELT**
- 10x Faster/Cheaper vs existing OSS Hudi, Iceberg, Delta pipelines
- Fully managed auto-scaling infrastructure w/ serverless experience in your VPC
- Simple UI + APIs for programmatic and templated CI/CD devops integration

Contact me: kyle@onehouse.ai

Thank You!