

APACHE XTABLE: INTEROPERABILITY BETWEEN OPEN TABLE FORMATS



Dipankar Mazumdar & Kyle Weller, [Onehouse.ai](https://onehouse.ai)
13th June 2024

SPEAKER BIOS

Dipankar Mazumdar



- ❑ Staff Developer Advocate @ Onehouse.ai
- ❑ Contributor @ Apache Hudi, XTable
- ❑ Prev: Data Architecture, Visualization, ML



in/dipankarmazumdar/



@dipankartnt

Kyle Weller



- ❑ Head of Product @ Onehouse.ai
- ❑ Prev: Product Manager - Azure Databricks, Azure ML, Bing Search



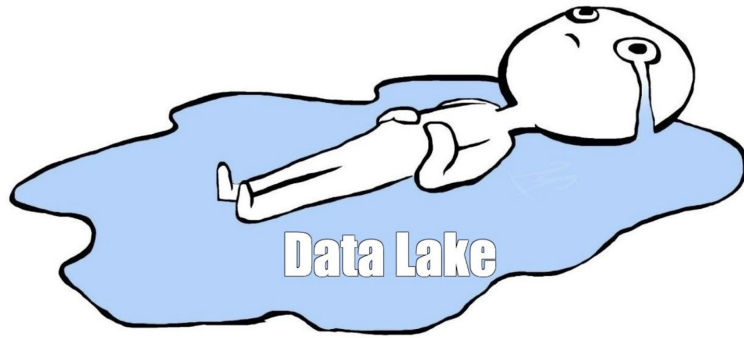
in/lakehouse/



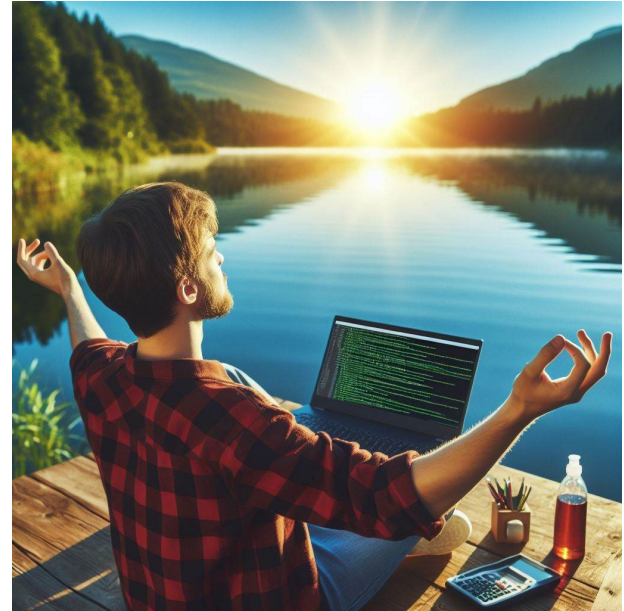
@KyleJWeller



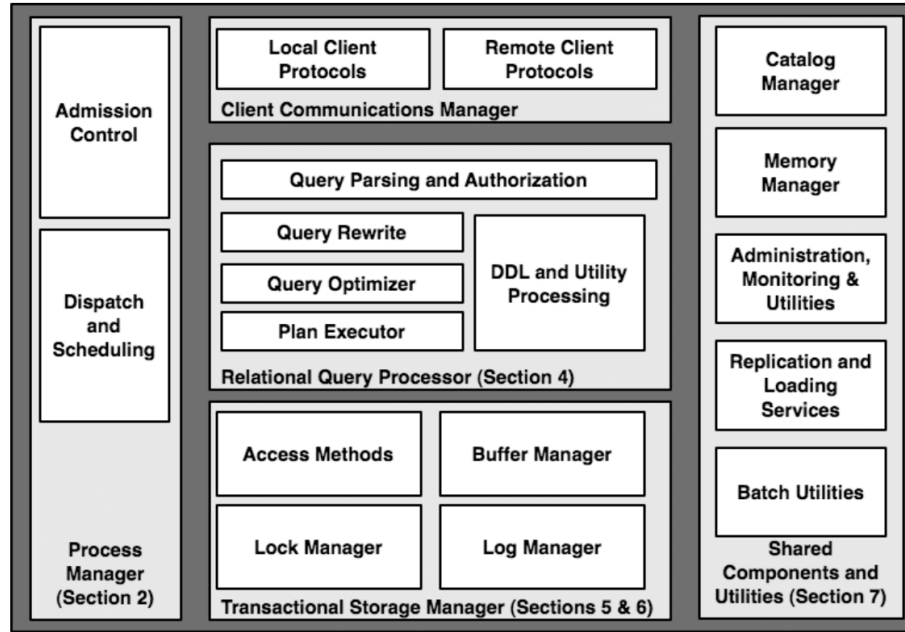
Data Lakes...



or



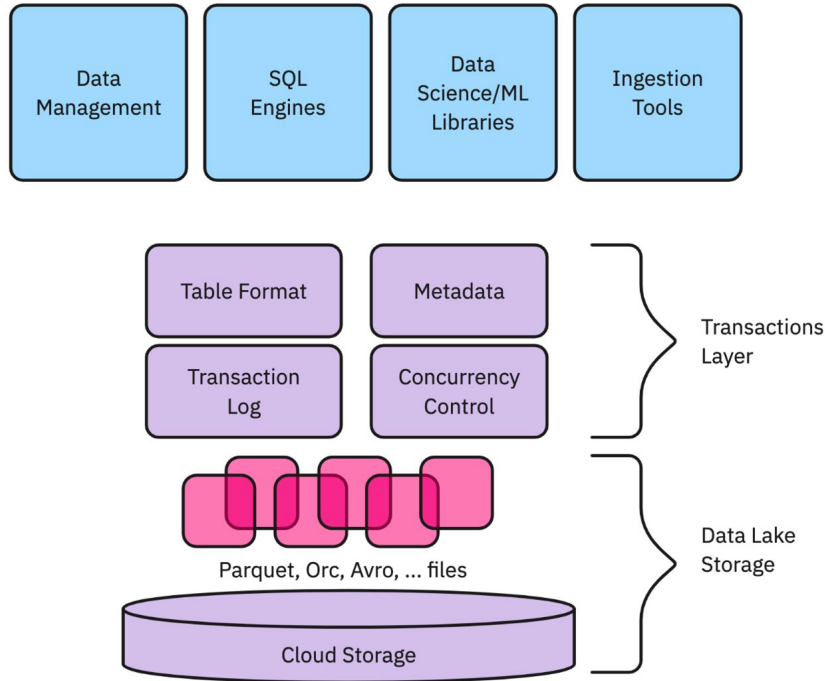
Architecture of a Database System



S3 Data Lake Storage



DATA LAKEHOUSE - UNBUNDLING OF THE DBMS



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 reduce 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 use 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, an text.
- **Support for diverse workloads:** including data science, machine learning, and SQL an 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



```

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



```

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:

```









(/ 4 4! +%| 4B? S ; L2 >P? LA? HNv


- 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

 Home

 Questions

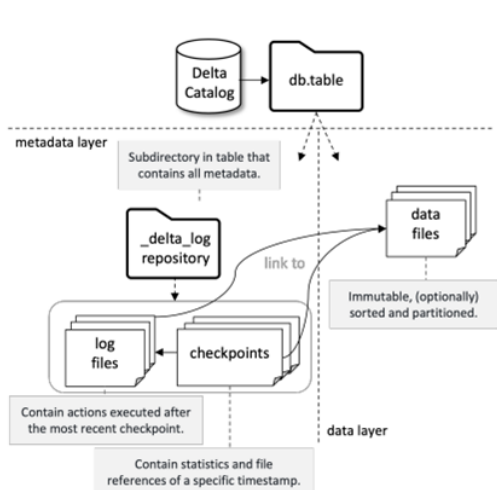
Thrift, Avro, Protocolbuffers - Are they all dead?
Asked 7 years, 3 months ago Modified 5 years, 11 months ago Viewed 40k times



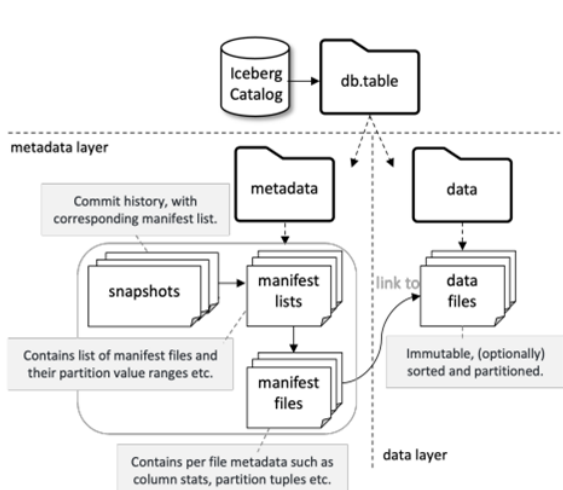


TECHNICAL FUNDAMENTALS

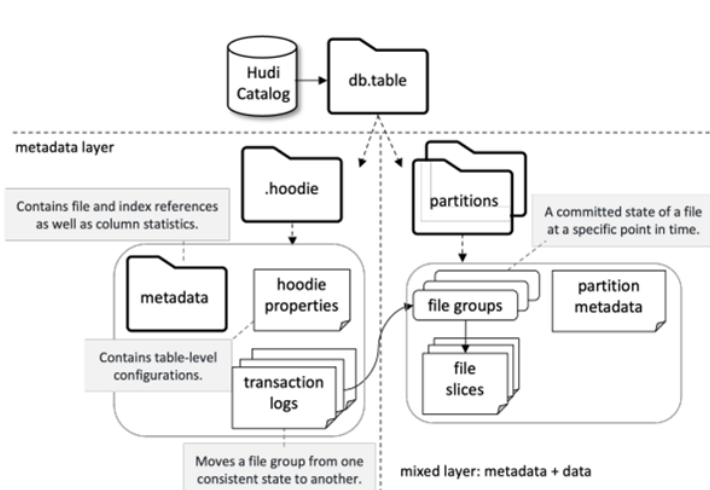
- Metadata abstractions on files in cloud object storage
- Tables with SQL semantics and schema evolution
- ACID transactions
- Updates and deletes (merge/upsert)
- Data layout optimizations for performance tuning



(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
```





WHICH FORMAT SHOULD I CHOOSE?

Choose  if:

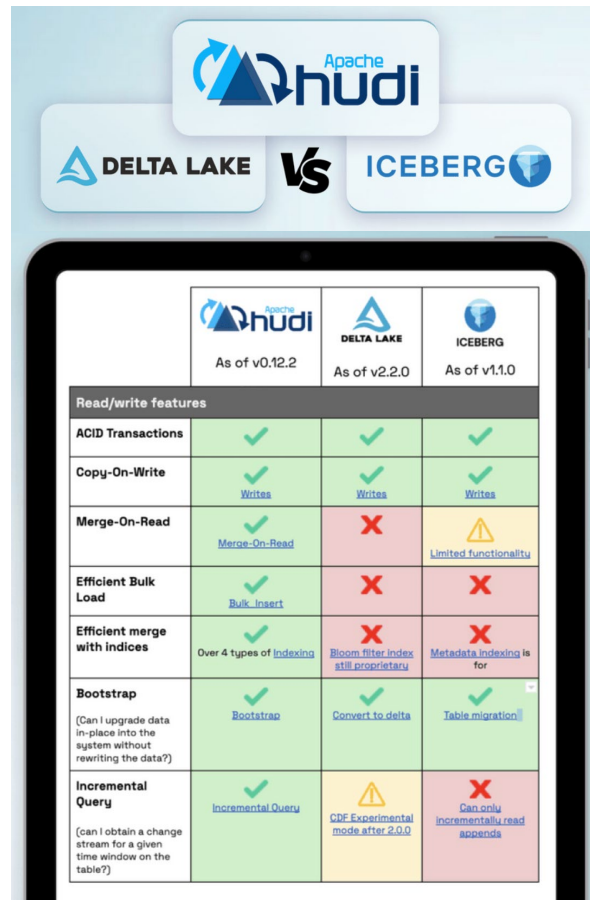
1. Mutable data - GDPR Deletes, Updates
2. CDC workloads
3. Low latency requirements
4. Large ETL pipelines - perf/cost w/ incremental ETL



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



	 As of v0.12.2	 As of v2.2.0	 As of v1.1.0
Read/write features			
ACID Transactions	✓	✓	✓
Copy-On-Write	✓ <small>Writes</small>	✓ <small>Writes</small>	✓ <small>Writes</small>
Merge-On-Read	✓ <small>Merge-On-Read</small>	✗	⚠ <small>Limited functionality</small>
Efficient Bulk Load	✓ <small>Bulk Insert</small>	✗	✗
Efficient merge with indices	✓ <small>Over 4 types of Indexing</small>	✗ <small>Bloom filter index still proprietary</small>	✗ <small>Metadata indexing is for</small>
Bootstrap <small>(Can I upgrade data in-place into the system without rewriting the data?)</small>	✓ <small>Bootstrap</small>	✓ <small>Convert to delta</small>	✓ <small>Table migration</small>
Incremental Query <small>(can I obtain a change stream for a given time window on the table?)</small>	✓ <small>Incremental Query</small>	⚠ <small>CDF Experimental mode after 2.0.0</small>	✗ <small>Can only incrementally load appends</small>



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:

2. Needs fastest performance in Spark with Photon
3. Wants an "easy-to-get-started" table format

Choose **ICEBERG** if:

1. Trino or Athena writes
2. Snowflake writes
3. Sensitive to performance
4. Partial table updates

What if you could work across all 3?

	Apache Hudi As of v0.12.2	DELTA LAKE As of v2.2.0	ICEBERG As of v1.1.0
Read/write features			
ACID Transactions	✓	✓	✓
Copy-On-Write	✓ <small>Writes</small>	✓ <small>Writes</small>	✓ <small>Writes</small>
Merge-On-Read	✓ <small>Merge-On-Read</small>	✗	⚠ <small>Limited functionality</small>
Efficient Bulk Load	✓ <small>Bulk Insert</small>	✗	✗
Efficient merge with indices	✓ <small>Over 4 types of Indexing</small>	✗ <small>Bloom filter index still proprietary</small>	✗ <small>Metadata indexing is for</small>
Bootstrap	✓ <small>Bootstrap</small>	✓ <small>Convert to delta</small>	✓ <small>Table migration</small>
Incremental Query	✓ <small>Incremental Query</small>	⚠ <small>CDF Experimental mode after 2.0.0</small>	✗ <small>Can only incrementally load appends</small>



EXAMPLE BENEFITS OF MIX-AND-MATCH

Writing

Choose  **Apache 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

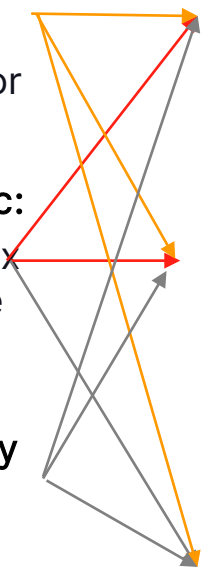
1. Get fastest queries with Photon acceleration
2. Great experience for Data Science

Choose  **ICEBERG** reading w/ Snowflake

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

Choose  **Apache HUDI** reading w/ DataProc (Spark)

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



INTRODUCING



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



The screenshot shows the GitHub repository for Apache XTable. The repository is public and has 591 stars, 82 forks, and 23 watchers. The commit history shows several commits by wuchunfu and jcamachor, all enforcing the ASF license header in non-Java files. The 'About' section describes XTable as a cross-table converter for lakehouse table formats, facilitating interoperability across data processing systems and query engines. It also provides links to the project website and related projects like Apache Iceberg, Delta Lake, and Apache Hudi.

>600
GH Stars ★

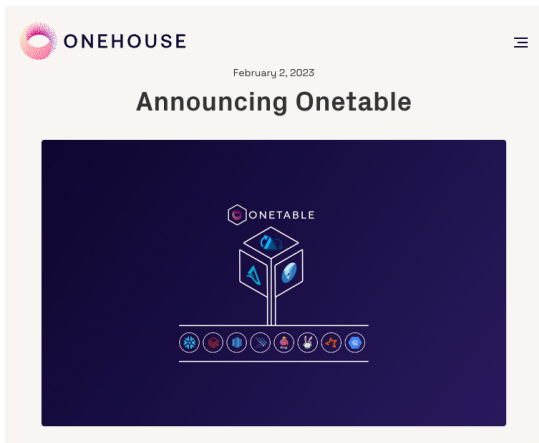
~] U
& LEM

OneTable
@OneTableOSS

Today we made it on the TOP PAGE of Github Repos Trending worldwide:
github.com/trending

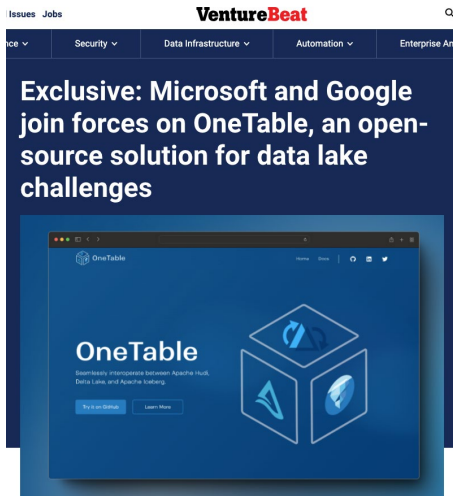
#OneTable #apachehudi #apacheiceberg #deltalake

Apache XTable™ TIMELINE



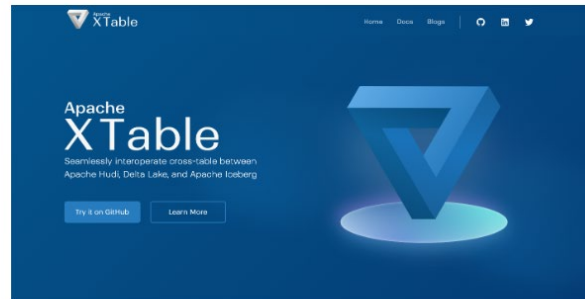
Onehouse announces
OneTable

Feb 2023



OSS Co-Launch with
Microsoft, Google, Onehouse

Nov 2023



OneTable is now "Apache XTable™ (Incubating)"

March 10, 2024

Dipankar Mazumdar, JB Onofré

Donation to ASF and
incubation as Apache XTable

Mar 2024

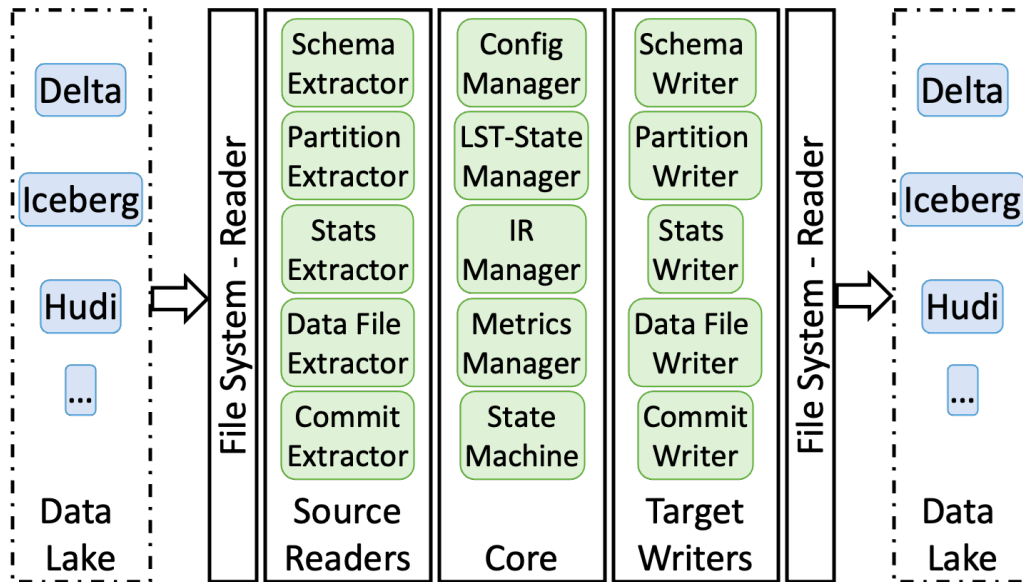




XTABLE ARCHITECTURE

Components:

1. Source Reader
1. Target Writer
1. Core Logic

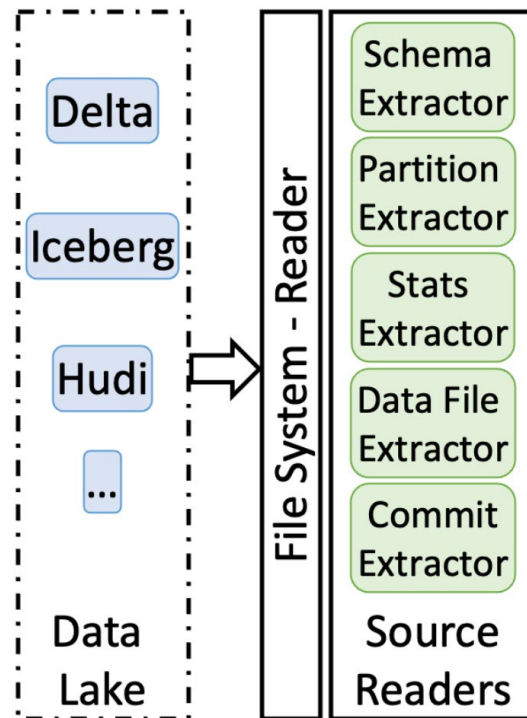


XTABLE ARCHITECTURE

Components:

1. Source Reader:

- LST-specific modules responsible for reading metadata from source table
- Operates using a pluggable file system
- Extracts info (schema, partition, transactions) & translates into **XTable's** internal representation

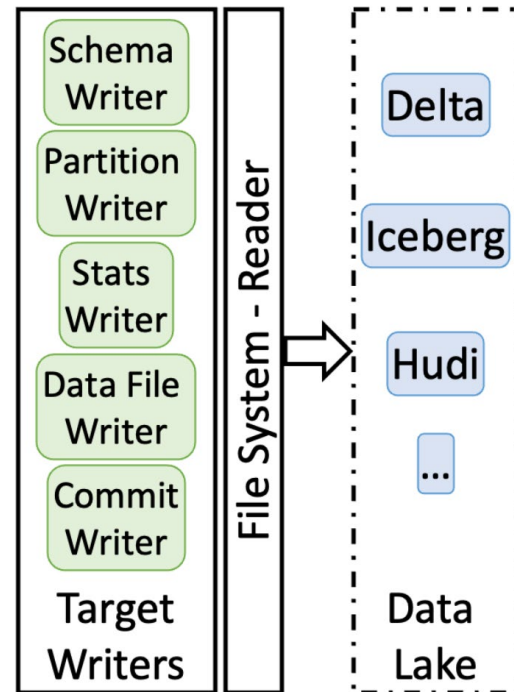


XTABLE ARCHITECTURE

Components:

2. Target Writer:

- takes the internal representation of metadata & accurately maps it to the target format's metadata structure
- includes re-creating schema, transaction logs & partition details in the new format

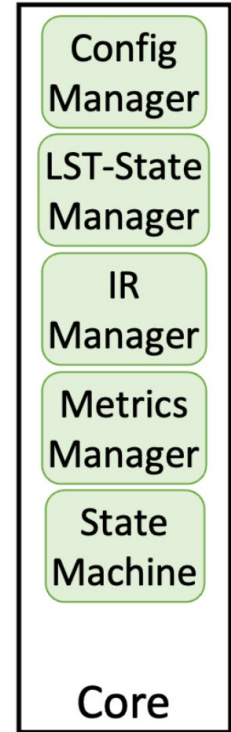


XTABLE ARCHITECTURE

Components:

3. Core Logic:

- CPU of XTable
- orchestrates the entire translation process
 - initialization of all components
 - managing sources/targets,
 - Handling tasks like caching for efficiency, state management for recovery & incremental processing, and telemetry for monitoring



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
yml
sourceFormat: HUDI
targetFormats:
  - DELTA
  - ICEBERG
datasets:
  -
    tableBasePath: s3://path/to/hudi-dataset/people # replace this with gs://path/to/hudi-dataset/people if y
    tableName: people
    partitionSpec: city:VALUE
```



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

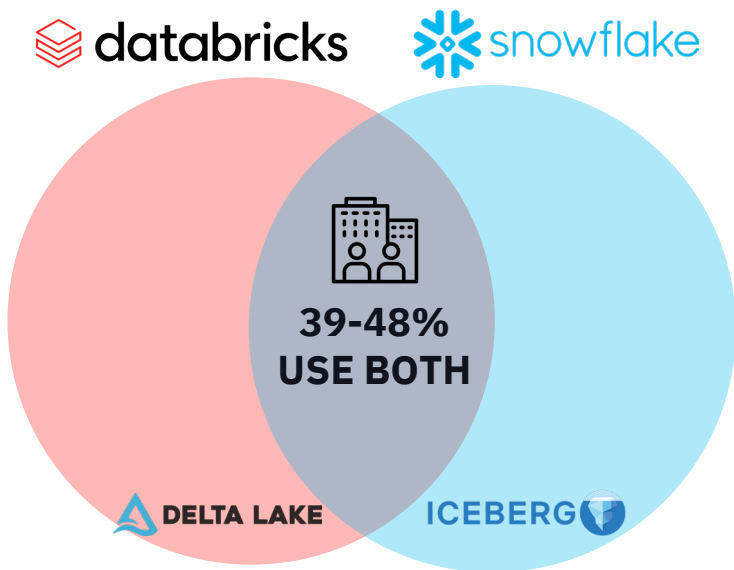


```
s3_bucket/my_table/
|- .hoodie/
|   |- hoodie.properties
|   |- metadata/
|- _delta_log/
|   |- 000000.json
|- metadata/
|   |- v1.metadata.json
|   |- snap-9fa1-2-16c3.avro
|   |- 0d9a-98fa-77.avro
|- file_1.parquet
|- file_2.parquet
|- file_N.parquet
```

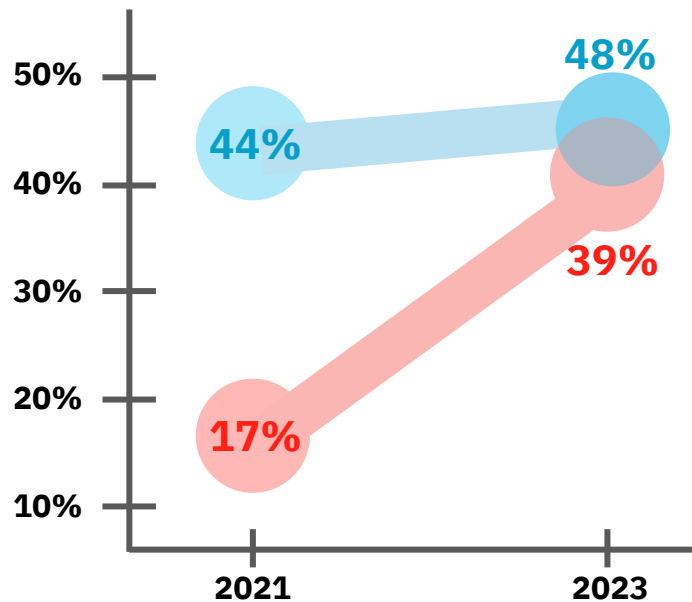
```
# any of these work on the same table
spark.read.format("hudi")
spark.read.format("delta")
spark.read.format("iceberg")
```



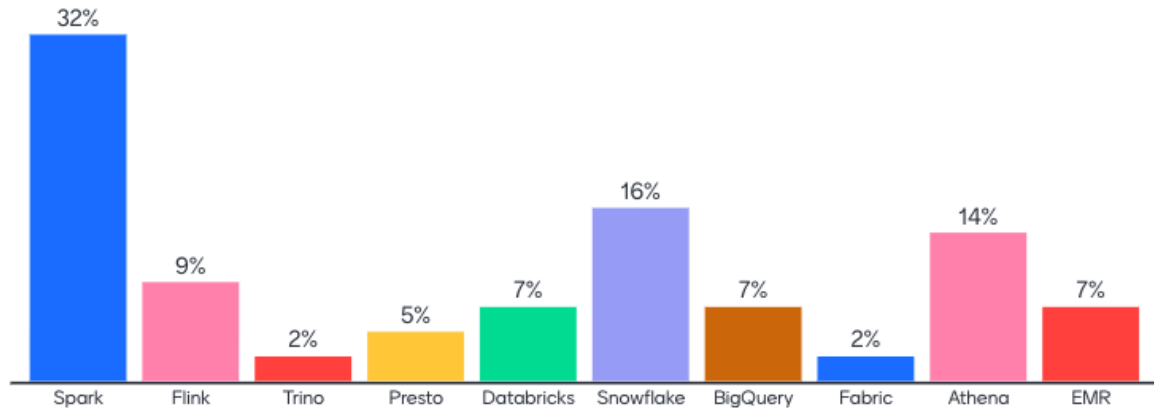
A TALE OF TWO...



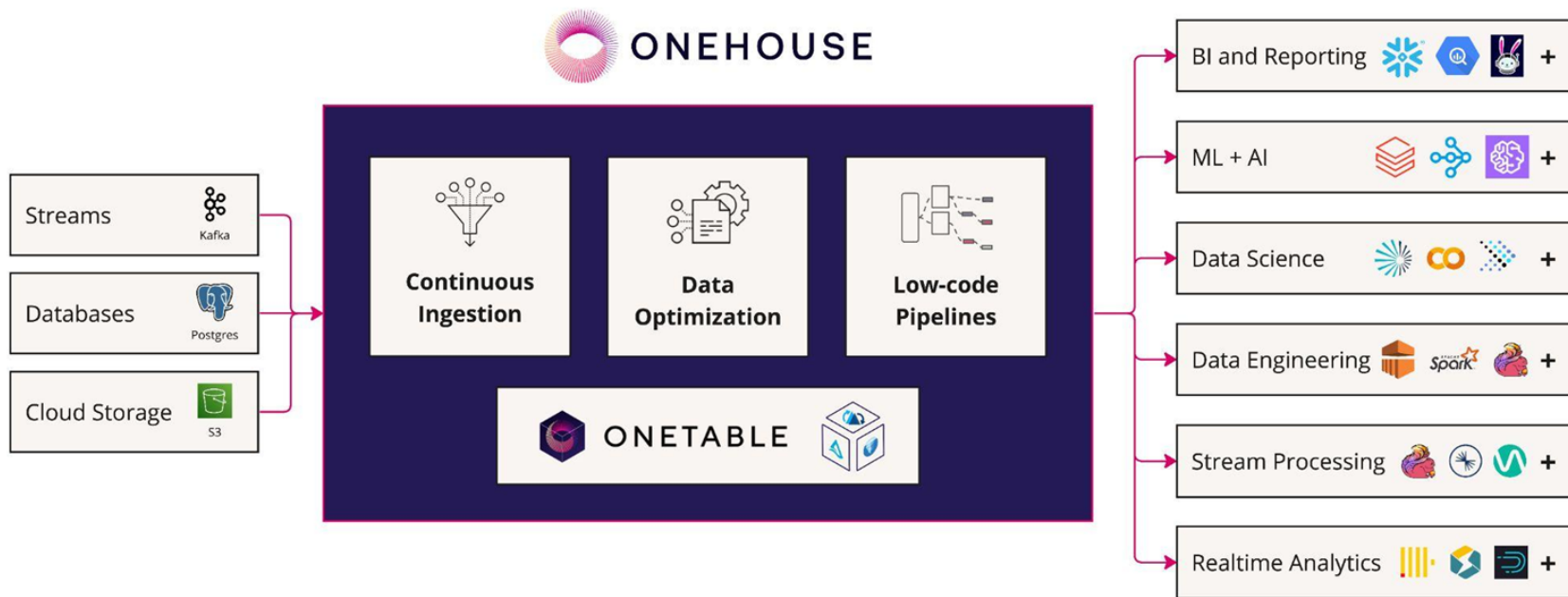
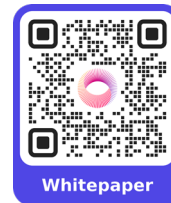
Overlap Growth 2021-2023



Which Query Engines Do You Use?



Apache XTable™ WHY BUILD IT?



Our Goal = Universal Data Lakehouse





Demo Time!



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

Initial Committers

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

Contributors, Users, Mentors



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
- Data Sharing
- Catalog



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/apachexable>



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



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



Thank You!

