

# ***APACHE ICEBERG***

***DATA***  ***LAKE***

***Which Framework Should You***

***Choose for Your Data Lake?***

# Why Data Lakes Need Reliable Frameworks

- Data lakes are crucial for storing massive datasets, but they often face issues like:
  - ✗ Messy organization
  - ✗ Inconsistent schemas
  - ✗ Slow performance

# How Iceberg and Delta Lake Help

Apache Iceberg  and Delta Lake 





improve data lakes with:

- ✓ ACID transactions for reliability
- ✓ Schema evolution for flexibility
- ✓ Time travel for analytics

# What Is an Open Table Format? 🤔

- An open table format is a blueprint for managing and querying data in data lakes.
- It adds structure and reliability on top of raw files.





# Features of Open Table Formats

-  Standardized metadata: Tracks schema and partitioning.
-  ACID compliance: Ensures reliable updates and deletes.
-  Compatibility: Works with Spark, Flink, Trino, etc.
-  Advanced features: Time travel, partition pruning.




# Apache Iceberg: The Flexible Organizer

- Built by Netflix for scalable data lakes.
- Works across multiple tools and cloud environments.

# Key Features of Apache Iceberg

-  Hidden Partitioning: Auto-organizes data for faster queries.
-  Schema Evolution: Update schemas without rewriting data.
-  Time Travel: Query historical data.
-  Engine Neutrality: Works with Spark, Flink, Hive, etc.

# Technical Design of Iceberg

-  Metadata Layer: Tracks structure and schema.
-  Snapshots: Version history for rollbacks and time travel.
-  Multi-Engine APIs: Standardized APIs for compatibility.



# Delta Lake: The Spark Powerhouse



- Created by Databricks for Apache Spark.
- Designed for real-time analytics and batch processing.

# Key Features of Delta Lake

- ✅ ACID Transactions: Ensures consistency during concurrent writes.
- ⚡ Batch + Streaming: Real-time and historical data in one table.
- 🕒 Time Travel: Rollback to previous versions.
- 🔗 Spark Optimization: Boosts performance with Spark.

# Technical Design of Delta Lake



Delta Log: Tracks every change to maintain ACID compliance.



Schema Enforcement: Prevents inconsistent data entry.



Partitioning and Z-Ordering: Speeds up queries with smart indexing.

# GitHub Repositories for Iceberg and Delta Lake






Apache Iceberg: <https://github.com/apache/iceberg>



Delta Lake: <https://github.com/delta-io/delta>

Visit these repositories for source code and community updates.



# Apache Iceberg Overview

-  Best for multi-engine setups (e.g., Spark, Flink).
-  Offers hidden partitioning and schema flexibility.
-  Ideal for multi-cloud environments.

# Delta Lake Overview 🔥

- 🔧 Best for Spark-first workflows.
- ⚡ Combines real-time streaming and batch processing.
- 💡 Ideal for Spark and Databricks ecosystems.

# How to Decide?

-  Choose Iceberg if:
  - - You use multiple tools (e.g., Flink, Trino).
  - - You need flexible partitioning.
-  Choose Delta Lake if:
  - - You rely heavily on Spark.
  - - You need real-time and batch capabilities.

# Supporting the Lakehouse Design



Lakehouse architecture combines data lakes with structured warehouse-like features.



Iceberg: Engine-agnostic and cloud-flexible.



Delta Lake: Optimized for Spark-heavy workflows.



# Key Takeaways

Both Apache Iceberg and Delta Lake are open table formats.

Iceberg excels in flexibility and multi-engine setups.

Delta Lake is ideal for Spark-focused, real-time workflows.

Evaluate tools and workloads to choose the best fit.

# Engage With Us!

#DataEngineering #ApacheIceberg #DeltaLake  
#OpenTableFormats #BigData #DataLakes