

08.10.2025

Apache Iceberg on AWS

David Greenshtein

Senior Solutions Architect – Analytics, AWS



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Agenda

Introduction to Open Table Formats

Apache Iceberg

- Overview
- Feature Highlights
- Table Strategy and table Maintenance

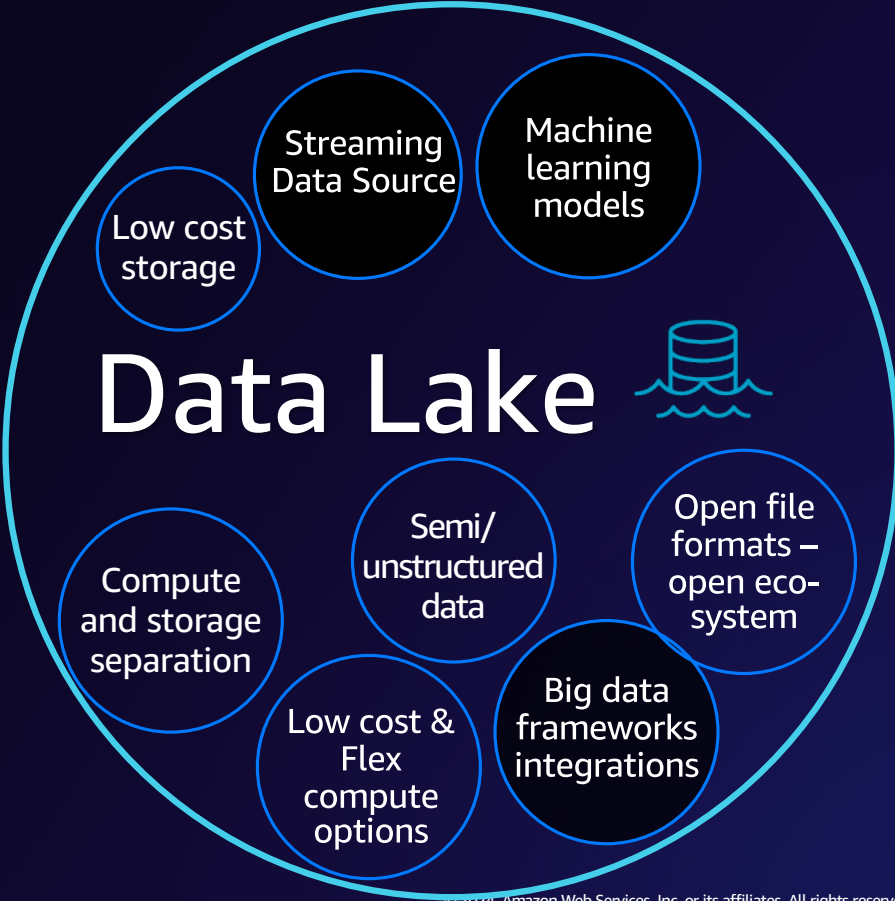
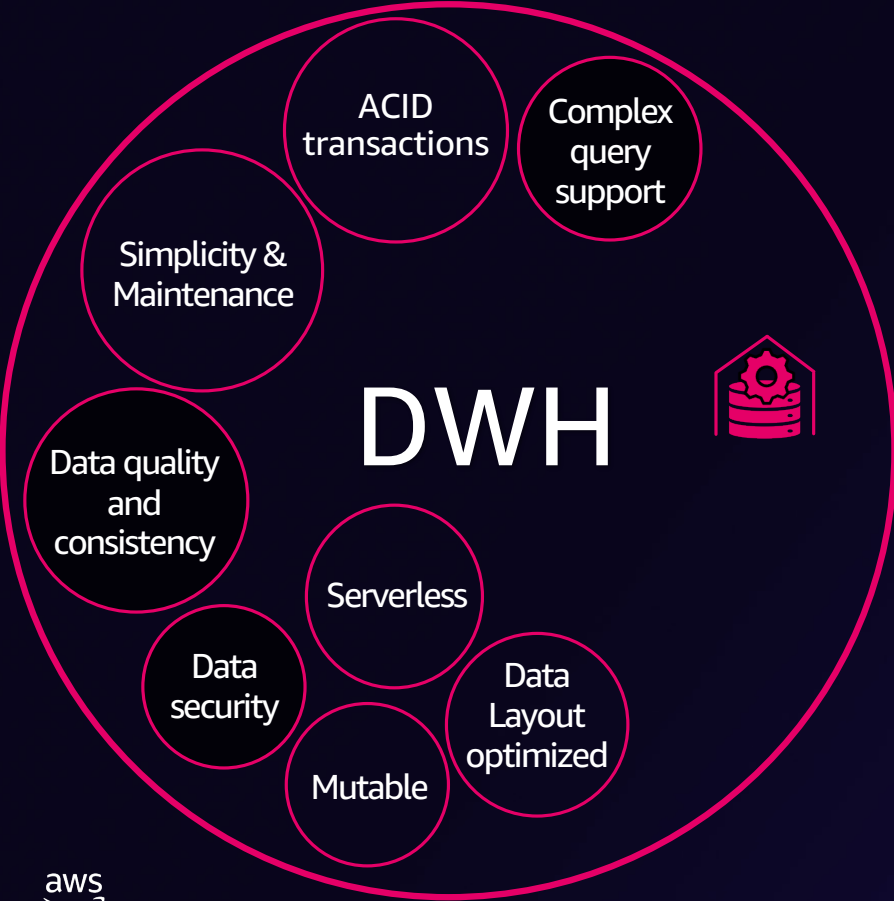


Open Table Formats



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Platforms are evolving



Customer requirements are evolving also

Can we do all of this with the same data governance and open standards?

Can we decouple storage from compute and support diverse consumers?

Can we bring the open source flexibility to the data warehouse?

DATA
WAREHOUSE

DATA
LAKE

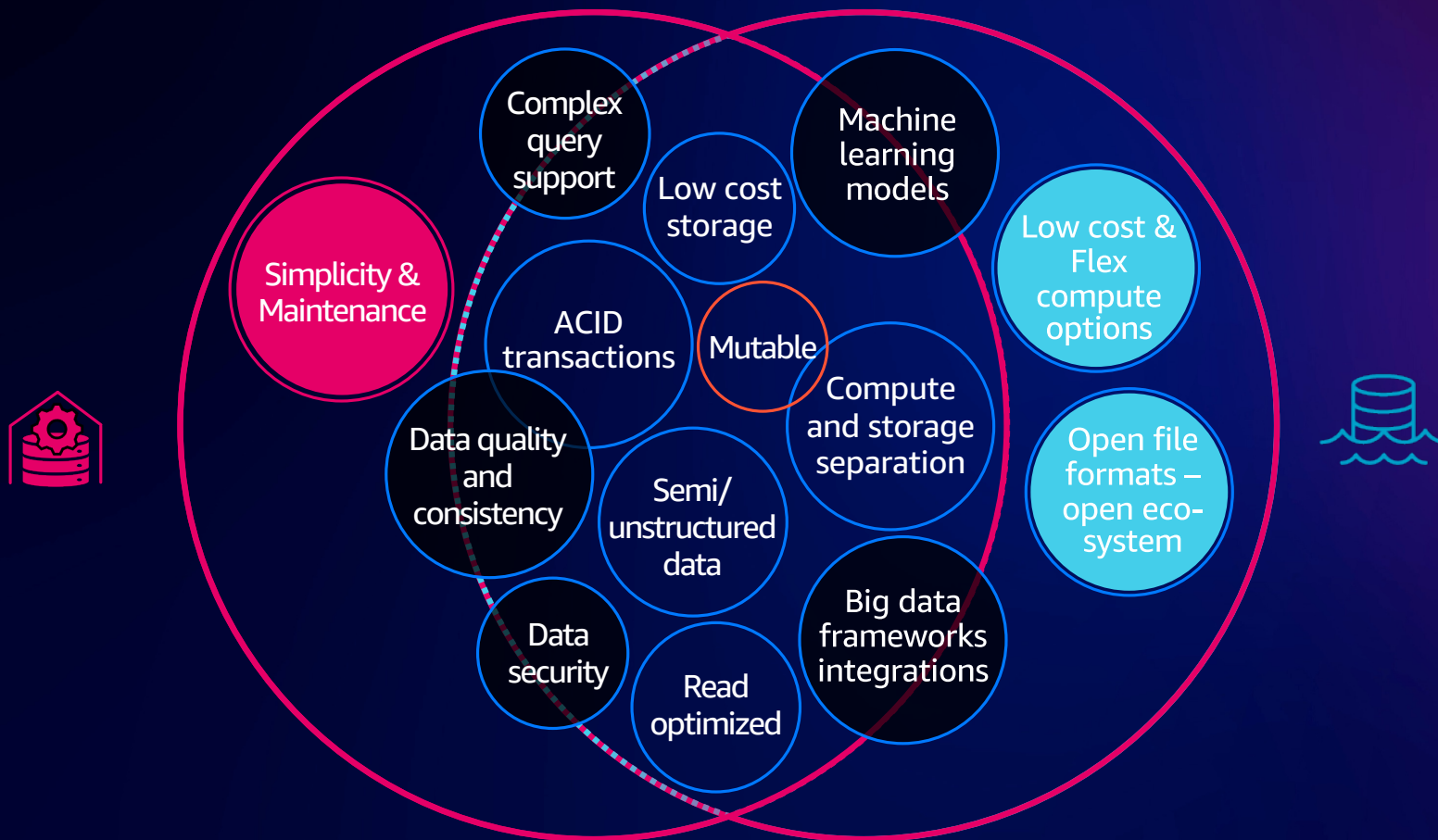
Can we deliver E2E Governance?

Can we bring the performance and strong ACID properties of data warehouse to data lakes?

Can you deliver best price/performance?



Modern Data Lakes using an Open Table Format



Modern Data Lakes

Open table formats (OTFs) provide transactional support and simplify data lake optimization and management



Apache Hudi



Apache Iceberg



Delta Lake

Apache Iceberg Overview



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

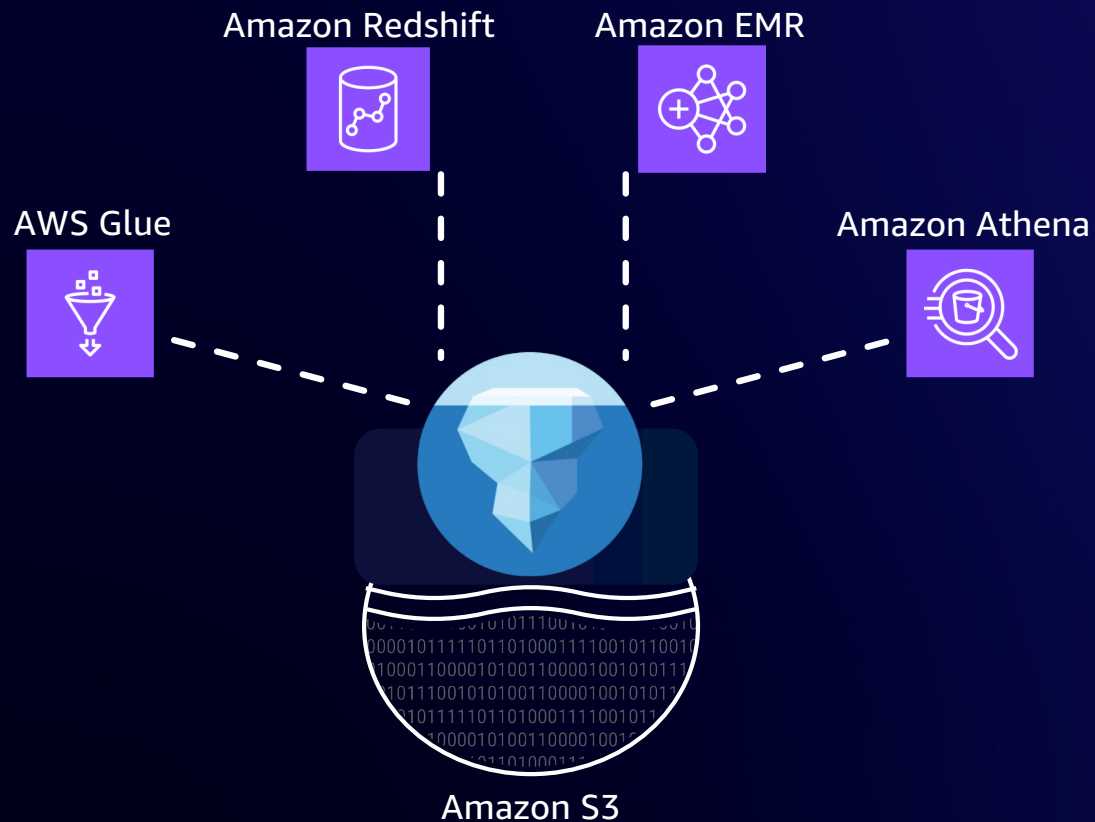
What is Apache Iceberg?

Apache Iceberg is an open table format for huge analytic datasets

- Open table format specification
- Set of libraries
- Java / Python / Rust APIs



AWS Apache Iceberg Integration

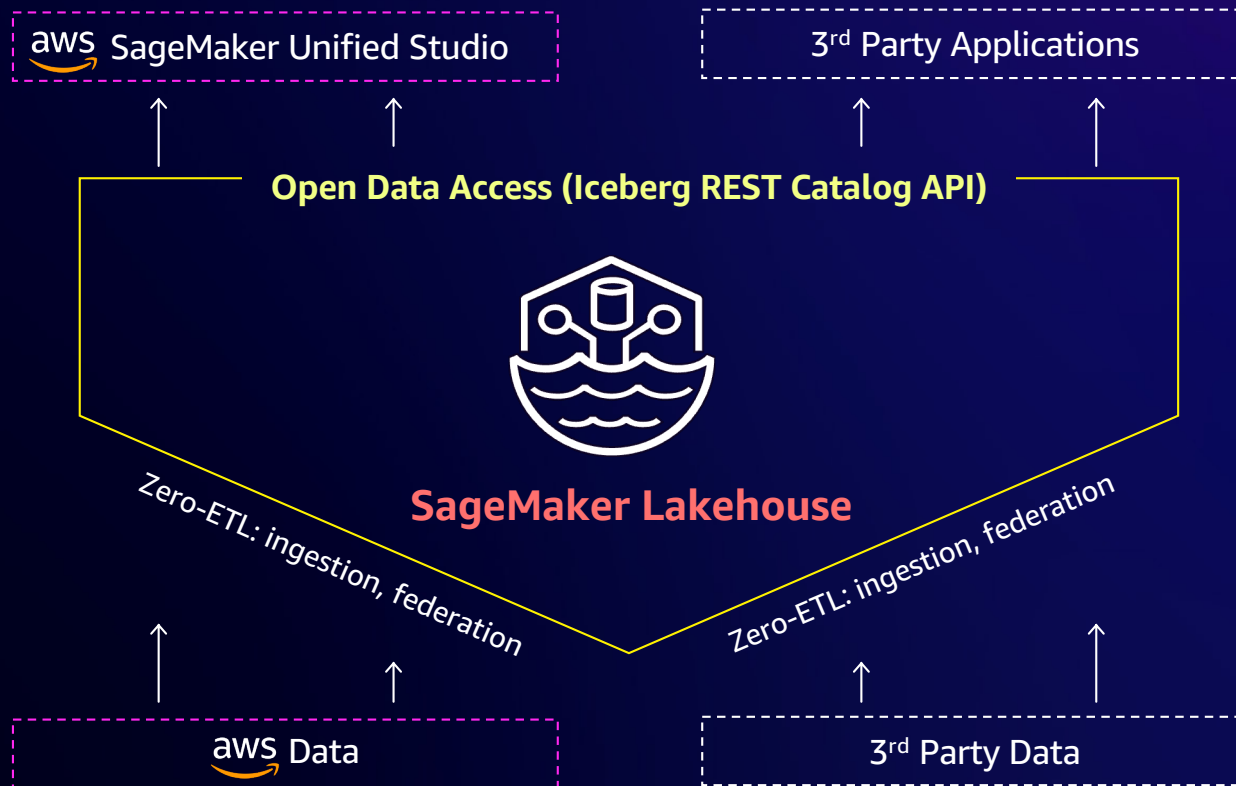


- EMR (v6.5+) - Spark, Flink, Trino, Presto
- Athena (v2 & v3)- DML, DDL, maintenance
- Glue - Interactive sessions, jobs, maintenance
- Redshift – Spectrum, Serverless



SageMaker Lakehouse

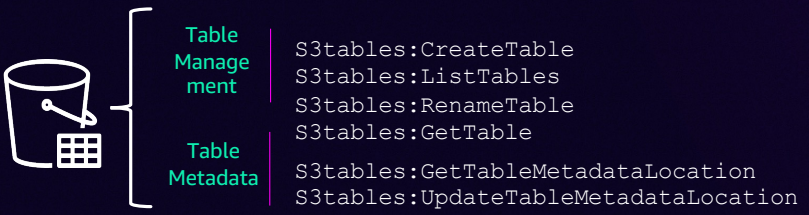
FIRST PARTY AND THIRD PARTY QUERY ENGINES AND TOOLS



NEW

Amazon S3 Tables

FULLY MANAGED STORAGE FOR APACHE ICEBERG DATA LAKE



New S3 storage class for Apache Iceberg data lakes

Amazon S3 APIs to read/write to S3 tables

Fully managed Iceberg table maintenance

Simple integration with the Lakehouse

10x Requests per second compared to standard Amazon S3 buckets



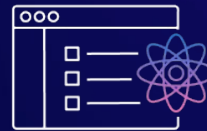
Apache Iceberg Features



Transactional updates



Time travel queries



Schema & partition evolution



Optimization & compaction



How Iceberg works



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Iceberg Table Anatomy

AWS Services



AWS Glue
Data Catalog

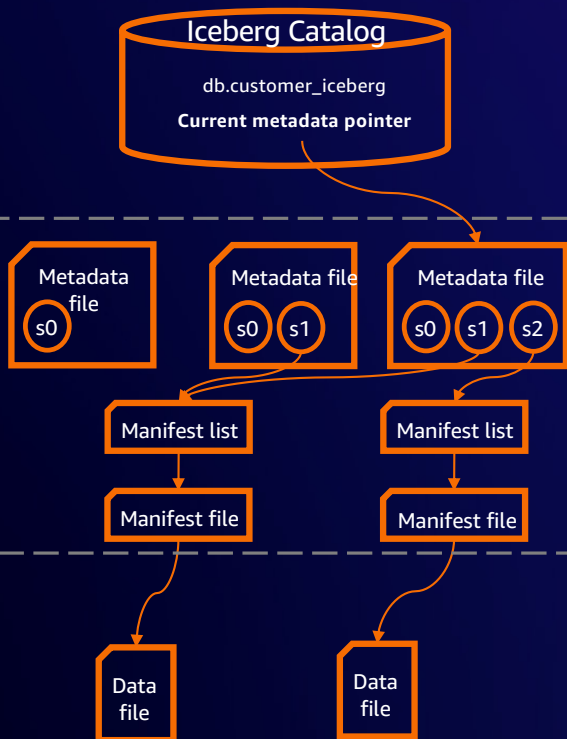


Amazon S3



Amazon S3

Iceberg Table Anatomy



Storage Layout

S3://datalake/iceberg_table/

```
| - metadata/  
|   | - v1.metadata.json  
|   | - ...  
|   | - v3.metadata.json  
|   | - snap....  
|   | - snap.0j2s.avro  
|   | - ...avro  
|   | - dg26.avro
```

```
| - data/  
|   | - partition_col=2023-08-10/  
|   |   | - 3dh3.parquet  
|   |   | - jd73.parquet  
|   |   | - gs62d.parquet
```



Iceberg Table Anatomy

Iceberg Table Anatomy

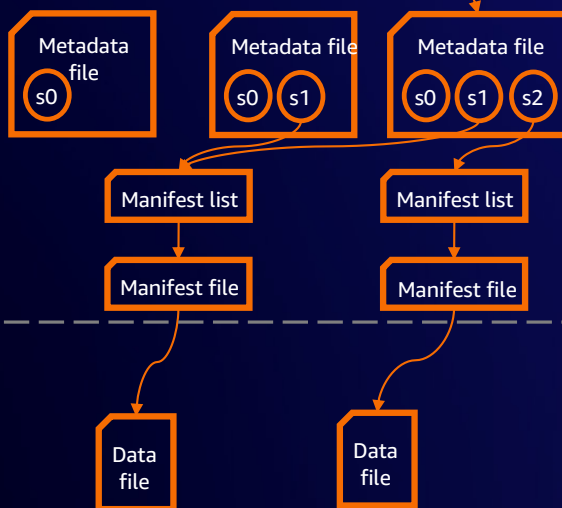
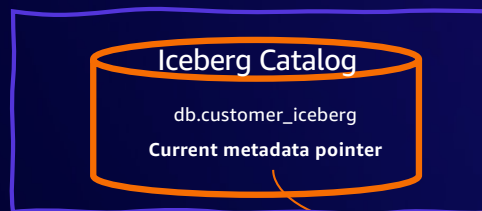
Storage Layout

Catalog

- Pointer to metadata file
- Unlike Hive Metastore, almost nothing else
- Supports atomic swap of the metadata files pointer

Metadata pointer

- S3 path to the latest metadata.json file



S3://datalake/iceberg_table/

```
| - metadata/  
|   | - v1.metadata.json  
|   | - ...  
|   | - v3.metadata.json  
|   | - snap....  
|   | - snap.0j2s.avro  
|   | - ...avro  
|   | - dg26.avro
```

```
| - data/  
|   | - partition_col=2023-08-10/  
|   |   | - 3dh3.parquet  
|   |   | - jd73.parquet  
|   |   | - gs62d.parquet
```



Iceberg Table Anatomy

Iceberg Table Anatomy

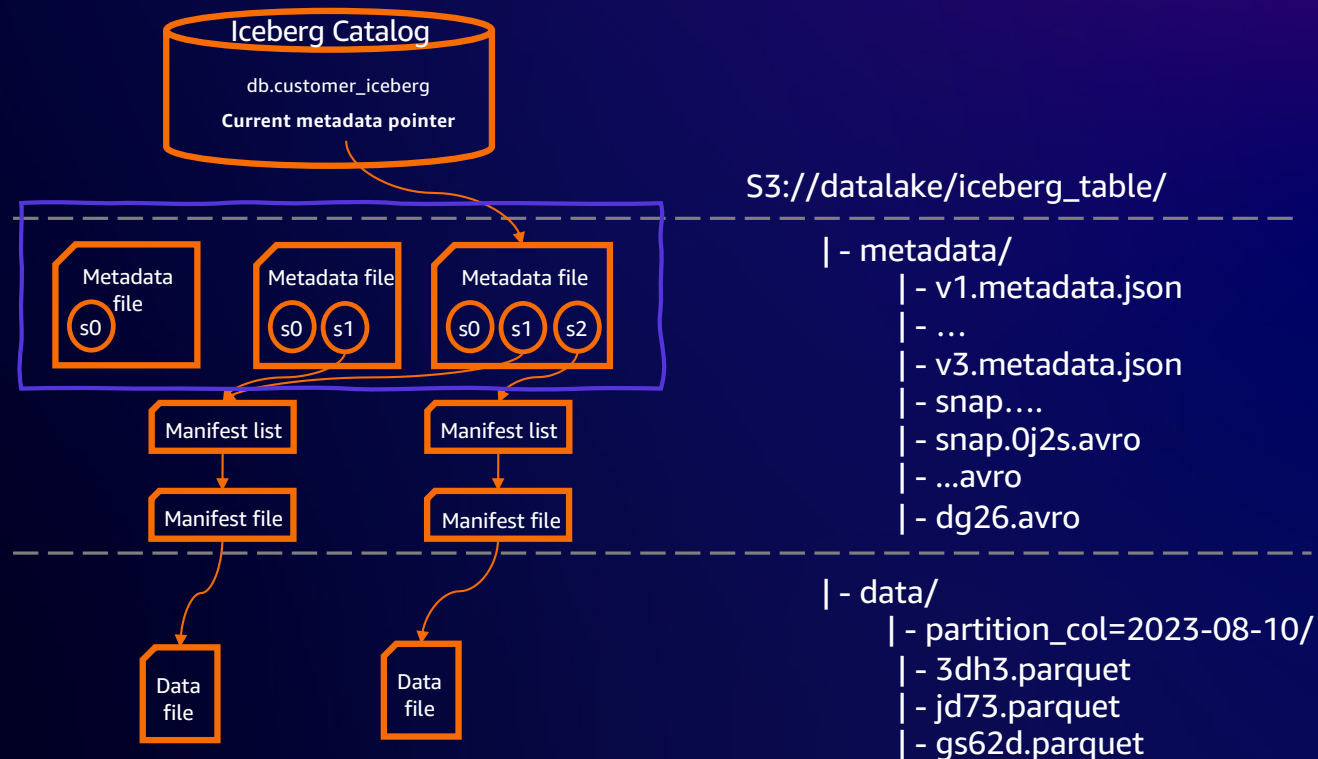
Storage Layout

Metadata file (json)

- “Root” of the table
- Contains table level metadata (schema, partitions, snapshots, properties)

Snapshot

- Represents the state of the table at a given point in time
- Points to the manifest list



Iceberg Table Anatomy

Iceberg Table Anatomy

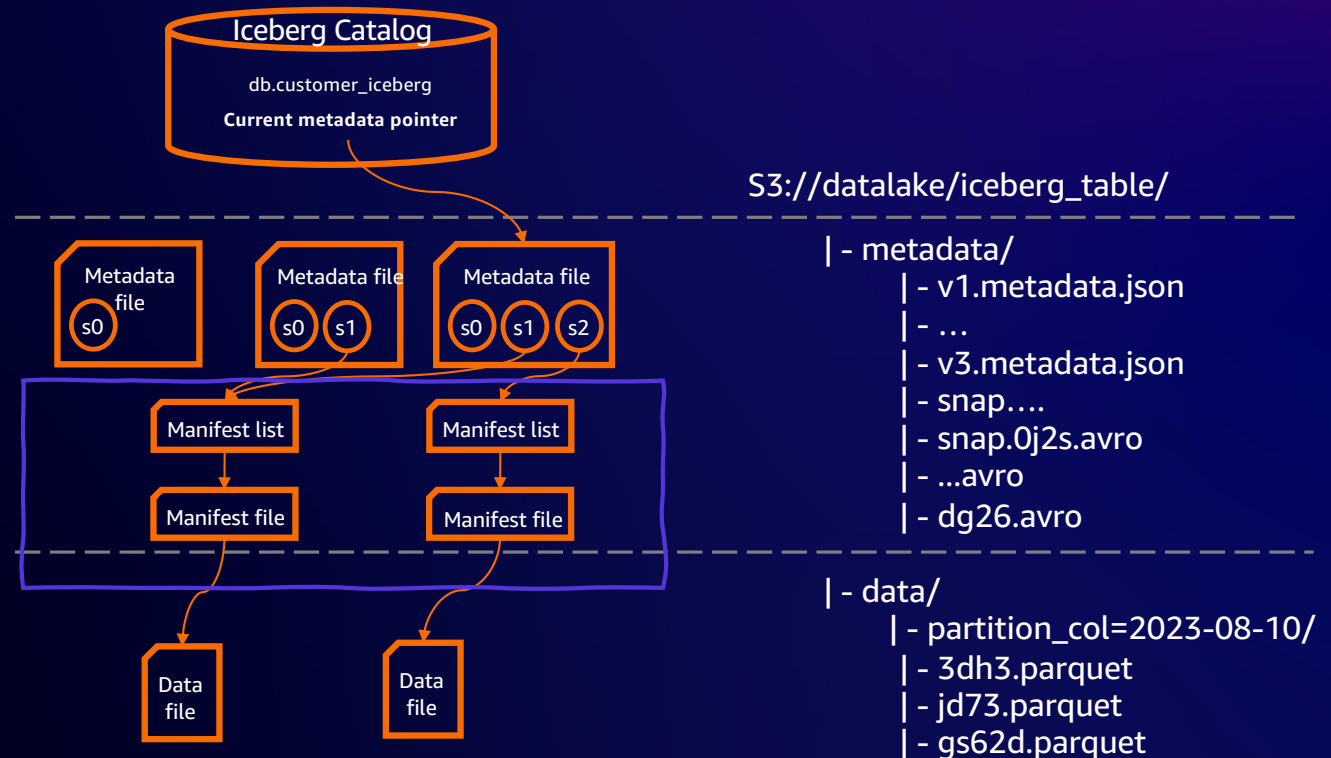
Storage Layout

Manifest list (avro)

- Pointers to Manifest files
- 1st level index of the table (i.e. high level statistics for query optimization)

Manifest file (avro)

- Pointers to data files
- 2nd level index of the table (i.e. data statistics for query optimization)



Quick Example

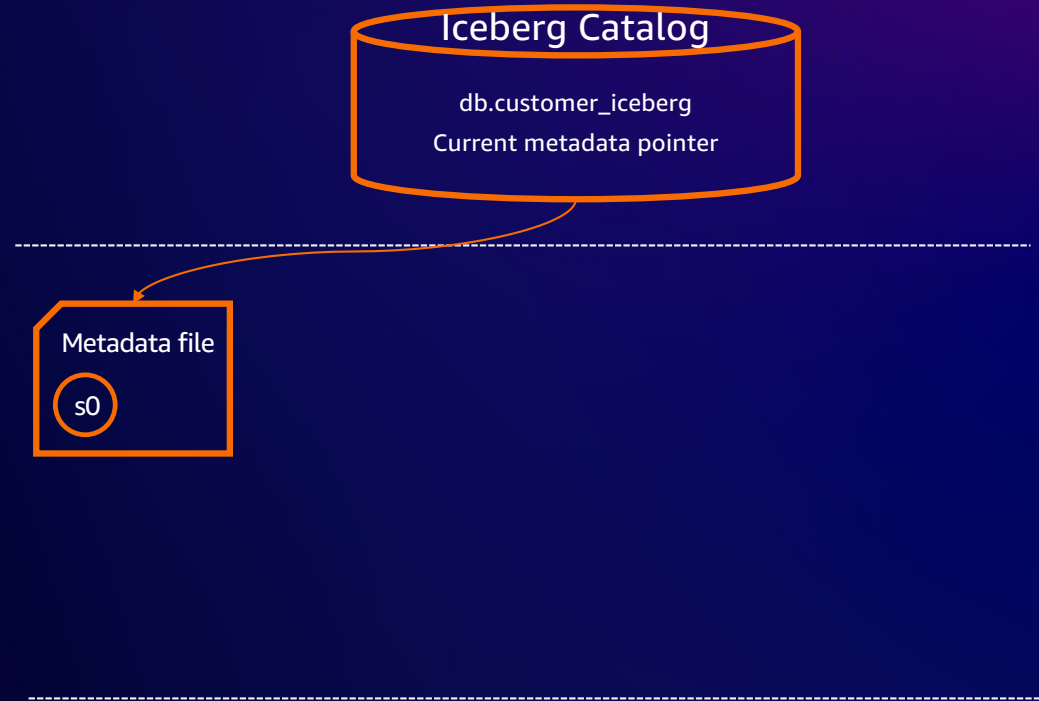
```
CREATE TABLE product_inventory (  
  product_id BIGINT,  
  product_name STRING,  
  category STRING  
)  
PARTITIONED BY (category)  
LOCATION 's3://my-bucket/product-inventory/'  
TBLPROPERTIES (  
  'table_type'='ICEBERG',  
  'format'='parquet'  
);
```

```
INSERT INTO product_inventory  
VALUES  
  (1, 'Microsoft Office', 'Software'),  
  (2, 'Adobe Photoshop', 'Software'),  
  (3, 'AutoCAD', 'Software'),  
  (4, 'Visual Studio', 'Software');
```



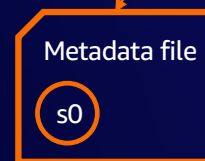
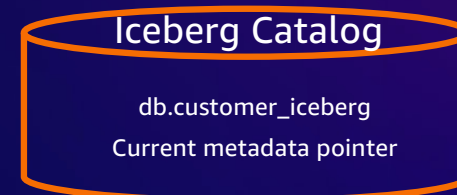
Create Table

S3://my-bucket/product-inventory/
/metadata/
 /v1.metadata.json
/data/



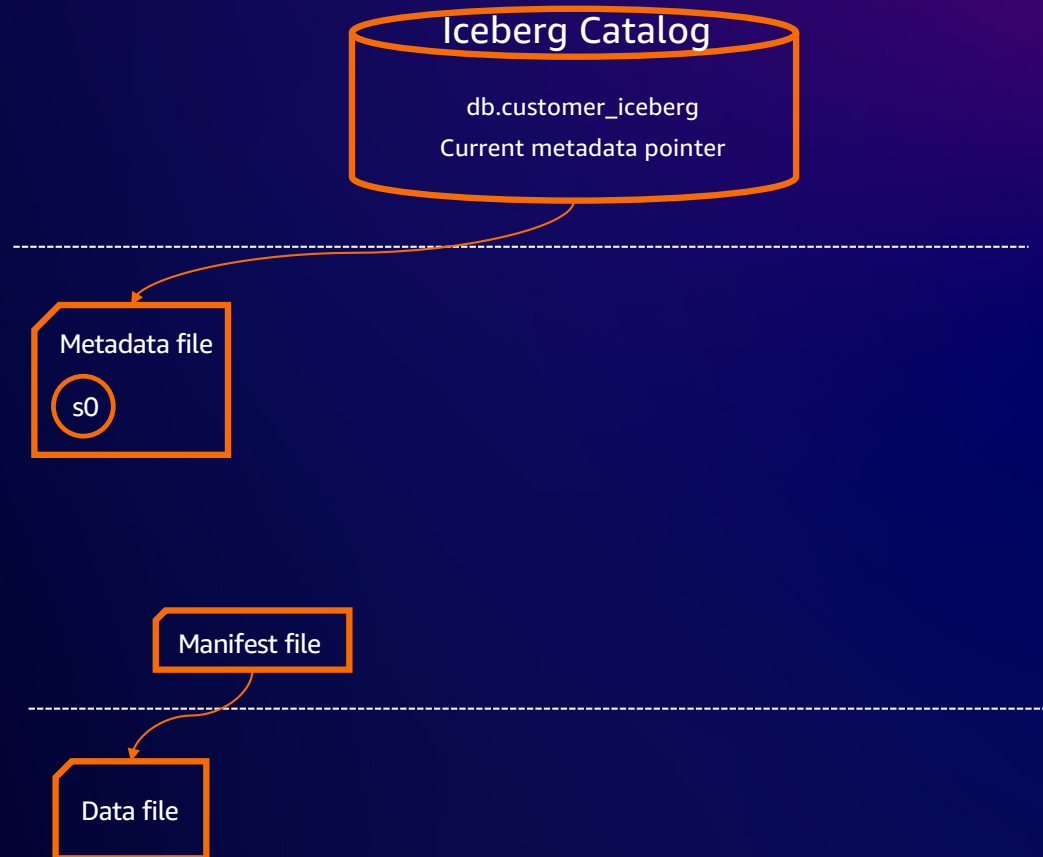
Insert Records

S3://my-bucket/product-inventory/
/metadata/
 /v1.metadata.json
/data/
 /category=Software/
 /3dh3.parquet



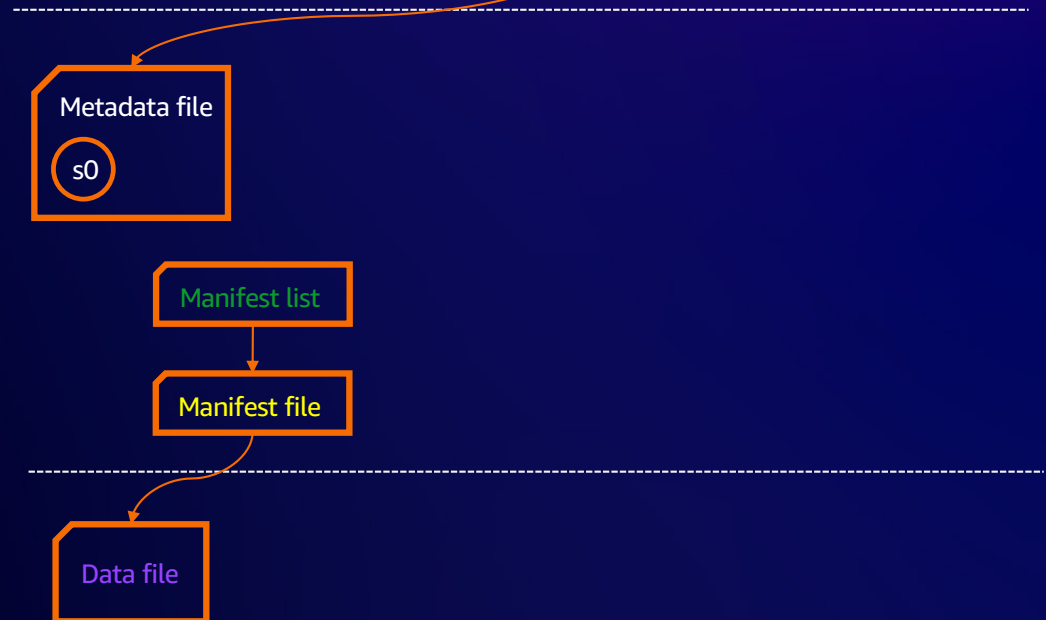
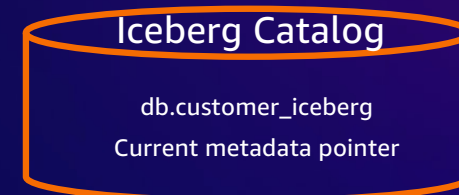
Insert Records

S3://my-bucket/product-inventory/
/metadata/
 /v1.metadata.json
 /dbda7.avro
/data/
 /category=Software/
 /3dh3.parquet



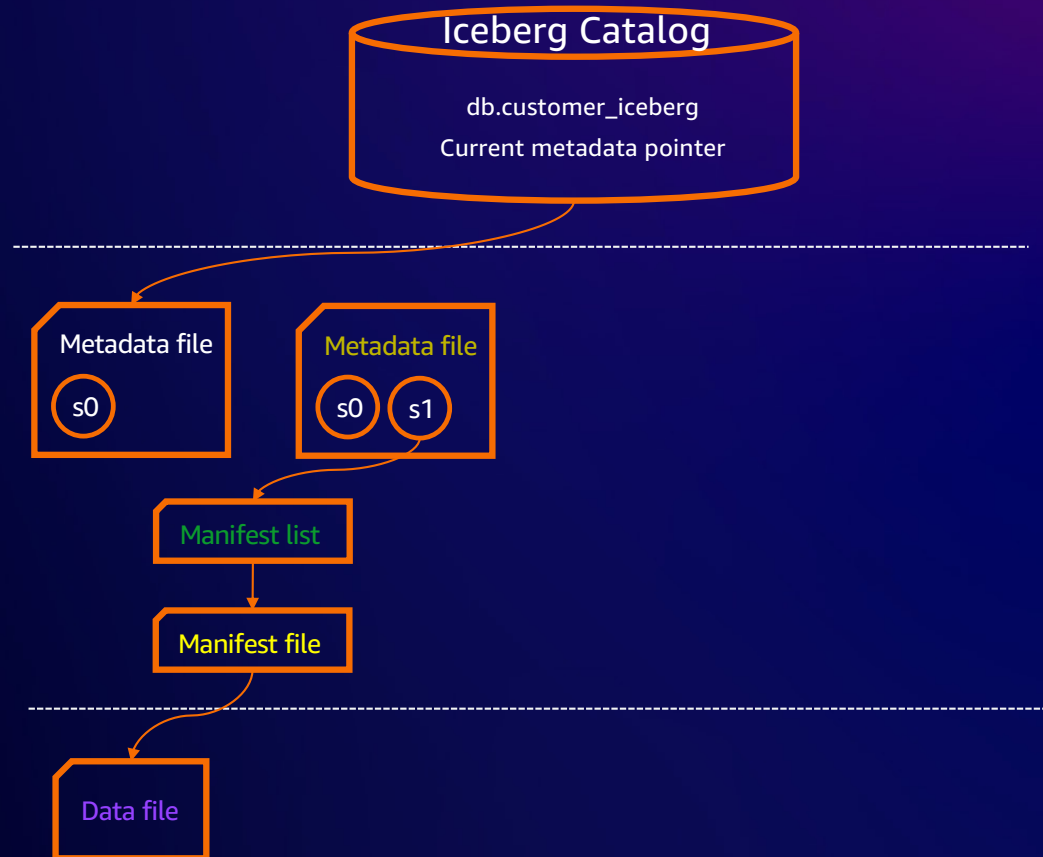
Insert Records

S3://my-bucket/product-inventory/
/metadata/
 /v1.metadata.json
 /snap.3ejn.avro
 /dbda7.avro
/data/
 /category=Software/
 /3dh3.parquet



Insert Records

S3://my-bucket/product-inventory/
/metadata/
 /v1.metadata.json
 /v2.metadata.json
 /snap.3ejn.avro
 /dbda7.avro
/data/
 /category=Software/
 /3dh3.parquet



Insert Records

S3://my-bucket/product-inventory/
/metadata/

/v1.metadata.json

/v2.metadata.json

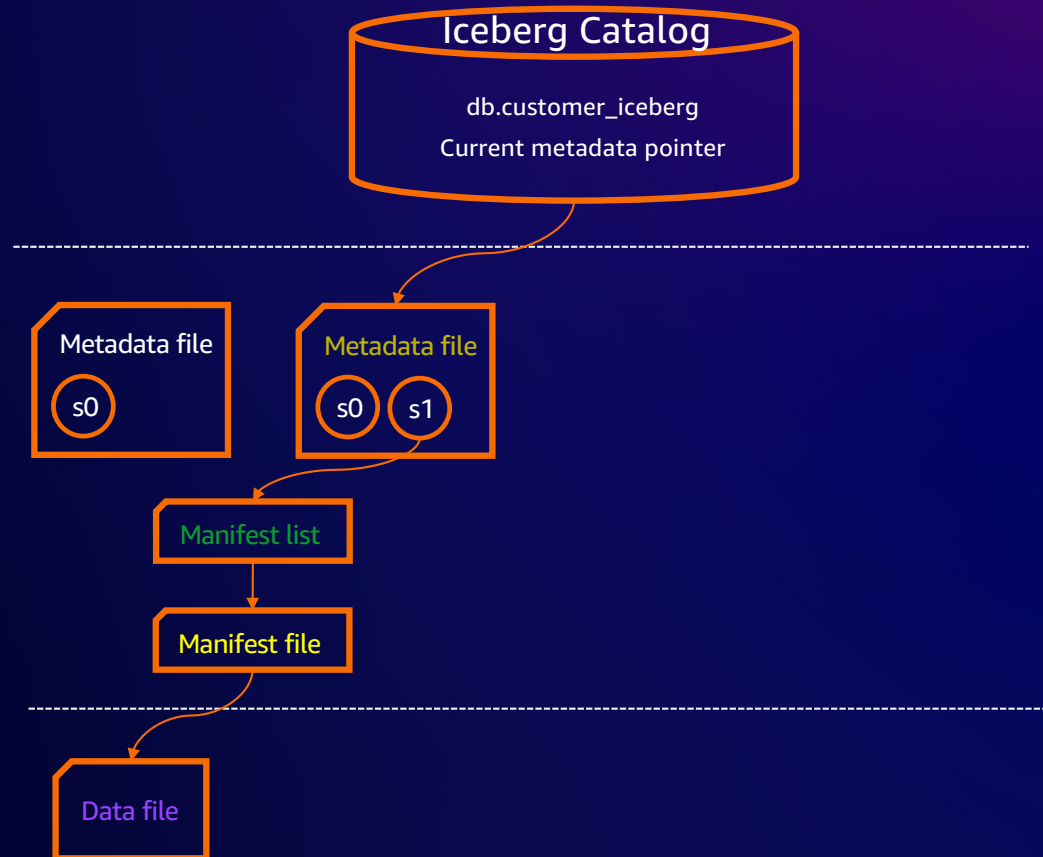
/snap.3ejn.avro

/dbda7.avro

/data/

/category=Software/

/3dh3.parquet



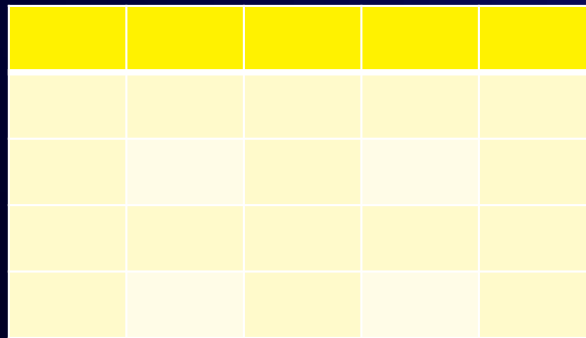
Apache Iceberg Feature Highlights



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

#1 Expressive SQL

Iceberg Table




```
UPDATE tablename SET xx=yy[,...]  
[WHERE predicate]
```

```
DELETE FROM tablename  
[WHERE predicate]
```

```
MERGE INTO tablename s  
USING targetable t  
ON s.x = t.x  
WHEN ...
```



#2 Schema & Partition Evolution

Schema evolution changes are **independent and free of side-effects** without rewriting files

Supported schema evolution changes (even in **nested** structures)

- **Add** column
- **Drop** column
- **Rename** column
- **Reorder** column
- **Update** column



#2 Schema & Partition Evolution

Iceberg supports changing partition layout

Table sales

Partitioned by month(date)

2023-01

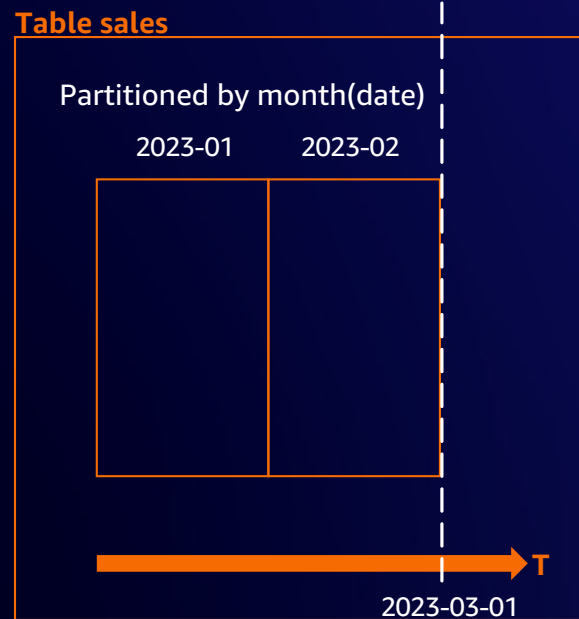
2023-02



#2 Schema & Partition Evolution

Just a metadata operation

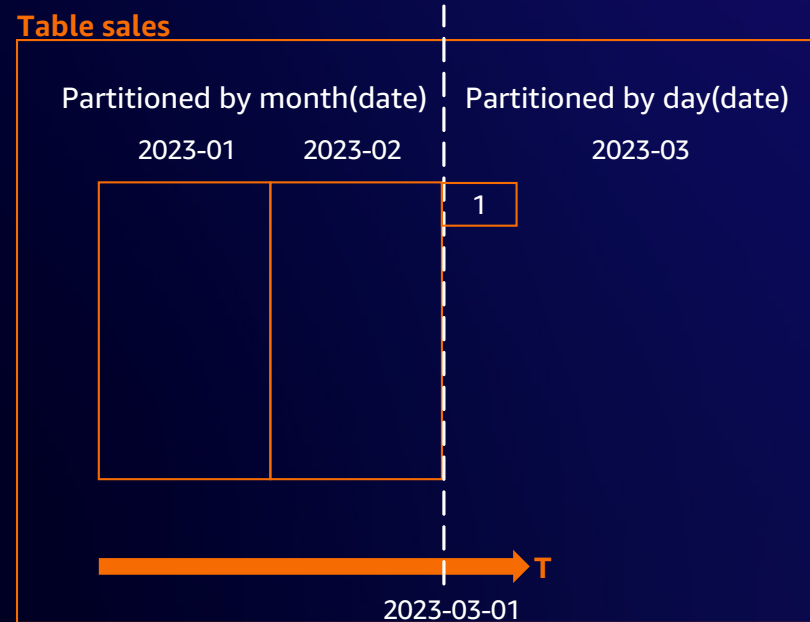
ALTER TABLE sales REPLACE FIELD
Month(date) WITH days(date)



#2 Schema & Partition Evolution

Just a metadata operation

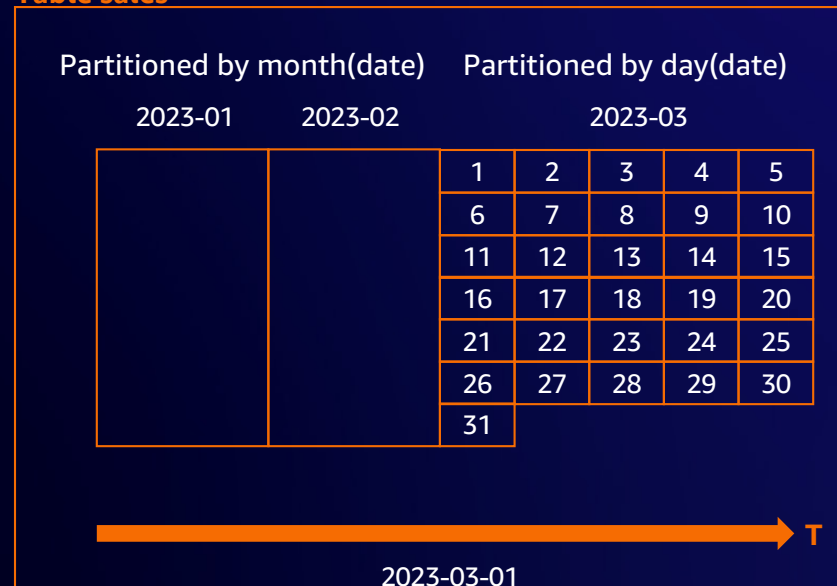
ALTER TABLE sales REPLACE FIELD
Month(date) WITH days(date)



#2 Schema & Partition Evolution

There are no additional costs such as moving data, migrating tables, or rewriting queries

Table sales



#2 Schema & Partition Evolution

Query example

```
SELECT * FROM sales
WHERE
  date > 2023-02-11 AND
  date < 2023-03-18
```

Table sales

Partitioned by month(date)		Partitioned by day(date)				
2023-01	2023-02	2023-03				
		1	2	3	4	5
		6	7	8	9	10
		11	12	13	14	15
		16	17	18	19	20
		21	22	23	24	25
		26	27	28	29	30
		31				

→ T

2023-03-01

#2 Schema & Partition Evolution

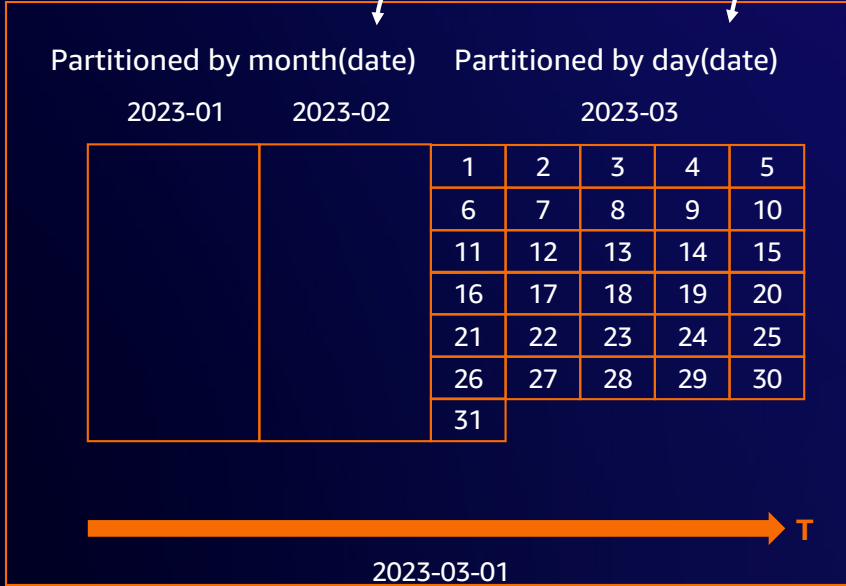
Query example

```
SELECT * FROM sales
WHERE
  date > 2023-02-11 AND
  date < 2023-03-18
```

Split plan 1

Split plan 2

Table sales



#2 Schema & Partition Evolution

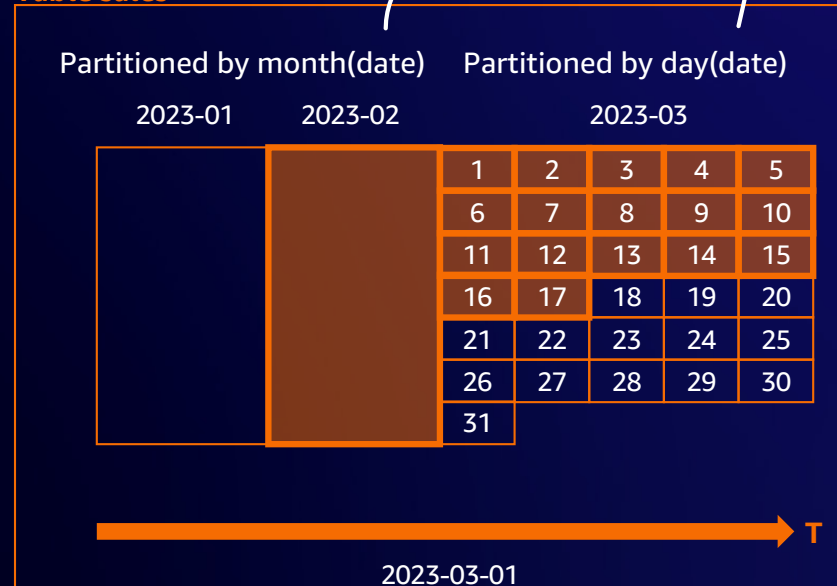
Query example

```
SELECT * FROM sales
WHERE
  date > 2023-02-11 AND
  date < 2023-03-18
```

Read data

Read data

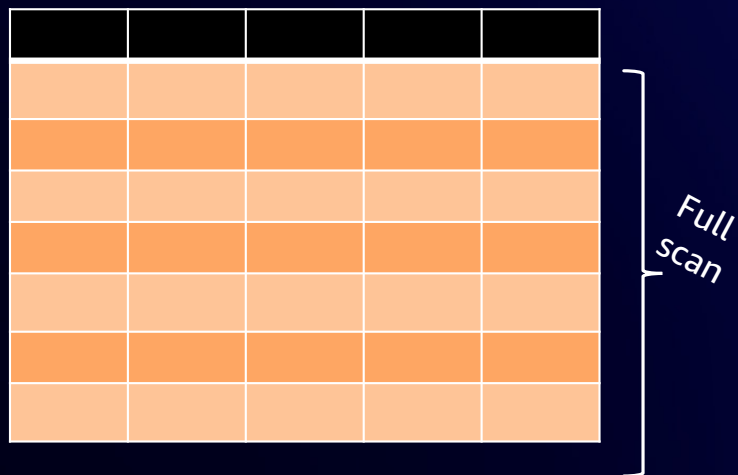
Table sales



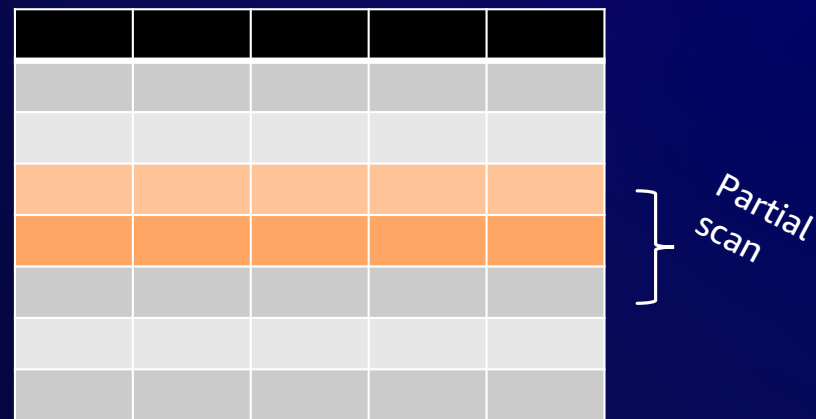
#3 Simplified Queries with Hidden Partitioning

Example of an *Hive* query

```
SELECT level, Count(1) as count
FROM logs
WHERE event_time BETWEEN
      '2018-12-01 10:00:00'
AND   '2018-12-01 10:10:00'
```



```
SELECT level, Count(1) as count
FROM logs
WHERE event_time BETWEEN
      '2018-12-01 10:00:00'
AND   '2018-12-01 10:10:00'
AND   event_date = '2018-12-01'
```



#3 Simplified Queries with Hidden Partitioning

Tables are created using "Partition Transforms"

```
CREATE TABLE catalog.db.logs (  
  id long,  
  log_event string,  
  level string,  
  event_time timestamp  
) PARTITIONED BY (day(event_time)) USING iceberg;
```

Example of an Iceberg query

```
SELECT level, Count(1) as count  
FROM logs  
WHERE event_time BETWEEN  
      '2018-12-01 10:00:00'  
AND   '2018-12-01 10:10:00'
```



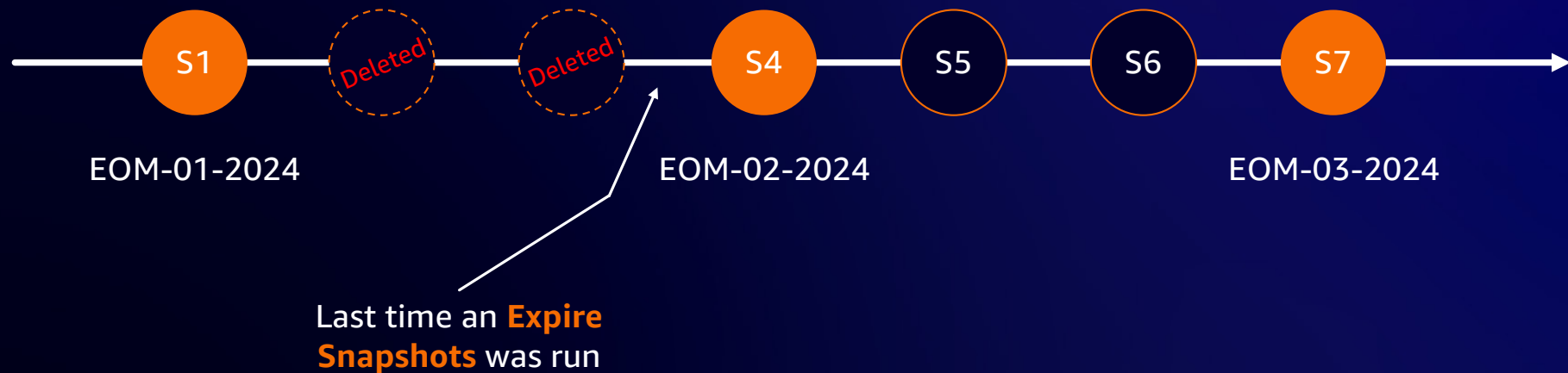
#4 Tagging and Branching for Compliance

Each transaction (insert/delete/update/compaction)
creates a new table snapshot



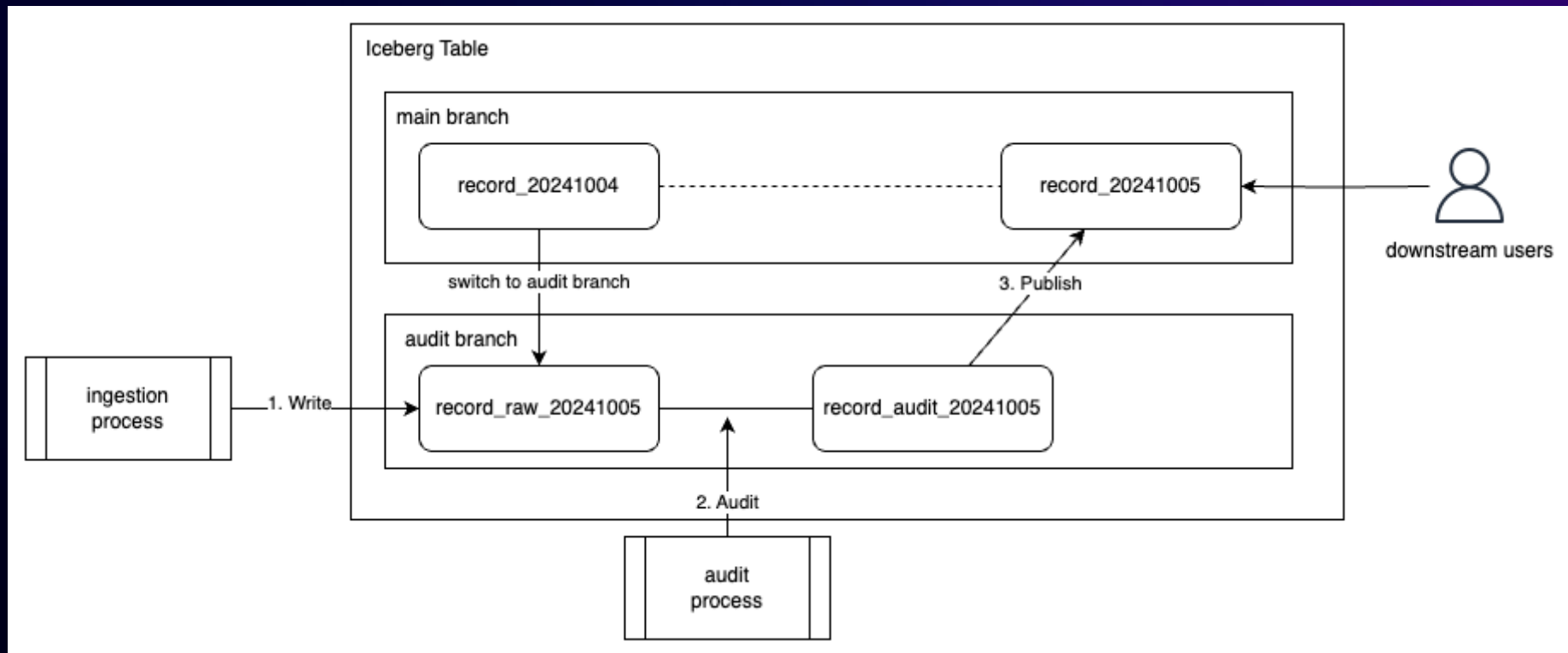
#4 Tagging and Branching for Compliance

Keep only the most relevant Table Snapshots for the business



#4 Tagging and Branching for DataOPS

Write – Audit- Publish (WAP) Pattern



Apache Iceberg Table Update Strategy



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Tables Update Strategies

Copy-On-Write



Merge-on-Read

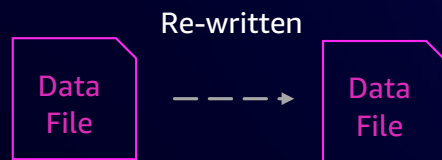


`write.update.mode` | `write.delete.mode` | `write.merge.mode`



Tables Update Strategies

Copy-On-Write – COW



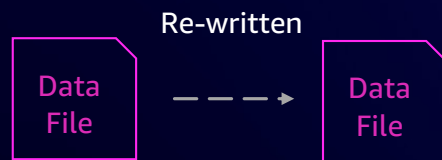
Writer/Reader Trade-off

Cost may be higher
for writes

No impact on the
readers side

Tables Update Strategies

Copy-On-Write – COW

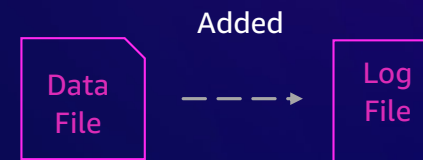


Writer/Reader
Trade-off

Cost may be higher
for writes

No impact on the
readers side

Merge-On-Read - MOR



Writer/Reader
Trade-off

Offloads writing
cost

Adds overhead
for reads

Apache Iceberg Table Maintenance



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Iceberg Table Maintenance

Iceberg tables may need additional maintenance to be “healthy” and “performant”

Iceberg provides a “Control Plane”

- Iceberg provides Procedures implemented in Spark, Flink, or Trino/ Athena → **table maintenance is managed by the user**
- AWS Glue Data Catalog provides a serverless experience for Iceberg tables → **table maintenance is managed by AWS**
- Fully managed Iceberg maintenance on a storage level using **S3 Tables**



Iceberg Table Maintenance

Most important actions

- **Compactions** – contains multiples operations:
 - combine small objects together
 - sort the data using hierarchal sorting
 - clustering the data using z-ordering
- **Expirations** – removes "old" metadata
- **Orphan reclamation** – removes unreferenced objects

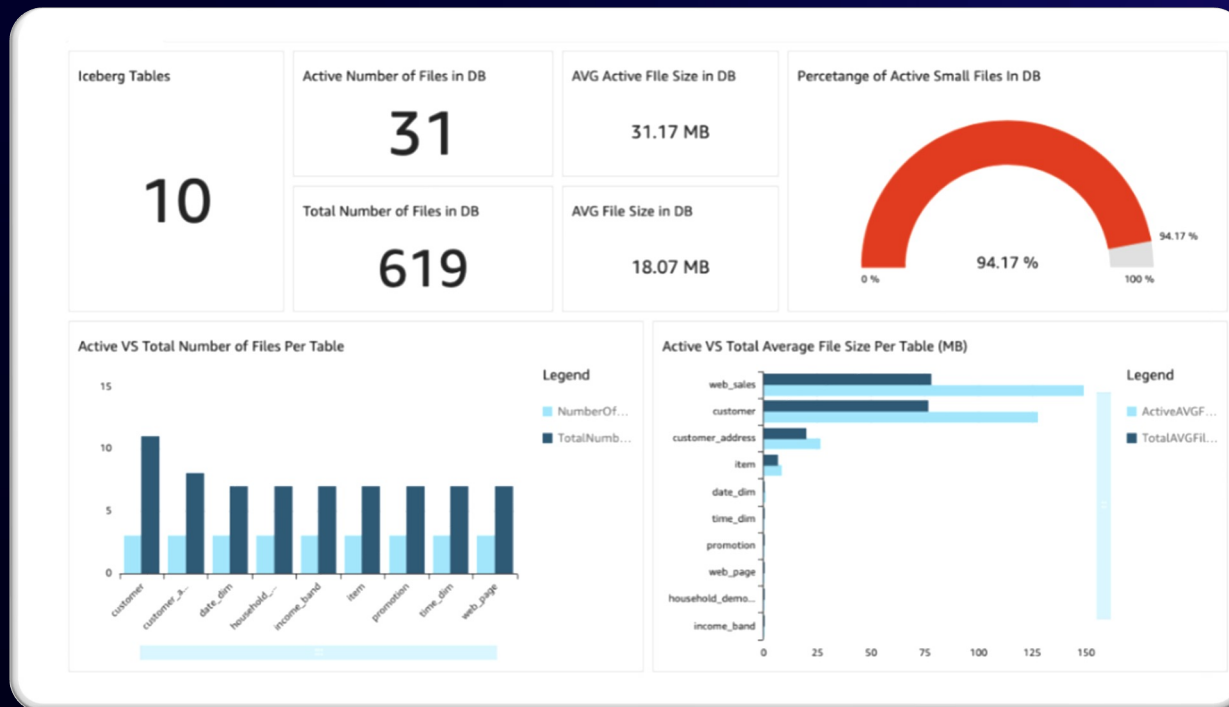
Maintaining Iceberg Tables

Monitor Iceberg table with metadata tables



Maintaining Iceberg Tables

Monitor Iceberg table with metadata tables



Key maintenance tasks:

- Rewrite Files (Data File Optimization)
- Rewrite Manifests (Metadata Optimization)
- Delete orphan files (Cleanup)

NEW

Table maintenance by S3

PERFORMANCE



Compaction: Consolidate small objects into larger ones to improve query performance.

Snapshot Retention: Remove unused snapshots





Key Takeaways

- Active Open-Source Community
- AWS Prioritizes Apache Iceberg
- Production-grade maturity

Migrate an existing data lake to a transactional data lake using Apache Iceberg



Amazon S3 Tables integration with Amazon SageMaker Lakehouse is now generally available



Thank you!

David Greenshtein
dgreensh@amazon.de



© 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.