# AWS re:Invent

**DECEMBER 2 – 6, 2024 | LAS VEGAS, NV**

STG370

# Simplify data management with Amazon S3 Tables

**David Lee**
Principal Product Manager,
Amazon S3
AWS

**Prathiban Mohanasundaram**
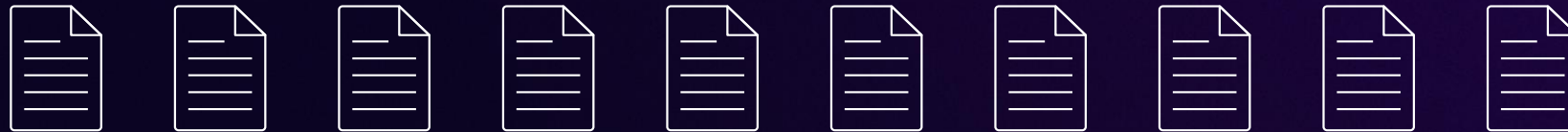Senior Software Development Manager,
Amazon S3
AWS

aws

# Agenda

**01**    Introduction to Parquet/Iceberg

**02**    Introducing Amazon S3 Tables

**03**    Data management with S3 Tables

**04**    S3 Tables demo

**05**    Q&A

# Today, Amazon S3 is also a *tabular* data store

10+ exabytes of Parquet data stored

Servicing 15 million requests per second

Transmitting hundreds of petabytes every day
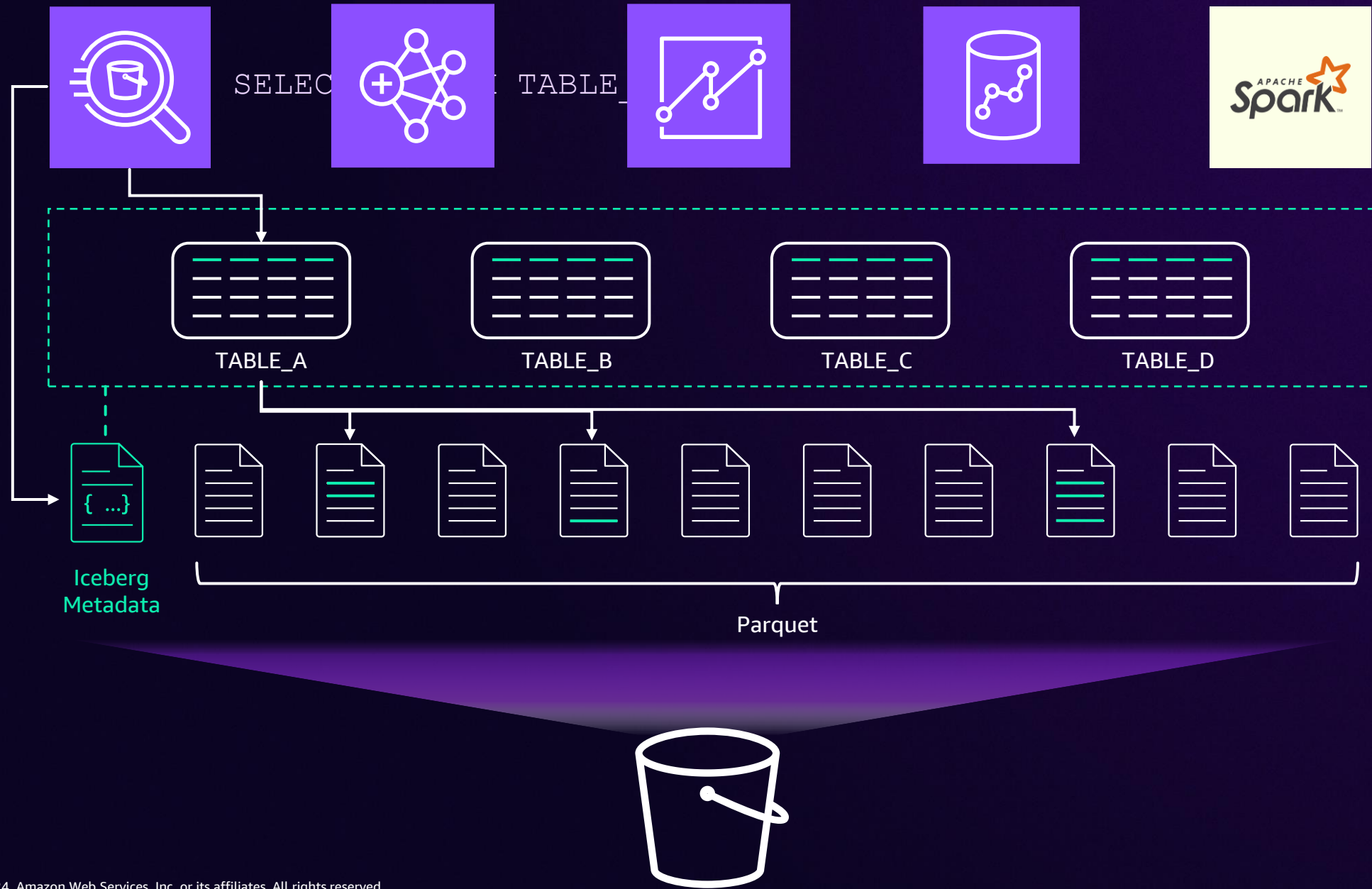
Parquet

# What is Apache Parquet?



Open source

Columnar data file format

Optimized for performance

SELEC... ... TABLE_...

TABLE_A        TABLE_B        TABLE_C        TABLE_D

Iceberg
Metadata

Parquet

# Customer problems to solve

Growing scale requires more and more performance

Enforcing table-level security and integrity is complex

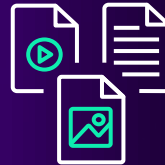Optimizing storage cost drives unexpected operational burden

**GA**    **Dec. 3, 2024**

# S3 Tables

Fully Managed Apache Iceberg Tables in S3

**Improved query performance** based on storage tuning and optimized data layout

**Simplified table security** and integrity controls

**Automated storage cost optimization** based on snapshot management and garbage collection

# Enhanced scalability and performance

**10x**

Transactions per second
(TPS)

**3x**

Improvement to
query performance

# Enhanced scalability and performance

**10x**

Transactions per second
(TPS)

Optimized key naming and layout

Amazon S3 tuned specifically for Iceberg workloads

Enables a higher starting point for S3 TPS scaling

*Learn more about Iceberg TPS scaling with Amazon S3 here:*

Enhanced scalability and performance

© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.

# Enhanced scalability and performance

**Automatic compaction of underlying Parquet files**

Fewer requests to S3

Better throughput

Better tail latencies

# Enhanced scalability and performance

**3x**

Improvement to query performance

# Seamless Integration



API

S3 Tables Catalog Implementation
for Apache Iceberg

SQL

Amazon
EMR

# Simplified security

The table is a first class AWS resource!

- Has an ARN

- Can take an Amazon S3 resource policy

- Has a dedicated endpoint:

   `s3tables.region.amazonaws.com`

# Fully managed

STORAGE COST OPTIMIZATION

Nightly maintenance runs

- Snapshot expiration
- Garbage collection

# Fully managed – 100% policy-driven table maintenance

Table management policy configuration examples

## Garbage collection

```
aws s3tables put-table-bucket-maintenance-configuration \
--table-bucket-arn "arn:aws:s3tables:us-east-2:4236238:bucket/customer-sales-prod \
--type icebergUnreferencedFileRemoval \
--value '{"status":"enabled",\
"settings":{"icebergUnreferencedFileRemoval":{"unreferencedDays":1, "nonCurrentDays":1}}}'
```

# Fully managed – 100% policy-driven table maintenance

Table management policy configuration examples

## Snapshot management

```
aws s3tables put-table-maintenance-configuration \
--table-bucket-arn "arn:aws:s3tables:us-east-2:423623854866:bucket/customer-sales-prod" \
--namespace customer-retail-sales \
--name customer-media-sales-table \
--type icebergSnapshotManagement \
--value '{"status":"enabled","settings": \
{"icebergSnapshotManagement":{"minSnapshotsToKeep":1, "maxSnapshotAgeHours":1}}}'
```

# Fully managed – 100% policy-driven table maintenance

Table management policy configuration examples

Compaction management

```
aws s3tables put-table-maintenance-configuration \
--table-bucket-arn "arn:aws:s3tables:us-east-2:423623854866:bucket/customer-sales-prod" \
--namespace customer-retail-sales \
--name customer-media-sales-table \
--type icebergCompaction \
--value '{"status":"enabled","settings":{"icebergCompaction":{"targetFileSizeMB":128}}}'
```

# S3 Table bucket APIs

Table operations

```
S3tables:ListTable
S3tables:CreateTable
S3tables:GetTableMetadataLocation
S3tables:UpdateTableMetadataLocation
S3tables:DeleteTable
```

Table management

```
S3tables:PutTablePolicy
S3tables:PutTableBucketPolicy

S3tables:PutTableMaintenanceConfig
S3tables:PutTableBucketMaintenanceConfig
```

# Simplify data management with Amazon S3 Tables



Improved performance

Simplified security

Seamless integration

Fully managed

# S3 Tables demo

# Create a table bucket

# Create a table and managing it using Amazon EMR

# Create a table and managing it using Amazon EMR

# Querying tables using Athena

# Q&A

# Thank you!

Please complete the session survey in the mobile app

**David Lee**

Principal Product Manager

Amazon S3

AWS

**Prathiban Mohanasundaram**

Senior Software Development Manager

Amazon S3

AWS