**Pachu**
**Staff Software Engineer**

Starburst

# Data ingestion is just a fact of life

Starburst

# Ingestion tailored for Iceberg data lakes

Starburst

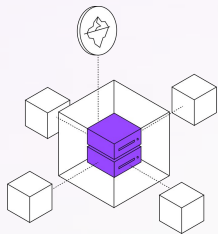# Starburst

**Starburst is a data platform** built on a data lakehouse architecture with Trino + Apache Iceberg and designed to accelerate analytic workflows from development to deployment.
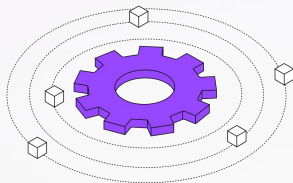
Starburst

**For an end-to-end data warehouse experience, five core capabilities must exist**

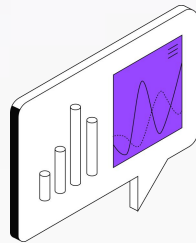# Starburst Galaxy Platform

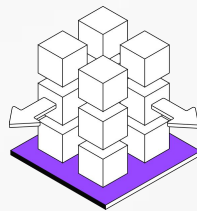## Data Ingestion

How data lands in the lake

## Query Engine

How data is processed for data pipelines and analyzed for consumption
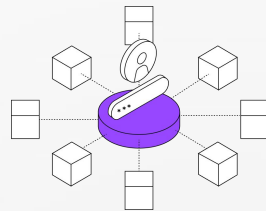
## Data Management

How tables are cleaned and optimized to maintain performance

## Automatic Capacity Management

How compute scales up and down to match demand

## Data Governance

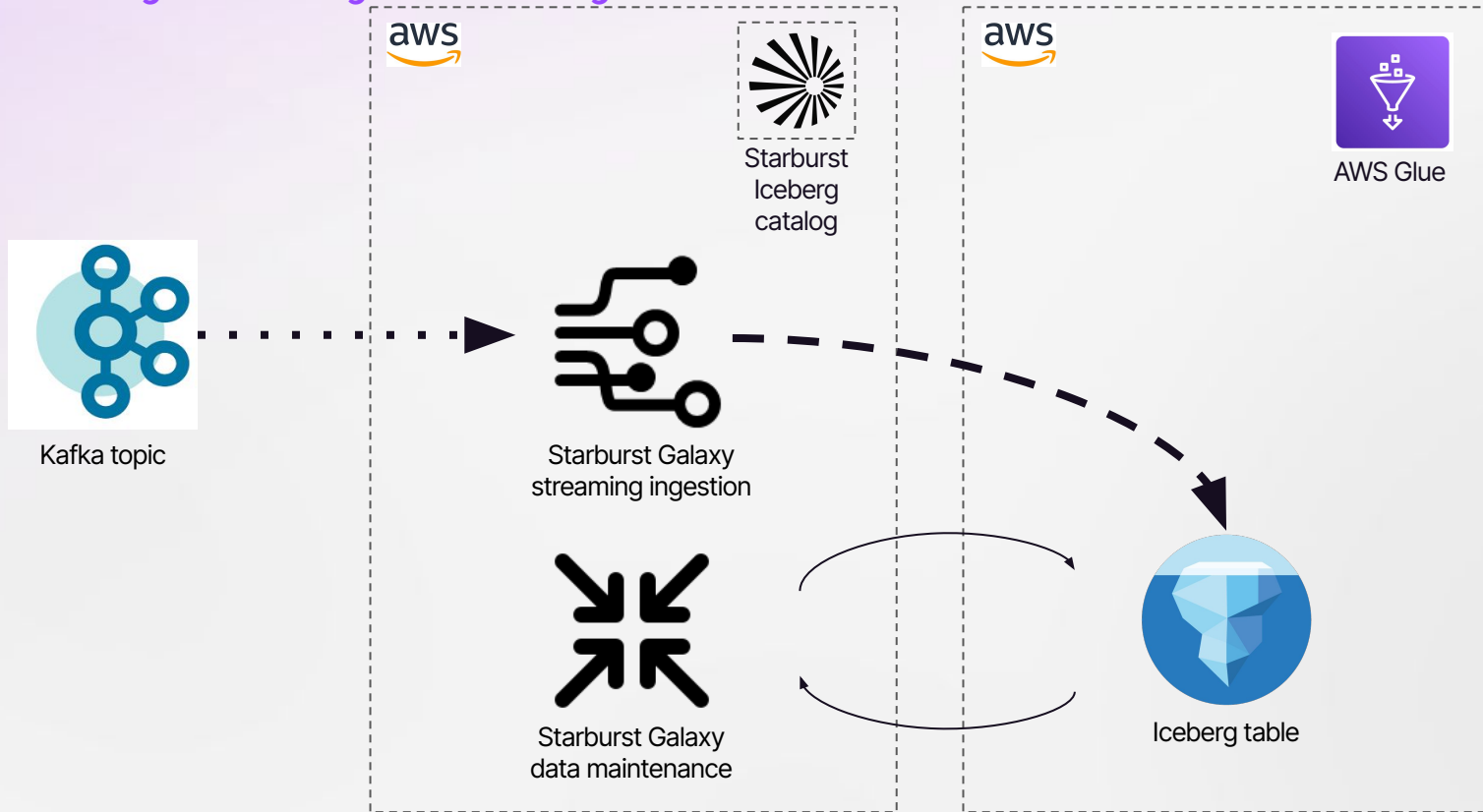How data is accessed and secured once in the lake

✺ Starburst

# Simplified Kafka-to-Iceberg ingestion

*Starburst ingestion: designed for Iceberg*

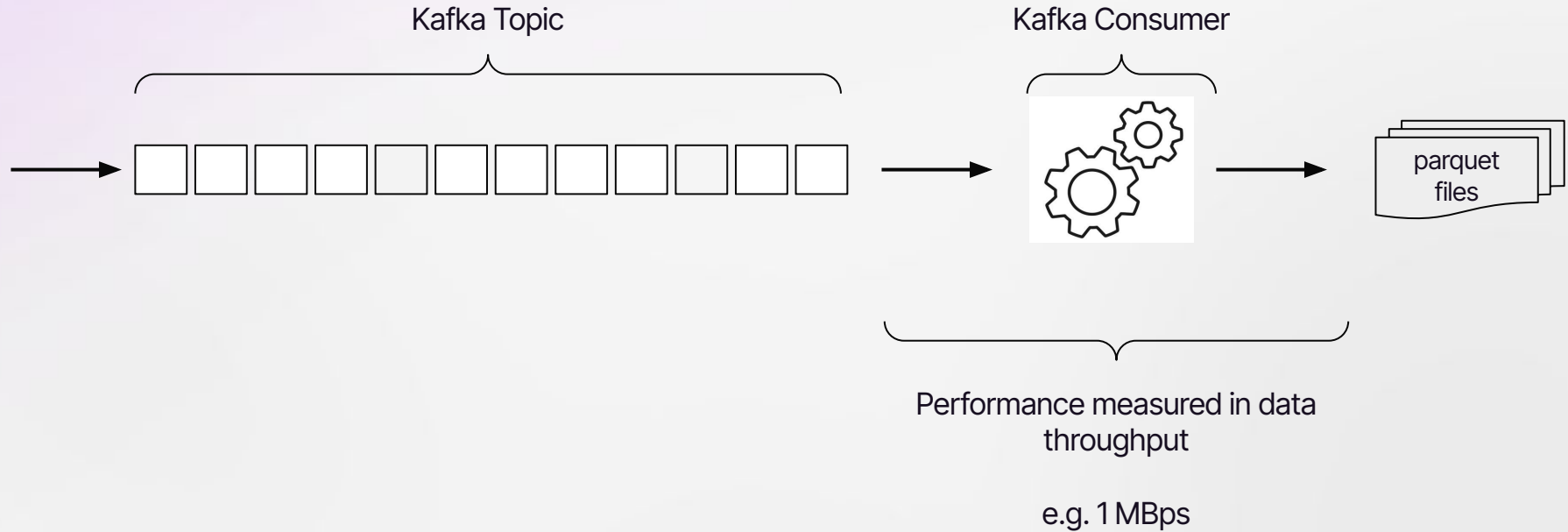# Unlocking Iceberg streaming ingestion

Building for Scale

Planning for Imperfection

Unified Data Platform

Starburst

# Kafka scaling 101

Kafka Topic

Kafka Consumer

parquet
files

Performance measured in data
throughput

e.g. 1 MBps

Starburst

# Kafka partitioning increases throughput

Kafka Topic

Kafka Consumers

Kafka Partitions

P1

P2

P3

parquet files

parquet files

parquet files

Performance measured in total data throughput
e.g. 1 MBps / partition * 3 partitions = 3 MBps

Starburst

# Are we done with data scaling concerns?

Starburst

# Challenge - parallelism vs compaction

Small files kill downstream query performance

# Challenge - parallelism vs compaction

# Too few consumers causes backlog

Starburst

# Challenge - parallelism vs compaction



Files fully compacted on write

Increasing data buffering time can increase file size, but also adds latency

Kafka partitions fully parallelized

Backlogged

File size

Parallelism

Not enough parallelism causes backlogging

# Challenge - parallelism vs compaction

How did we solve this?

# Dynamic load coordination

Starburst

# Starburst Solution - Dynamic load coordination

# Starburst Solution - Dynamic load coordination

# Optimizes query performance on near real-time data

Starburst

# Then there's Iceberg commit contention problem

Starburst

# Challenge - commit contention

# Challenge - commit contention



Kafka Topic

Kafka Consumers

Kafka Partitions

P1

Failed commit

P2

Failed commit

P30

metadata file

metadata file

manifest list

manifest list

manifest file

manifest file

manifest file

parquet files

parquet files

parquet files

Starburst

How did we solve this?

# Iceberg commit coordination

Starburst

# Starburst Solution - Iceberg commit coordination

**Proven internet scale throughput**

# Tested up to 100 GB per second being written to a single Iceberg table.

Starburst

# Unlocking Iceberg streaming ingestion

**Building for Scale**

**Planning for Imperfection**

**Unified Data Platform**

☀ Starburst

# When things go wrong
## *(and they will)*

Starburst

# Invalid data

# Challenge - dealing with invalid data

# Challenge - dealing with invalid data



Kafka Topic

Kafka Consumer

ETL Pipeline

???

Invalid records

parquet files

Raw / Bronze

parquet files

Structured / Silver

parquet files

Reporting / Gold

Starburst

# Starburst solution - transactional dead lettering

# Misconfigured Schema

# Challenge - Misconfigured schema

**Configured Schema**

```
"side" -> VARCHAR
"symbol" -> VARCHAR
"quantity" -> INTEGER
"price" -> BIGINT
"account" -> BIGINT
```

**Actual message**

```
{
    "side": "SELL",
    "symbol": "SNOW",
    "quantity": 1453,
    "price": 25.10,        Wrong type
    "account": 4567,
    "stop_loss": 26.00     New column
}                          omitted
```

Kafka Topic

Kafka Consumer

???

Mismatched type appears

New field appears

Starburst

# Starburst Solution - Error detection, classification, and notification



Kafka Topic

Starburst Streaming Ingestion

Mismatched type appears

New field appears

parquet files

parquet files

Dead Letter / Errors Table

Notification Service

Starburst Galaxy Web UI

Open Standards

*(Upcoming)*

OpenTelemetry

Monitoring Tools

*(Upcoming)*

DATADOG

Starburst

# Starburst solution - reset and replay

## 1. Adjust Schema

### Original Schema

```
"side" -> VARCHAR
"symbol" -> VARCHAR
"quantity" -> INTEGER
"price" -> BIGINT
"account" -> BIGINT
```

### New Schema

```
"side" -> VARCHAR
"symbol" -> VARCHAR
"quantity" -> INTEGER
"price" -> DOUBLE
"account" -> BIGINT
"stop_loss" -> DOUBLE
```

## 2. Choose Reset Point

### Backfill options

Would you like to manage backfill for this table ?

○ Apply changes without backfill

● Rewind table to savepoint and backfill

**Select savepoint**

Savepoints available for the last 30 days.

Select date
18 Sep 2024 📅 ?

;

Savepoint *
Sep 18, 2024, 11:30:00 PM ⌄ ?

Cancel    Save

## 3. Automatically Replay Data Ingestion with New Schema

New Timeline    Old Timeline



Starburst

# Demo

Starburst

# Data ingest

## Create your first ingest source

Continuously stream data from Kafka or ingest files from Amazon S3. Learn more about data ingest. 

1. Connect to a source

2. Land your data in a raw table

3. Create transform table and schematize data

**Connect new source**

# Unlocking Iceberg streaming ingestion

**Building for Scale**

**Planning for Imperfection**

**Unified Data Platform**

Starburst

# Disparate tools     vs.     Unified platform

**Disparate tools**

- Integration: seems easy, but headache to debug

- Configure and maintain multiple tools

**Small frictions, but they add up**

**Unified platform**

- One product for ingest, transform, and analytics
- Faster feedback loop from raw data to decisions
- Simplified operations, stronger governance, and happier teams

**Focus on insights, not integration**

≋ Starburst

# Leading price performance

**Benchmark cost study of 8mb/s throughput (*us-east-1*)**

| Kafka Ingestion Solution | Monthly Cost | Annual Cost |
|---|---|---|
| Starburst Streaming Ingest | $47.16 | $565.92 |
| Vendor 1 | $79.20 (1.68x) | $950.40 |
| Vendor 2 | $224.01 (4.75x) | $2,688.12 |
| Vendor 3 | $531.23 (11.26x) | $6,374.76 |

❋ Starburst

# What about other ingestion sources?

Starburst

# Iceberg file loader

## Continuous S3 to Iceberg ingestion along similar principles



Raw files on
AWS S3

Iceberg Tables on
AWS S3

Starburst

# Starburst Galaxy Ingestion - Summary

## Building for Scale

- Dynamic load coordination optimizes query performance

- Commit coordination enables massive scale

## Plan for Imperfection

- Transactional dead lettering
- Error detection, classification, notification
- Reset and replay

## Unified Data Platform

- End-to-end: from ingest to analytics
- Custom architecture yields price performance advantages

✳ Starburst

Sign up for our newsletter
and enter to win a limited-edition
Starburst community T-shirt

Join your local
Meetup Group

Try Galaxy Free
$500 CREDIT

Catalogs

Each catalog contains configuration for Starburst Galaxy to access a data source. Configure catalogs and use them in clusters to query data sources in Starburst Galaxy.

Create catalog

60 catalogs    Search catalogs

| Name | Index status | Kind | Description | Cloud | Region | Tags |
|------|--------------|------|-------------|-------|--------|------|
| amit_json | ✓ | | - | aws | US East (Ohio) | No tags assigned. |
| az_lakehouse | ✓ | | Galaxy metastore with Azure Pay-As-You-Go Blob Storage for Iceberg | A | (US) East US | No tags assigned. |