# Rewriting History

**Migrating petabytes of data to Apache Iceberg using Trino**

**Marc Laforet**
Senior Data Engineer
@ Shopify

# whoami

📍 Toronto, Canada

📊 Worked in data space for ~ 6 years
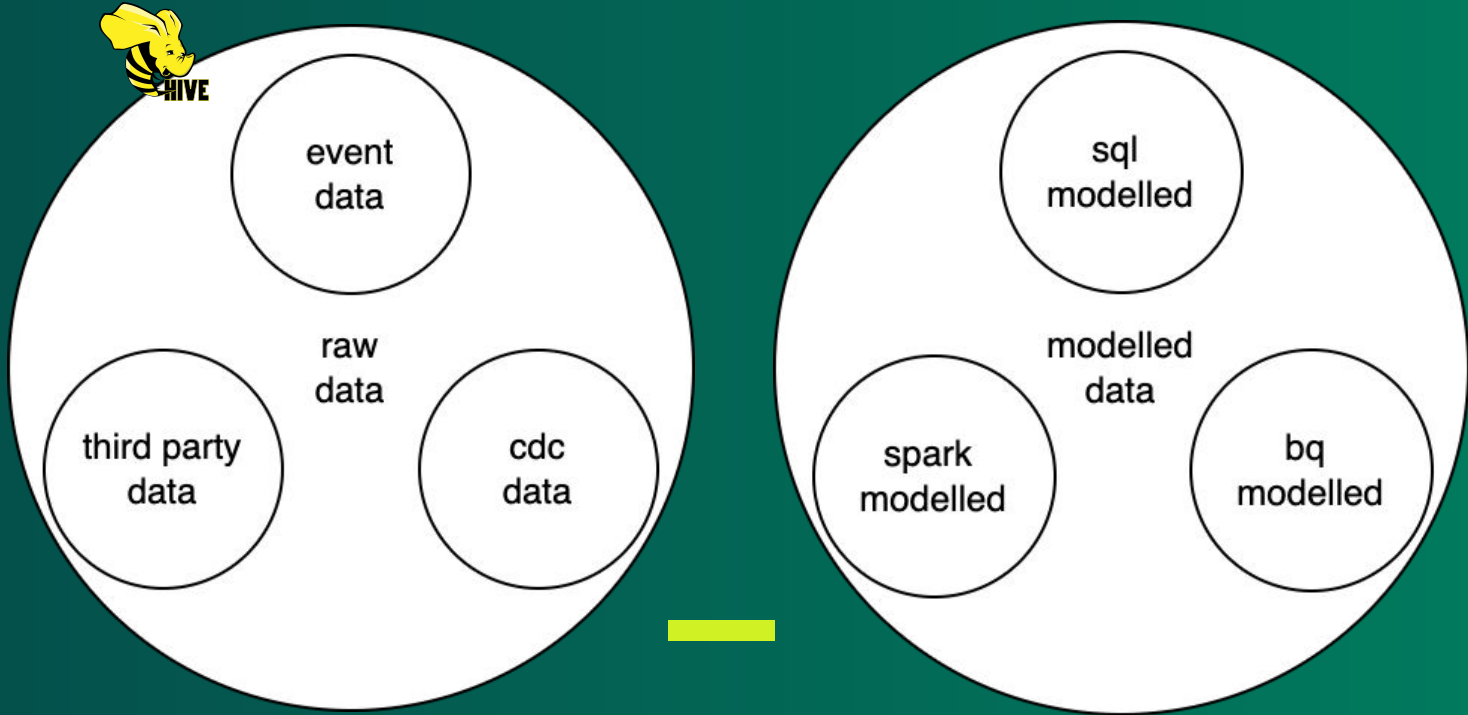
🏊 Swim in my free time

# Agenda

# Agenda

# Lakehouse Team

**Ingest 1st party data**

---

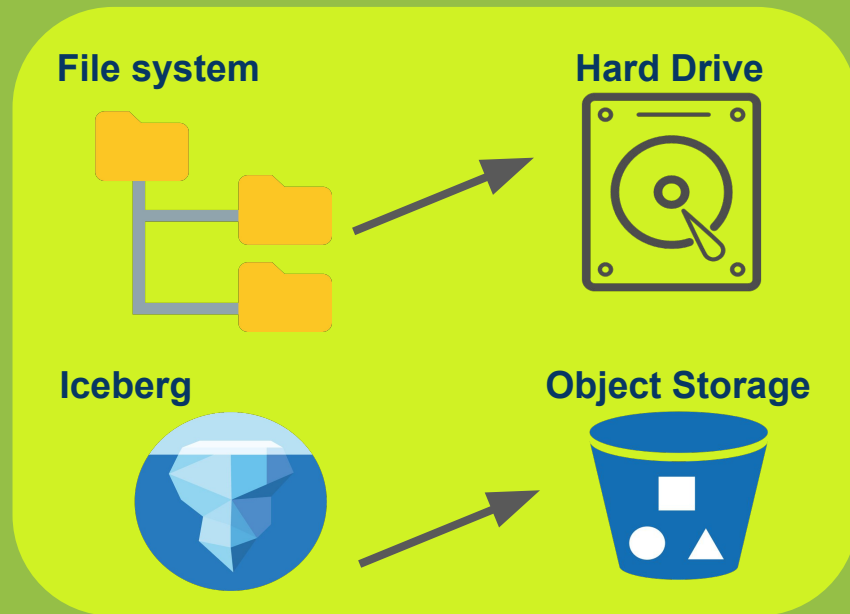**Storage and retrieval of 1st party data**
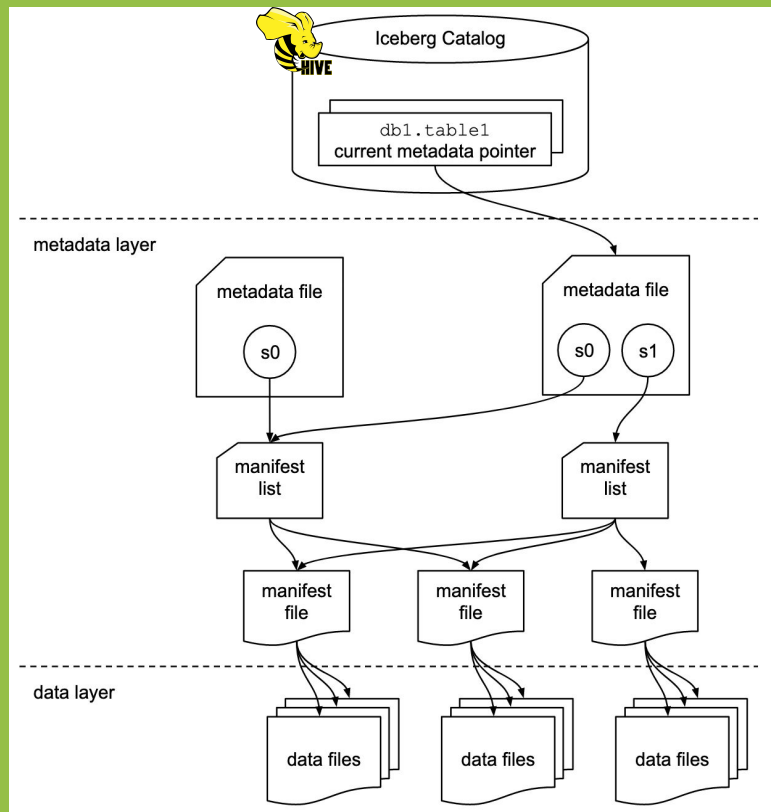
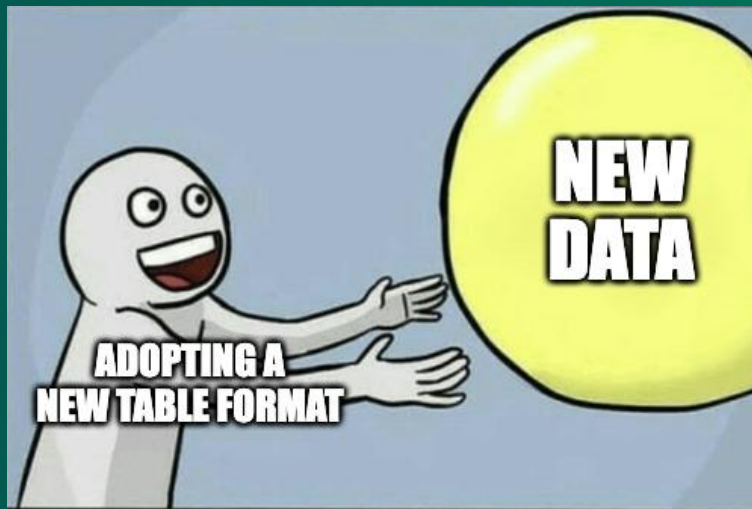# Lakehouse Datasets

# Lakehouse Architecture

# Open Table Formats
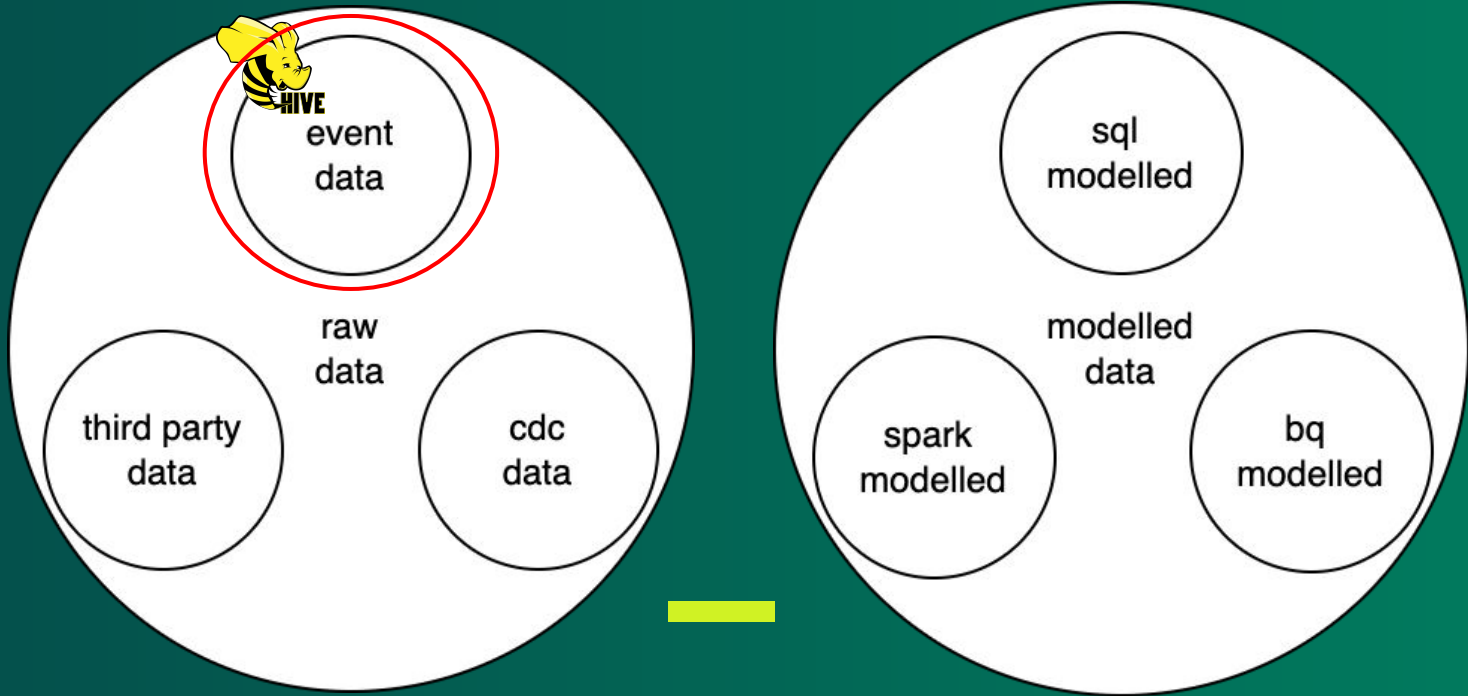
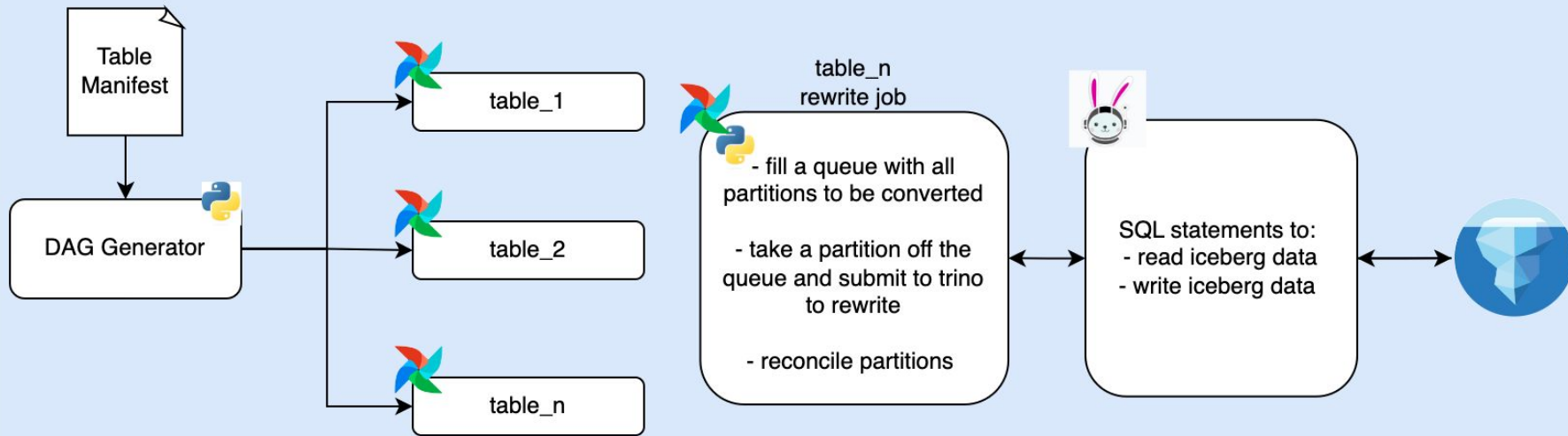Libraries that keep track of files making up a table

# Iceberg

# Lakehouse Datasets

# Migration Strategies

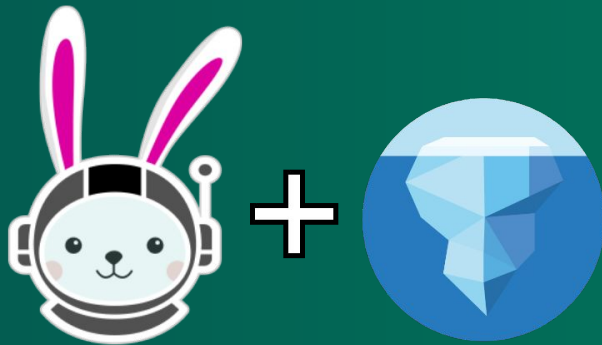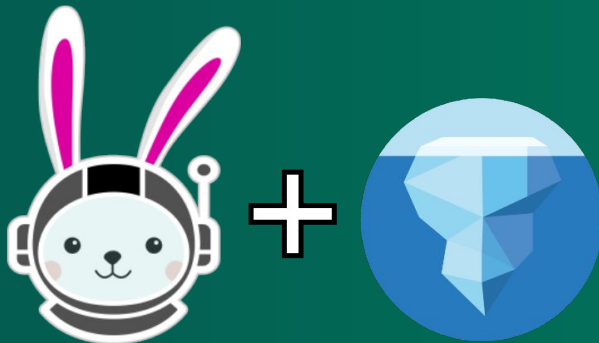| Current State | Strategy |
| --- | --- |
| Data is not in a supported file format | Rewrite everything |
| Data is in a supported file format but you want to change the partioning strategy | Rewrite everything |
| Data is in a supported file format and properly partitioned | Generate manifest files only |

# System Design

# Agenda

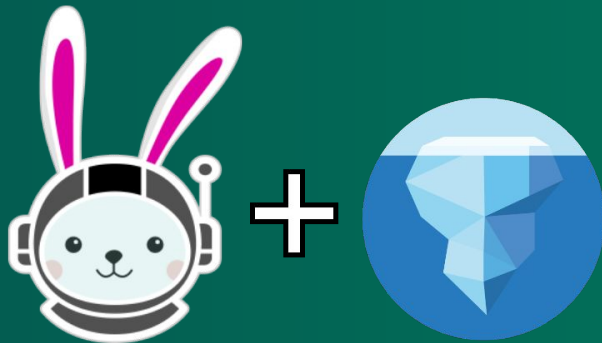# Iceberg tables partitioned on timestamp with time zone

**Unknown type: class io.trino.spi.type.LongTimestampWithTimeZone java.lang.IllegalArgumentException:**
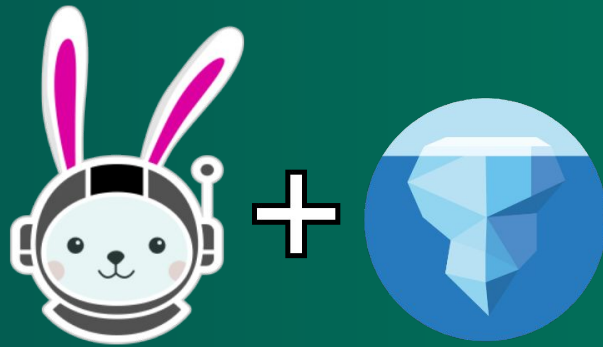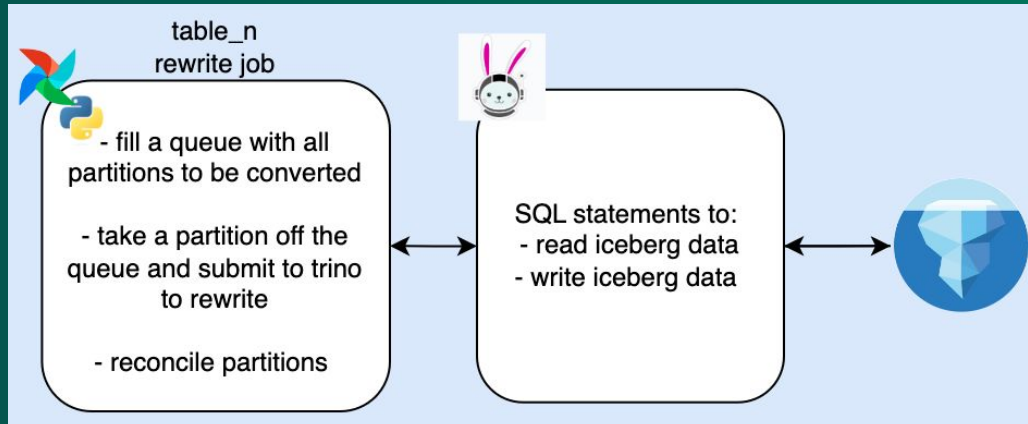Unknown type: class io.trino.spi.type.LongTimestampWithTimeZone

# PR-9757

## Return timestamp time zone partitions in UTC
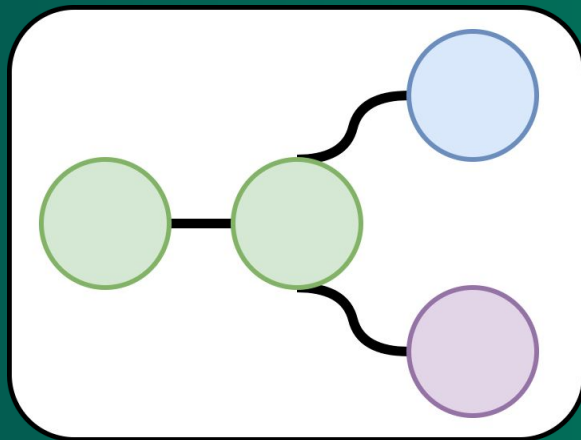
# Scale Iceberg Writes

**org.apache.iceberg.exceptions.CommitFailedException:**
Metadata location
[<bucket-name>/<namespace>/<table-name>/metadata/203942-8936d
432-7ef0-4bfe-8098-bf9ca06580fe.metadata.json] is not same as table
metadata location
[<bucket-name>/<namespace>/<table-name>/metadata/203943-cf4a61
94-1ee7-4483-846c-f495ab237e78.metadata.json] for
<namespace>.<table_name>

**PR-11886**

**Throw proper error to make use of iceberg's commit retries**

# Missing File Problem

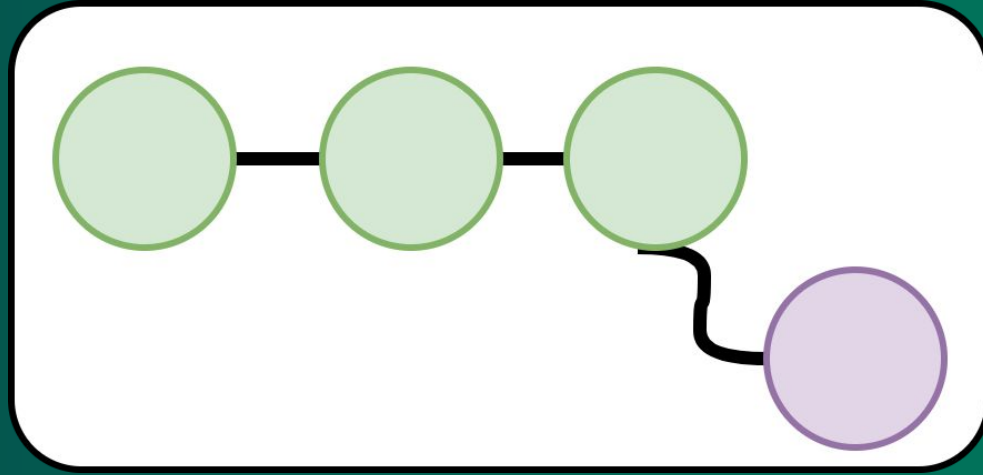**org.apache.iceberg.exceptions.NotFoundException:**
Failed to open input stream for file:
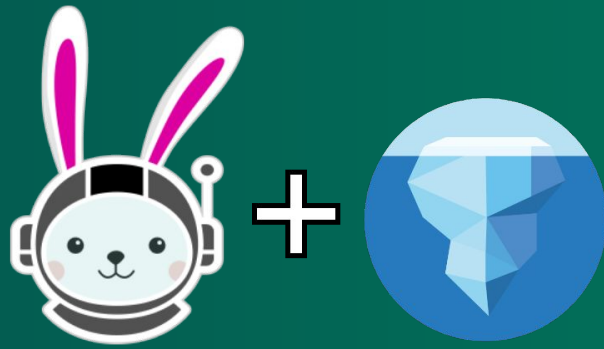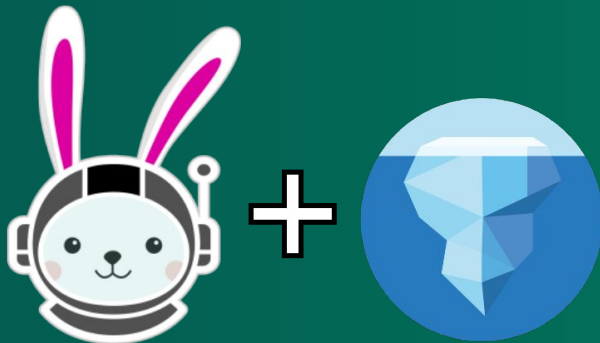[<bucket-name>/<namespace>/<table-name>/metadata/snap-8501054
657342128291-2-8eeb7ee8-54c7-4e08-a2db-f5cbdd855835.avro

**PR-14118**

Avoid pro-active file cleanup on HMS time-out

# Agenda

# Benchmarks

# Benchmarks

## Planning Time



time (ms)

| hive + json.gzip | iceberg + parquet |
|---|---|
| 891.12 | 41.9 |

## Cumulative User Memory

memory_(tb)

| hive + json.gzip | iceberg + parquet |
|---|---|
| 596 | 21.6 |

## Execution Time

time_(min)

| hive + json.gzip | iceberg + parquet |
|---|---|
| 27.99 | 1.84 |

# Benchmarks

## Planning Time

| | hive + json.gzip | iceberg + parquet |
|---|---|---|
| time (ms) | 4250 | 33.58 |

## Cumulative User Memory

| | hive + json.gzip | iceberg + parquet |
|---|---|---|
| memory (tb) | 8190 | 63.5 |

## Execution Time

| | hive + json.gzip | iceberg + parquet |
|---|---|---|
| time (min) | 1838 | 1.78 |

# Reflections & Advice

- Using a modern table and file format can really improve query performance

- Huge benefit using open source software

- Write new and old data to the same table

- Construct a migration priority list

- Airflow offered nice out of the box task scheduling but wasn't great at managing long-running tasks

# Project Contributors

Sam
Wheating

Victoria
Bukta

# Thank you

# The Iceberg Lakehouse

Data Scientist: Hey! How do I use my output from this tool as input to this other tool?

# How many table formats are in use?

## Hive Tables

- Standard hive tables

- Synthetic partition key column

- Columnar data being stored in json.gzip format

## Proprietary Table Formats

- Datasets migrated into bigquery

## Custom Table Formats

- Teams invented their own table formats for their tool

- No ACID properties

- Full drop every time

Dynamic Writer Scaling

https://github.com/trinodb/trino/pull/10614

File pruning using table statistics

https://github.com/trinodb/trino/pull/9326
https://github.com/trinodb/trino/issues/4115

Stitching with a view

https://github.com/trinodb/trino/pull/8621
https://github.com/trinodb/trino/pull/8540

Utilities

https://github.com/trinodb/trino/pull/10480 (vicky)
https://github.com/trinodb/trino/pull/10459 (austen)

Monitoring

https://github.com/trinodb/trino/pull/12026 - update
https://github.com/trinodb/trino/issues/7933 - merge