# Data lineage and observability with OpenLineage
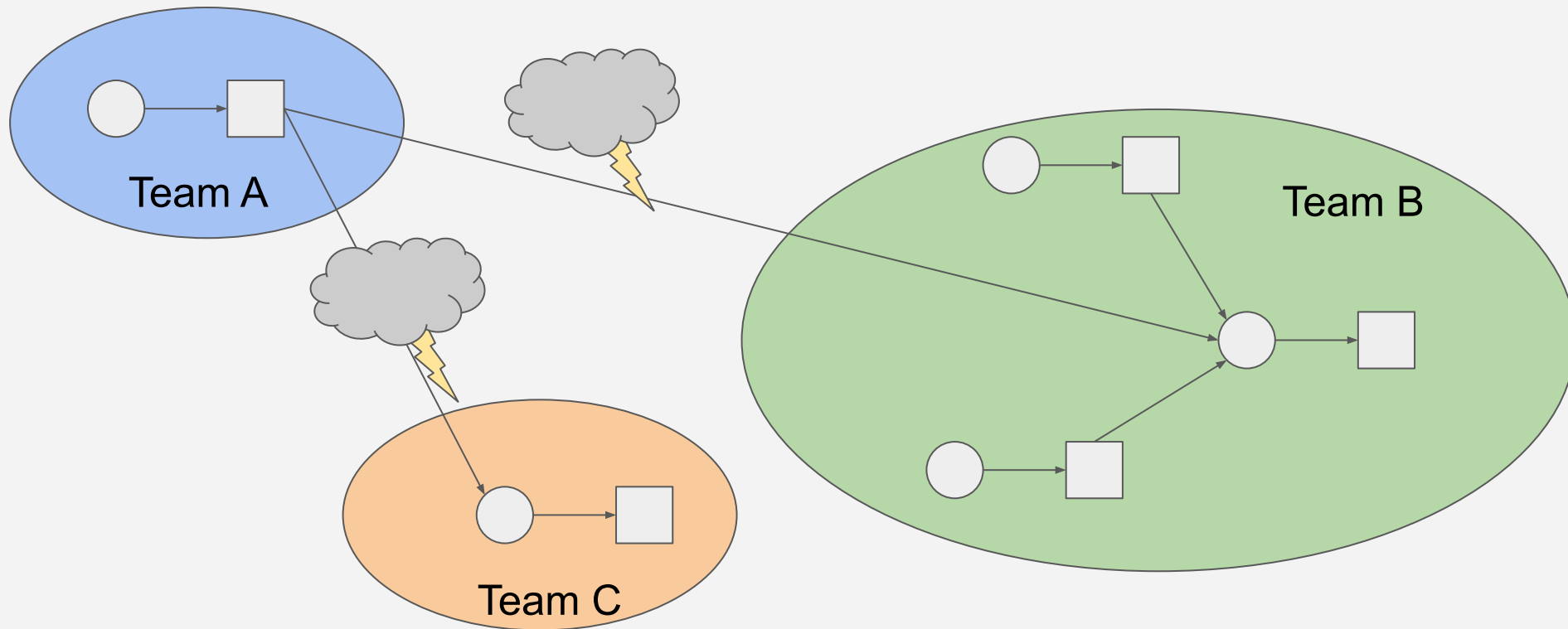
Julien Le Dem, CTO and Co-Founder Datakin | Mai 2021

# AGENDA

- The need for metadata
- **OpenLineage** - the open standard for lineage collection - and **Marquez**, its reference implementation
- Spark observability with OpenLineage

datakin

# The need for Metadata

datakin

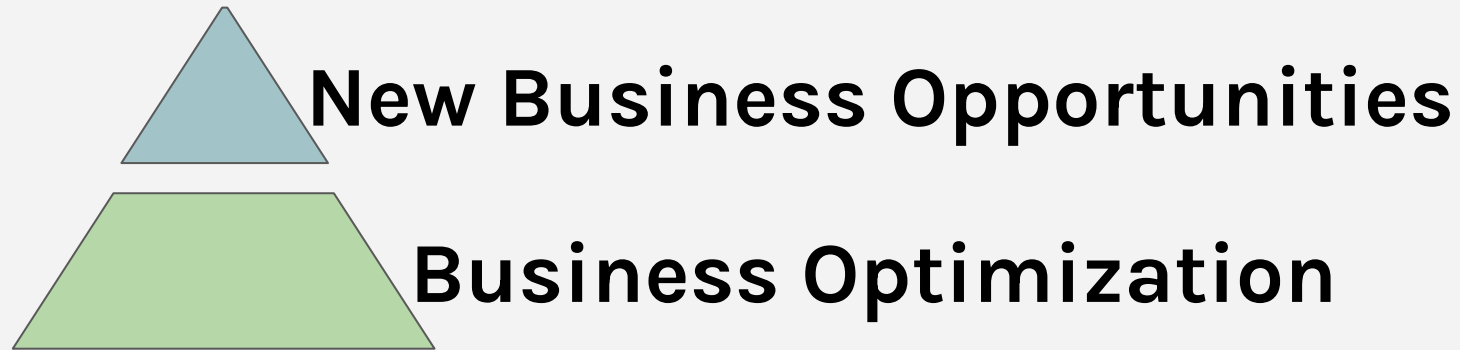# Building a healthy data ecosystem

# Today: Limited context

**DATA**

- What is the data source?
- What is the schema?
- Who is the owner?
- How often is it updated?
- Where is it coming from?
- Who is using the data?
- What has changed?

datakin

# ~~Maslow's~~ Data hierarchy of needs

New Business Opportunities

Business Optimization

Data Quality

Data Freshness

Data Availability

# OpenLineage

# OpenLineage contributors

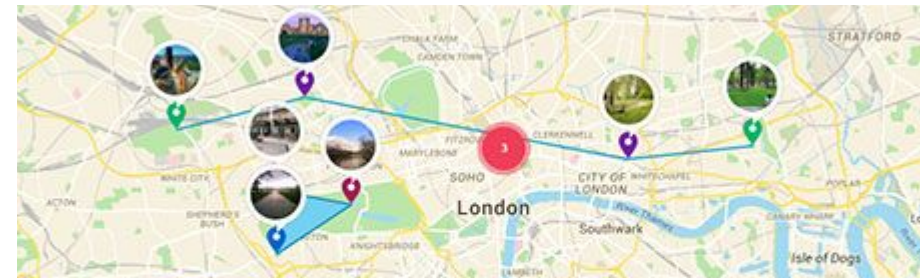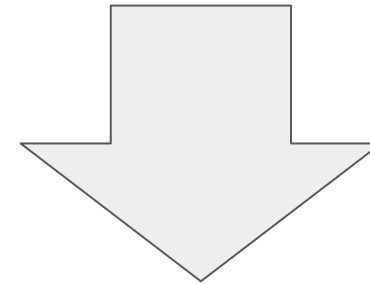Creators and contributors from major open source projects involved

# Purpose:

Define an Open standard for metadata and lineage collection by instrumenting data pipelines as they are running.
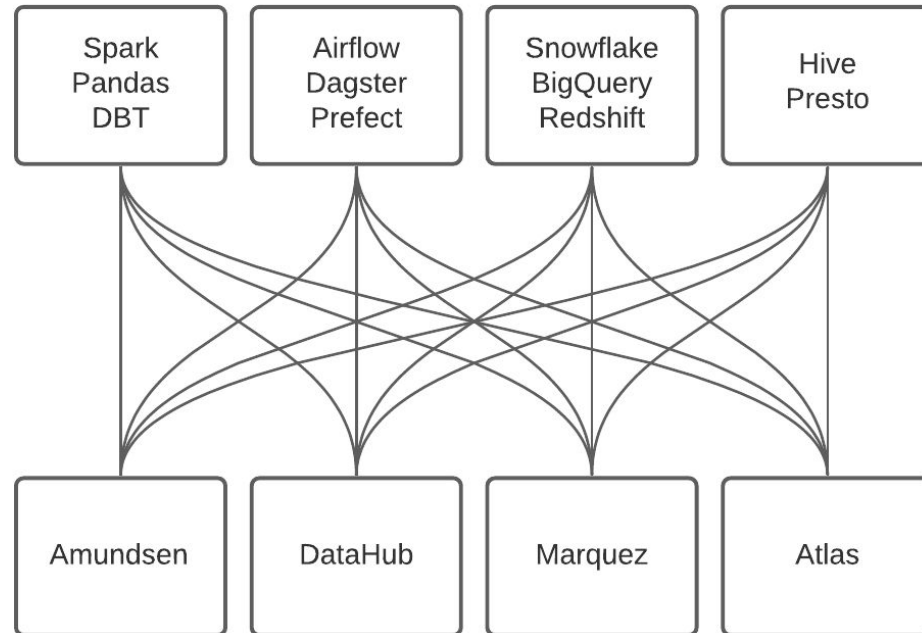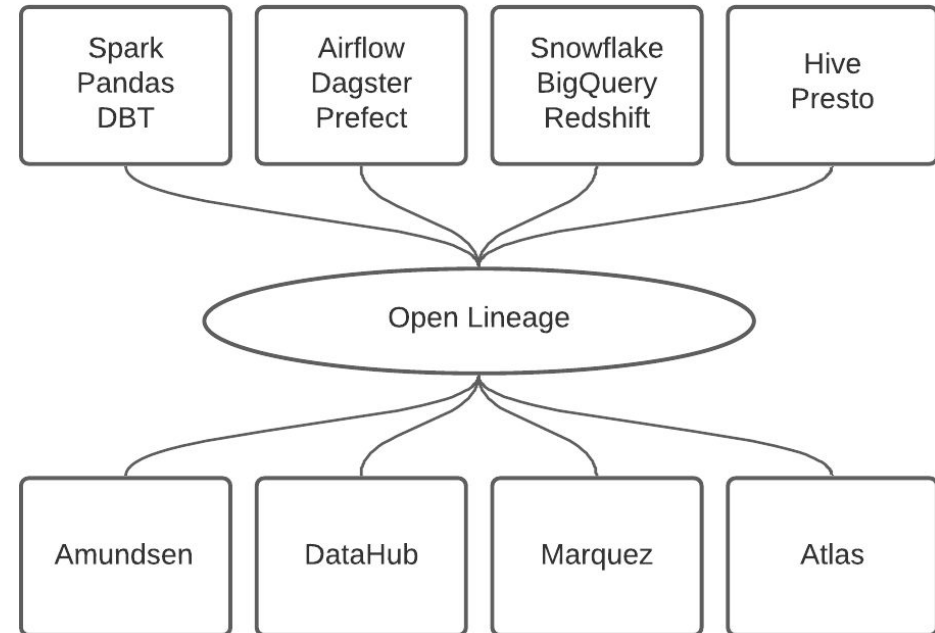
# Purpose:
# EXIF for data pipelines

# Problem

**Before:**

**With Open Lineage**



- Duplication of effort: Each project has to instrument all jobs
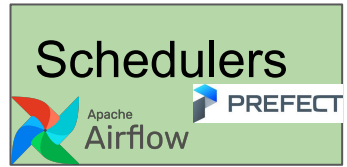- Integrations are external and can break with new versions

- Effort of integration is shared
- Integration can be pushed in each project: no need to play catch up

datakin

# Open Lineage scope | Not in scope

Integrations



Warehouse

Schedulers



...

Metadata and lineage collection standard

Backend

HTTP client

Kafka client

GraphDB client

...

Consumers



Kafka topic



Graph db

# Core Model:
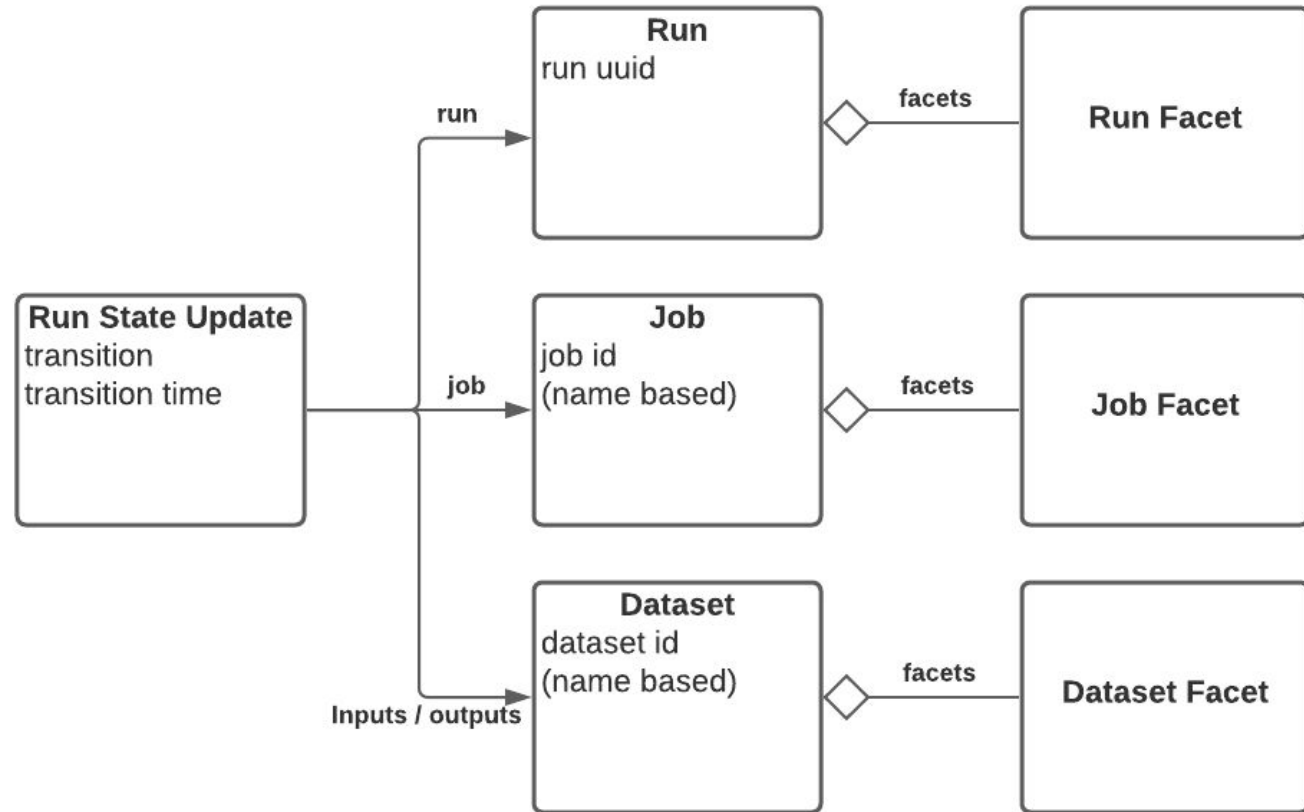
- JSONSchema spec

- Consistent naming:
Jobs:
    scheduler.job.task
Datasets:
    instance.schema.table

# Protocol:

- Asynchronous events:

  Unique run id for identifying a
  run and correlate events

- Configurable backend:
  - Kafka
  - Http

Examples:

- **Run Start event**
  - source code version
  - run parameters


- Run **Complete** event
  - input dataset
  - output dataset version and schema

# Facets

- **Extensible:**

  Facets are atomic pieces of metadata identified by a unique name that can be attached to the core entities.

- **Decentralized:**

  Prefixes in facet names allow the definition of Custom facets that can be promoted to the spec at a later point.
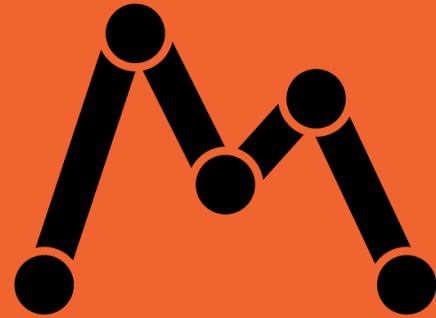
# Facet examples

**Dataset:**
- Stats
- Schema
- Version
- Column level lineage

**Job:**
- Source code
- Dependencies
- params
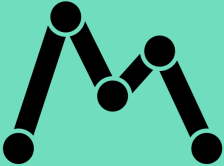- Source control
- Query plan
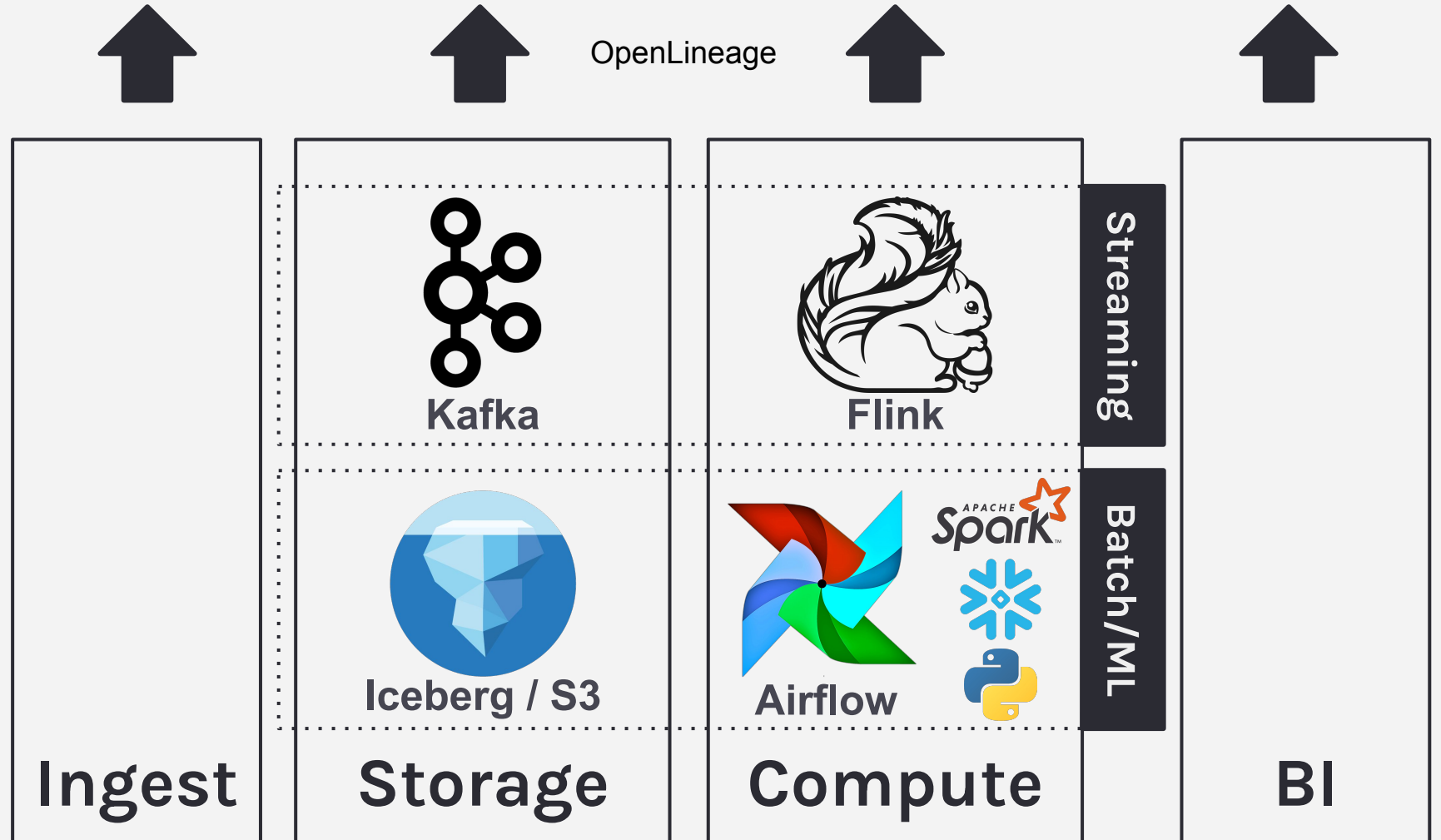- Query profile
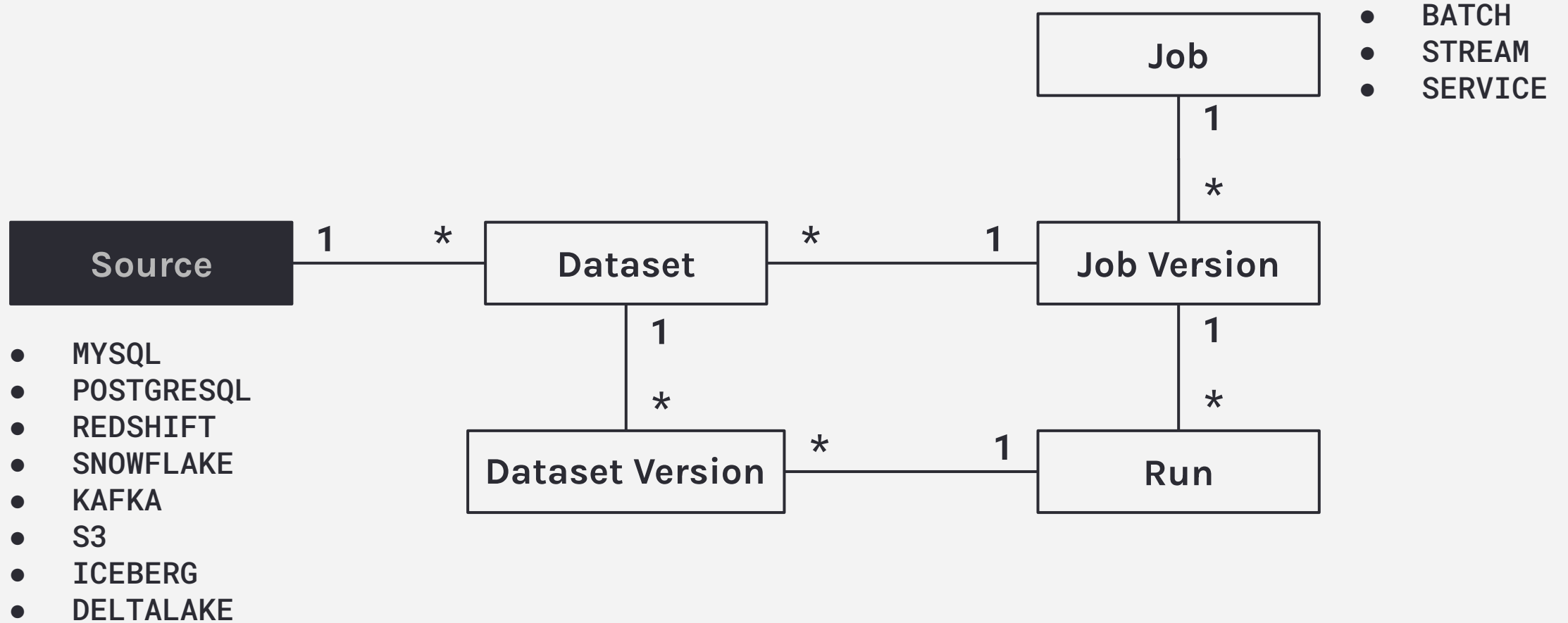
**Run:**
- Schedule time
- Batch id

datakin

MARQUEZ

# Marquez: Data model

# Datakin leverages Marquez metadata



- **Open Lineage and Marquez standardize metadata collection**
  - Job runs
  - Parameters
  - Version
  - Inputs / outputs

- **Datakin enables**
  - Understanding operational dependencies
  - Impact analysis
  - Troubleshooting: What has changed since the last time it worked?

# Spark observability with OpenLineage

datakin

# Spark java agent

**spark.driver.extraJavaOptions:**

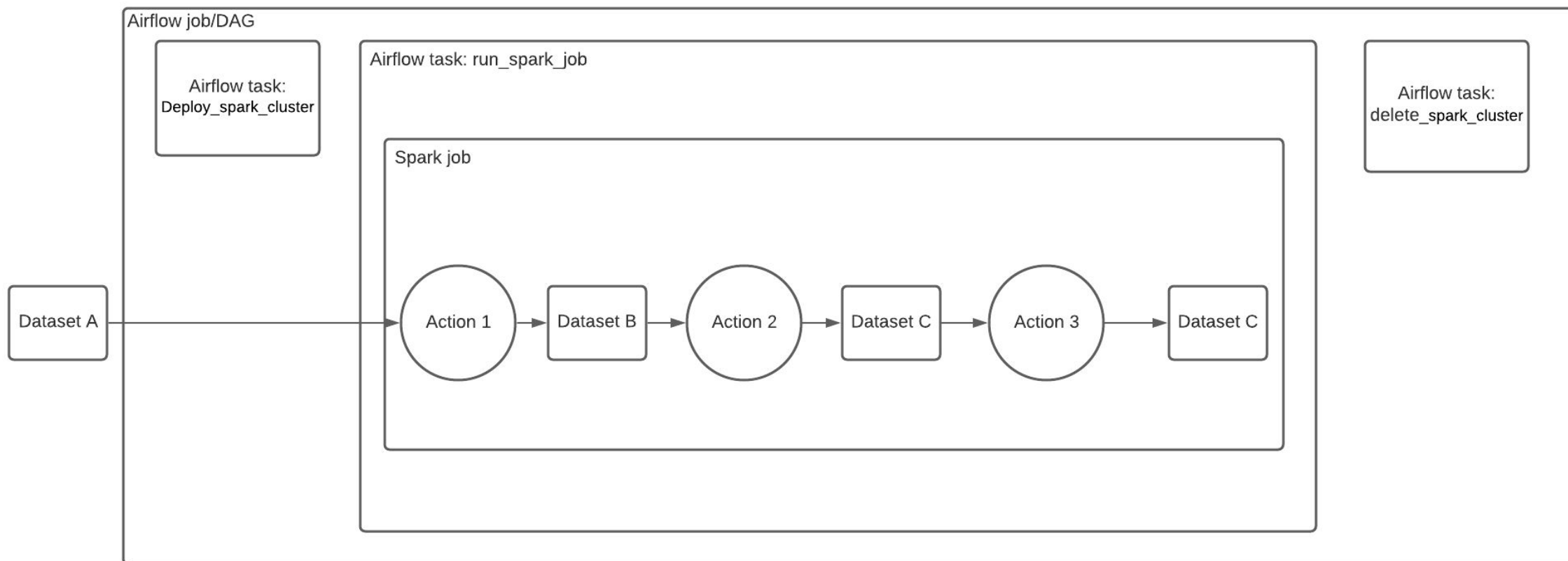-javaagent:marquez-spark-0.13.1.jar={argument}

# Metadata collected

**Lineage:** inputs/outputs

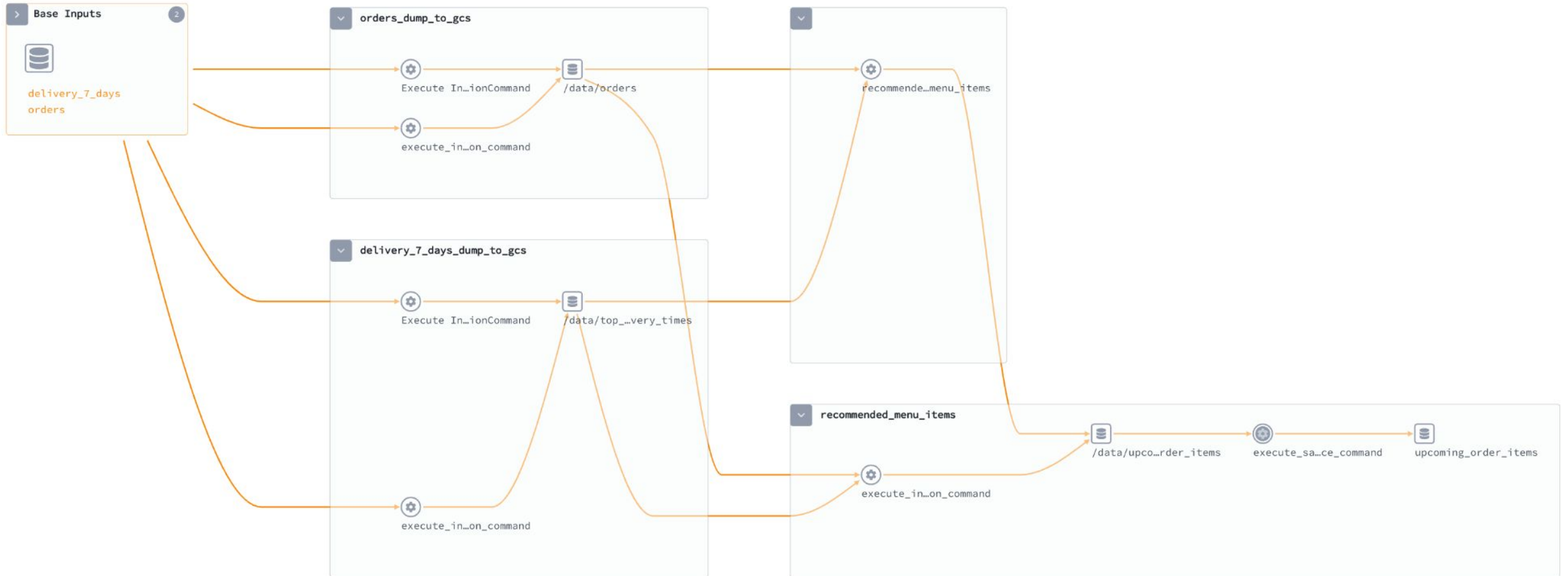**Data volume:** row count/byte size

**Logical plan**
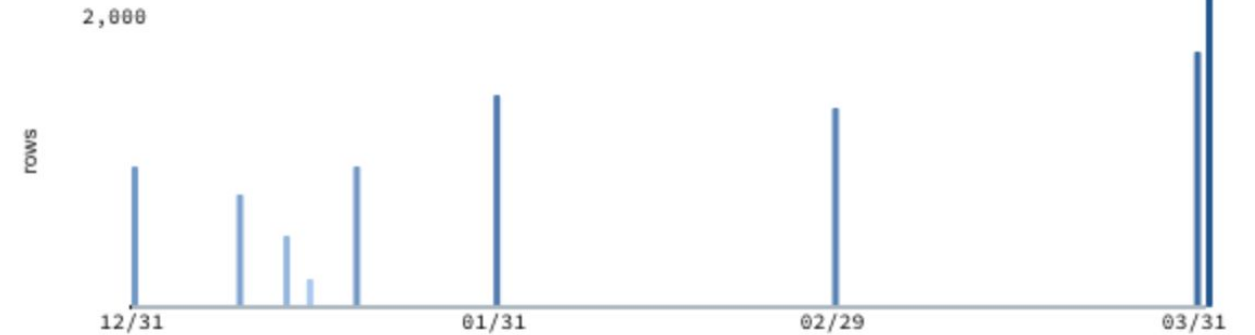
# Lineage model
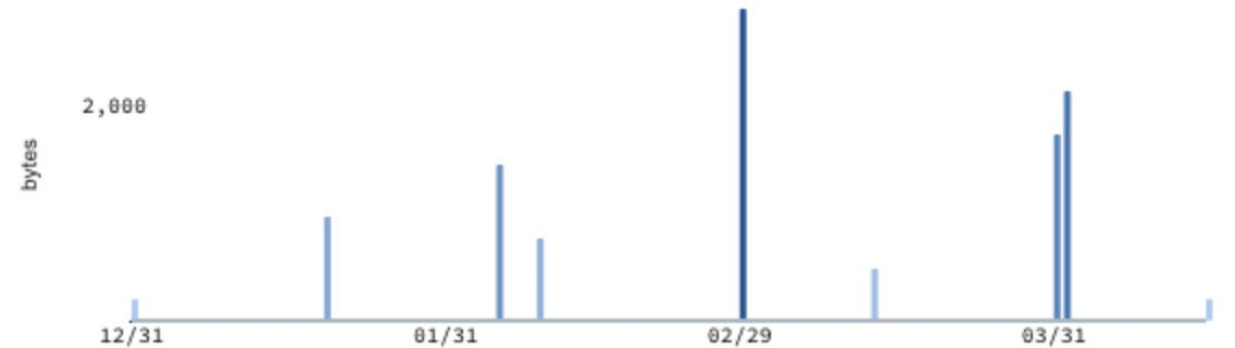
# Lineage Example across jobs

# Example of OpenLineage metadata usage:

# Data volume evolution



Row count

Byte count

# Join the conversation

**OpenLineage:**

    **Github:** github.com/OpenLineage ★

    **Slack:** OpenLineage.slack.com

    **Twitter:** @OpenLineage 🐦

    **Email:** groups.google.com/g/openlineage

**Marquez:**

    **Github:** github.com/MarquezProject/marquez ★

    **Slack:** MarquezProject.slack.com

    **Twitter:** @MarquezProject 🐦

datakin