

# Data lineage and observability with OpenLineage



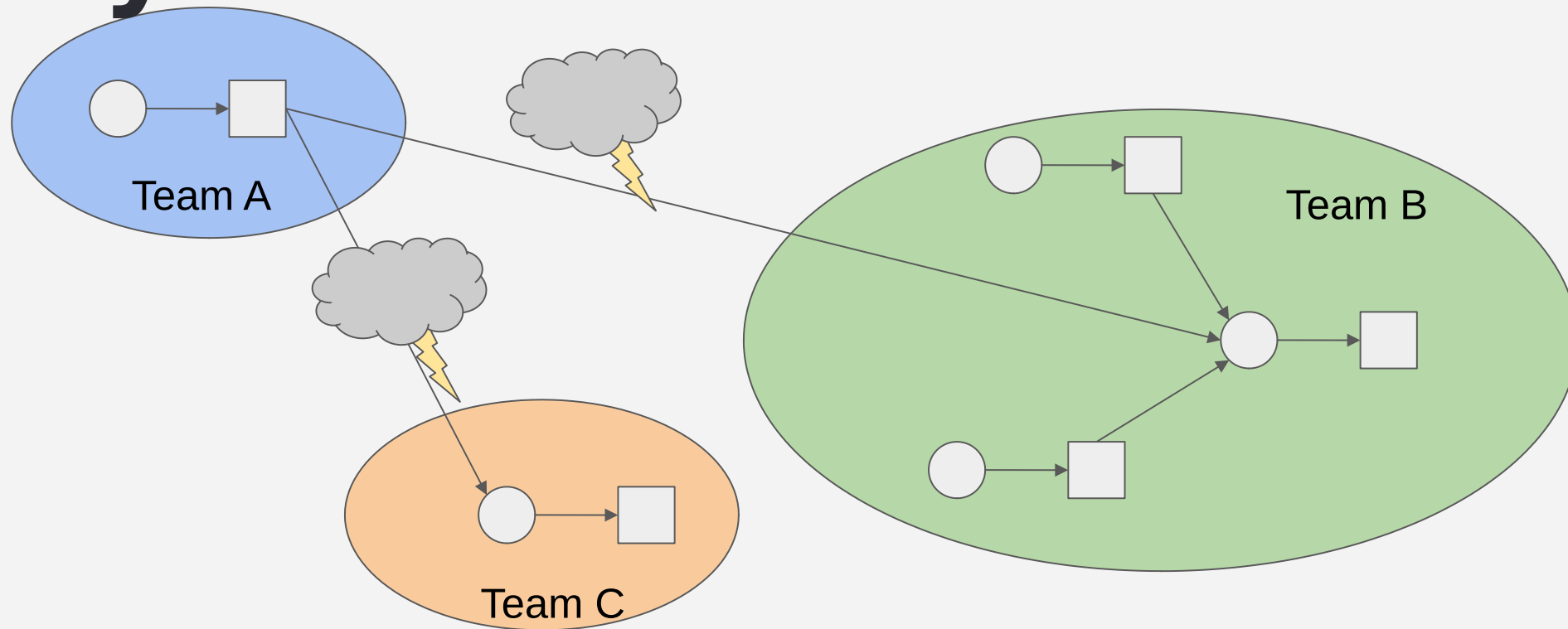
Julien Le Dem, CTO and Co-Founder Datakin | June 2021

# AGENDA

- The need for metadata
- **OpenLineage** - the open standard for lineage collection - and **Marquez**, its reference implementation
- Data observability in practice

The need for Metadata

# 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?



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



**Data Quality**

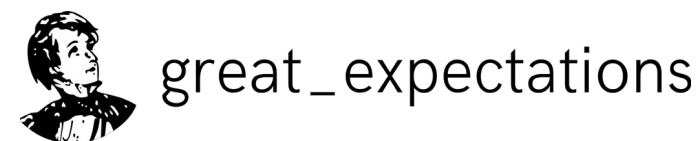
**Data Freshness**

**Data Availability**

OpenLineage

# OpenLineage contributors

Creators and contributors from major open source projects involved



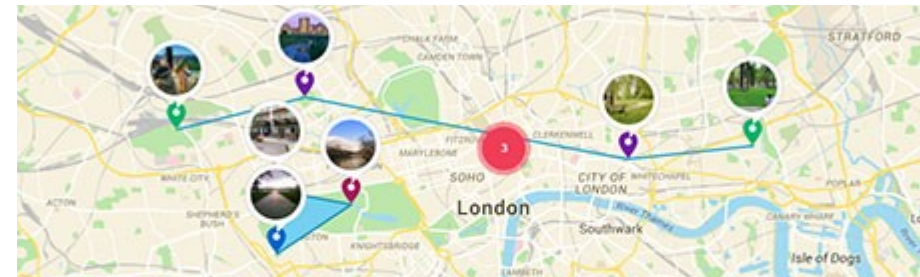
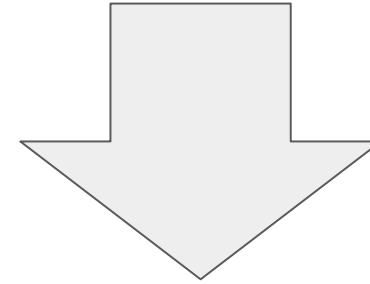


## Purpose:

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

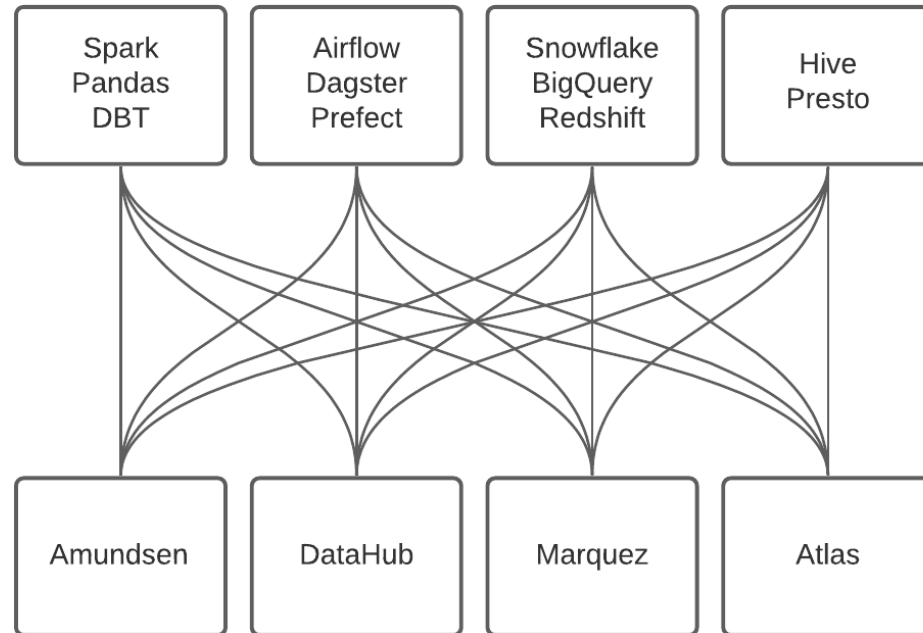


# Purpose: EXIF for data pipelines



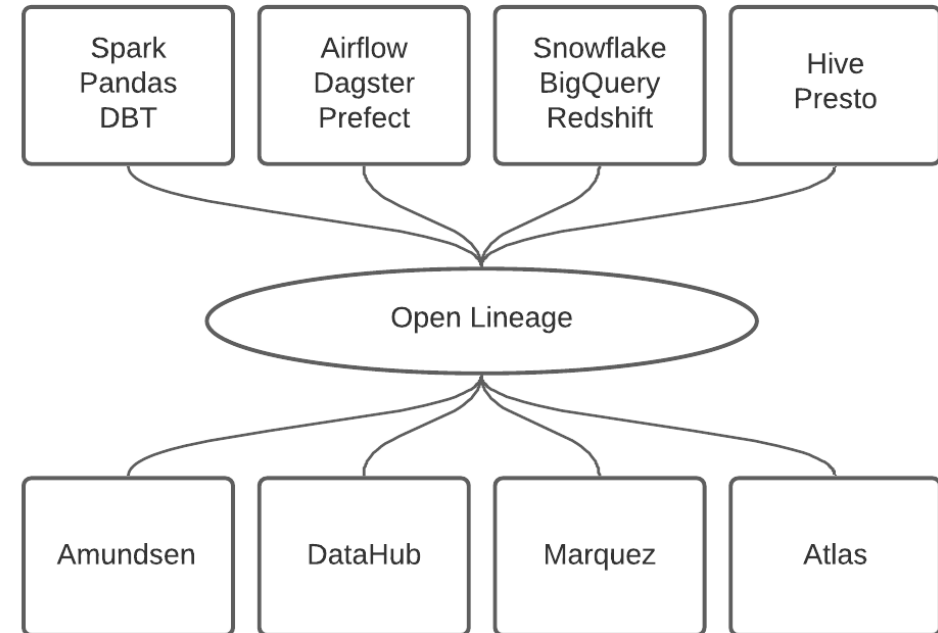
# Problem

## Before:



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

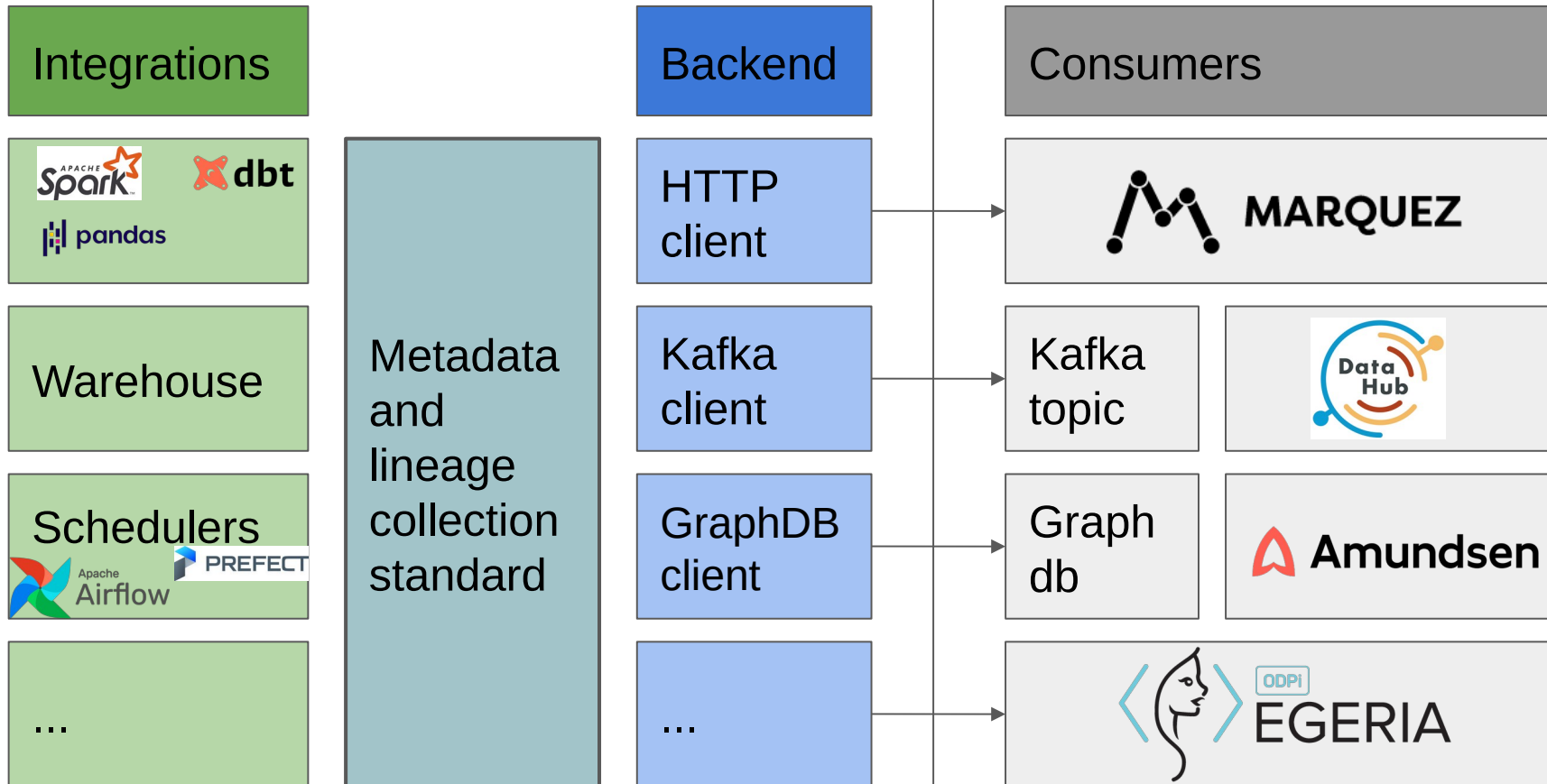
## With Open Lineage



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

# Open Lineage scope

# Not in scope



# Core Model:

JSONSchema spec

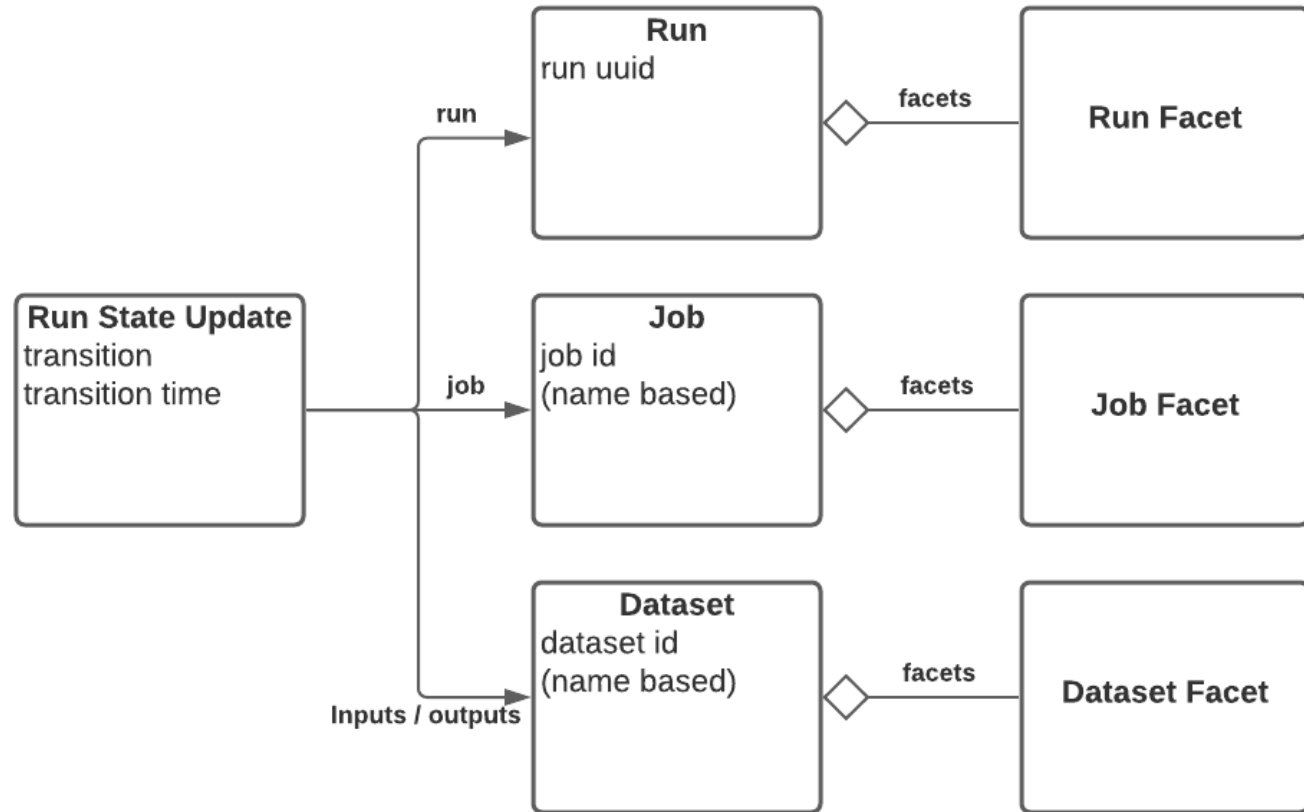
Consistent naming:

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
- Source control
- Query plan

## **Run:**

- Schedule time
- Batch id
- Query profile
- Parameters



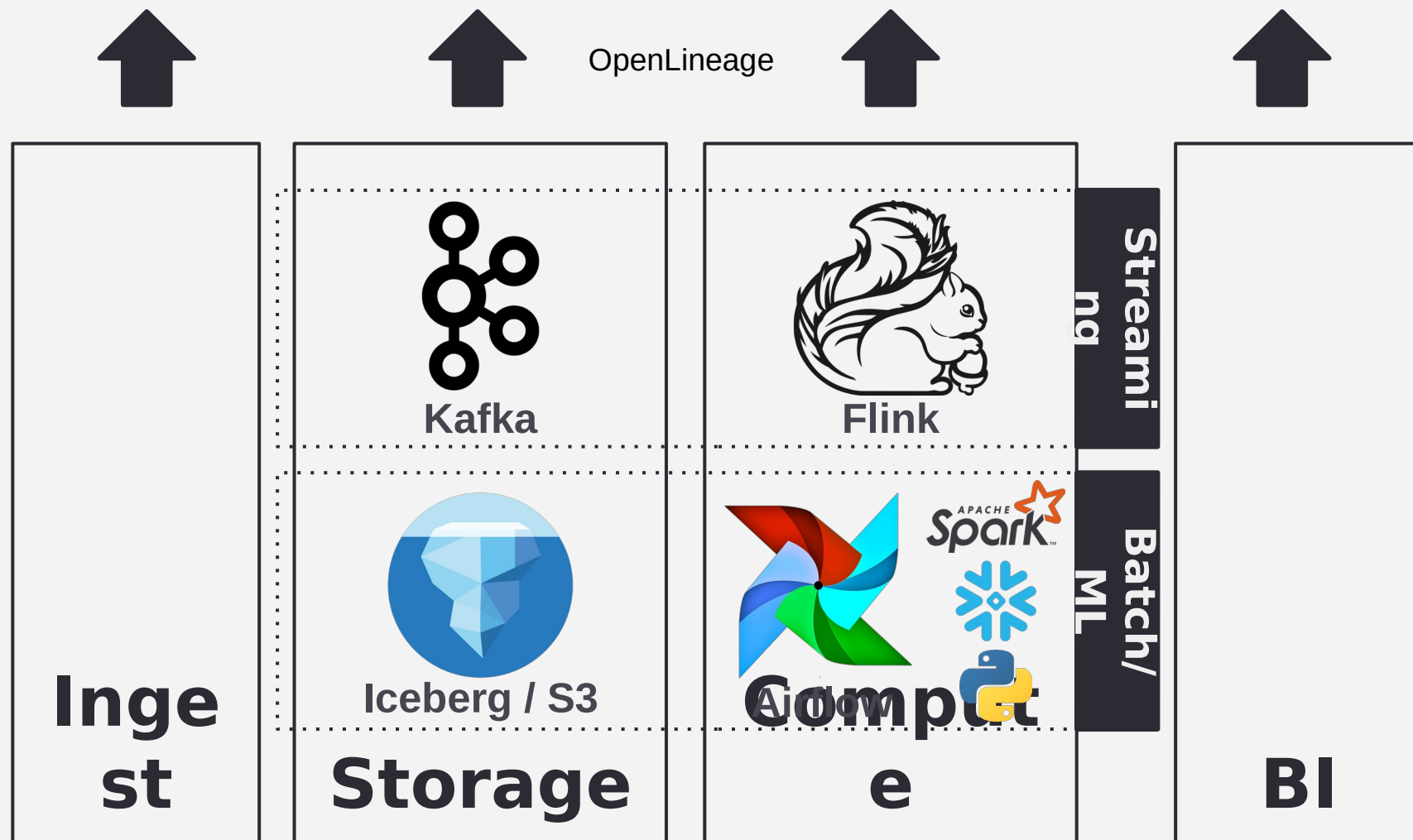


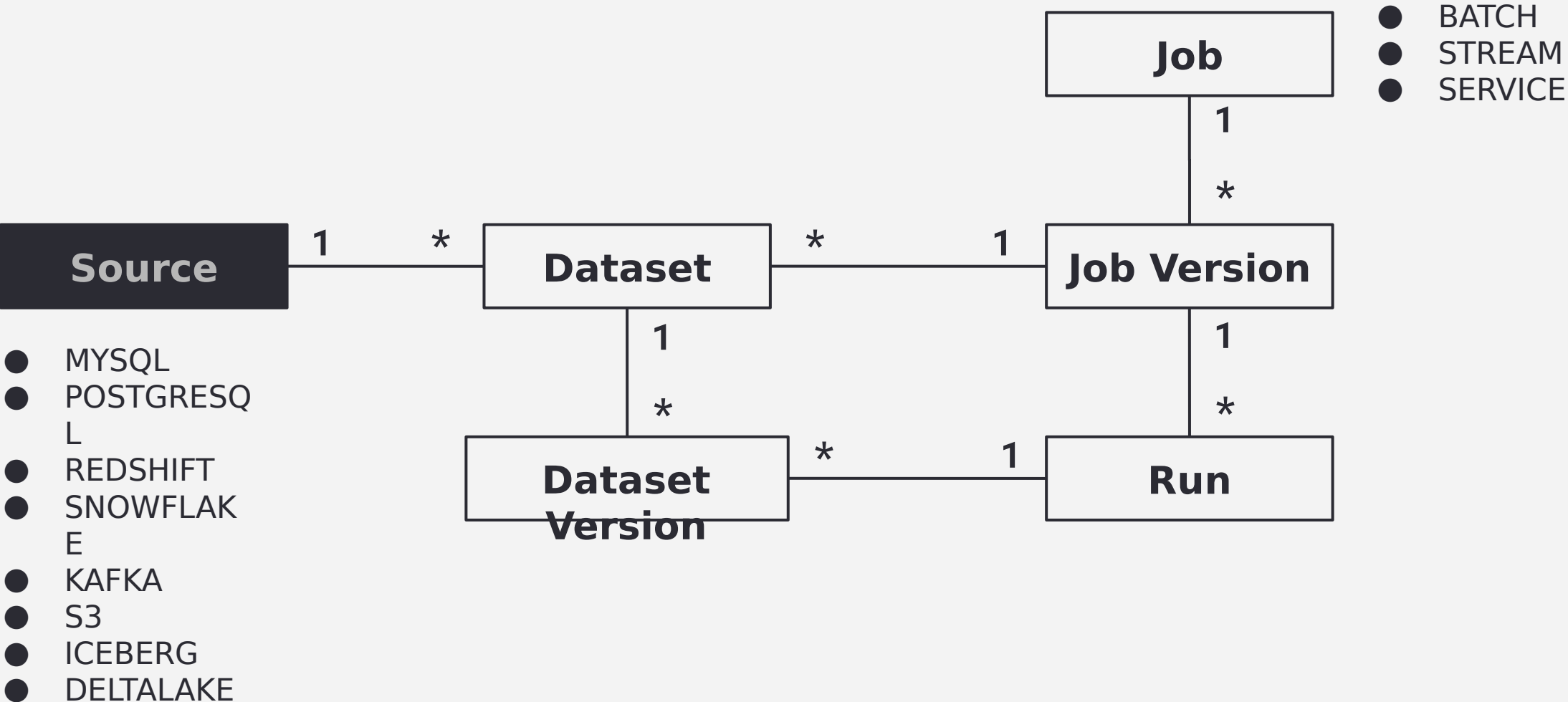
**MARQUEZ**



# Metadata: MARQUEZ

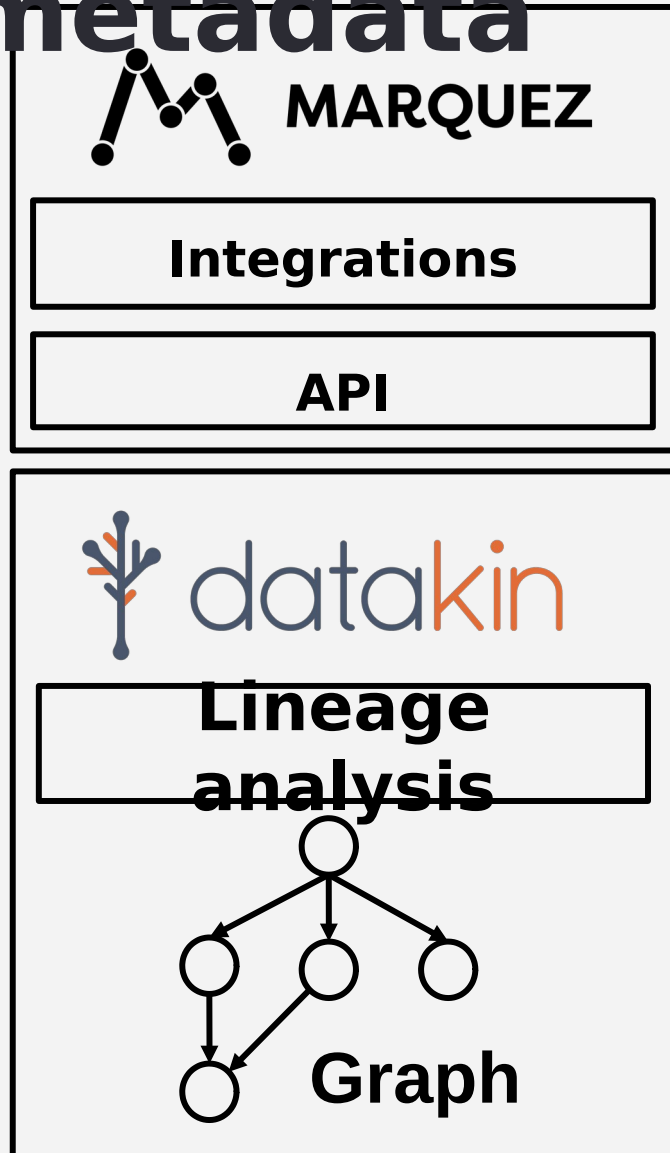
- **Data Platform** built around **Marquez**
- **Integrations**
  - Ingest
  - Storage
  - Compute





# 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?

# Data observability in practice

# Integrations

## Covered:



Google BigQuery



## Beta:



great\_expectations



# Airflow integration

```
pip3 install marquez-airflow
```

```
MARQUEZ_BACKEND=HTTP
```

```
MARQUEZ_URL=http://  
marquez.example.com
```

```
MARQUEZ_NAMESPACE=my_namespace
```



# Spark java agent

**spark.driver.extraJavaOptions:**

**-javaagent:**marquez-spark-{version}.jar=  
http://marquez.example.com/api/v1/  
namespaces/my\_namespace/jobs/  
{job\_name}/runs/{uuid}





# Metadata collected

**Lineage:** inputs/outputs

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

**Logical plan**

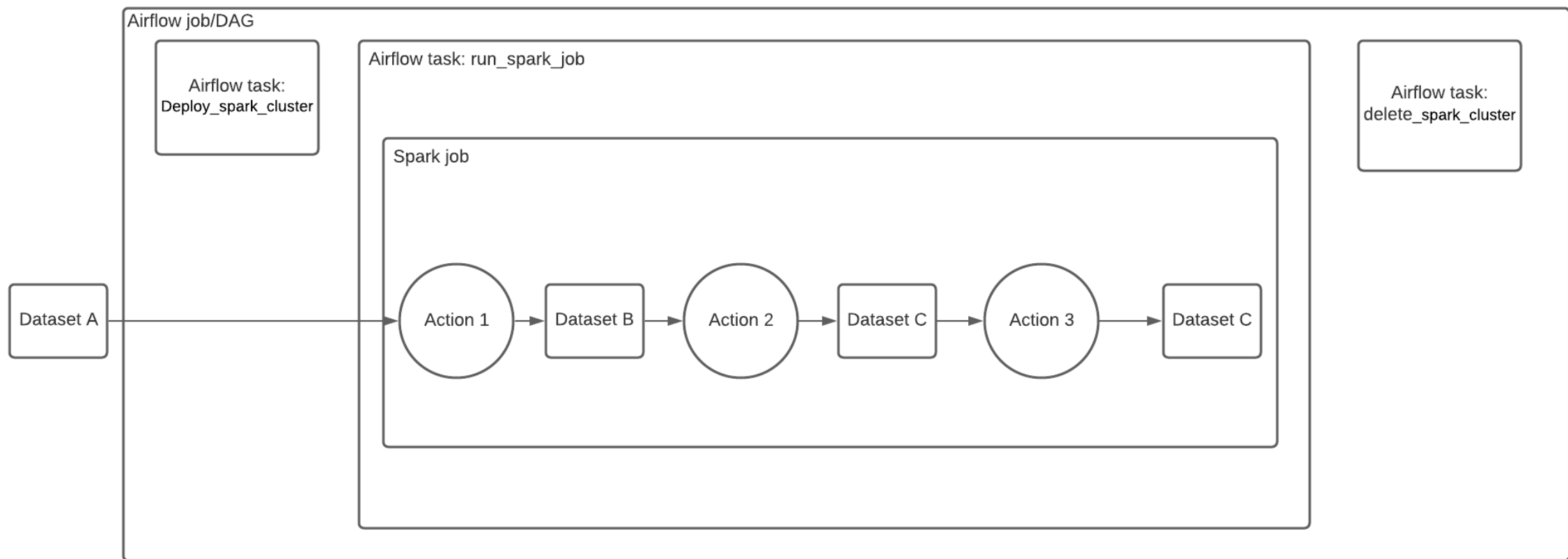
**source code location / version**

**schema**

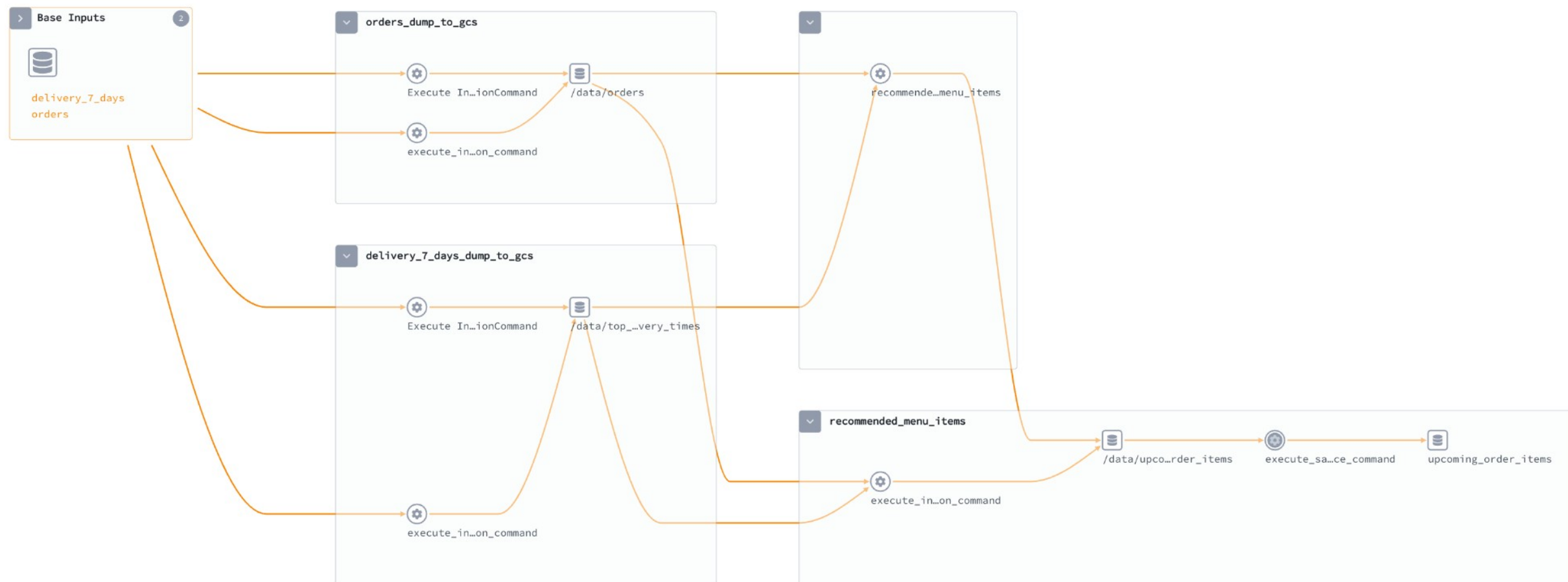
**processing time**



# Lineage model



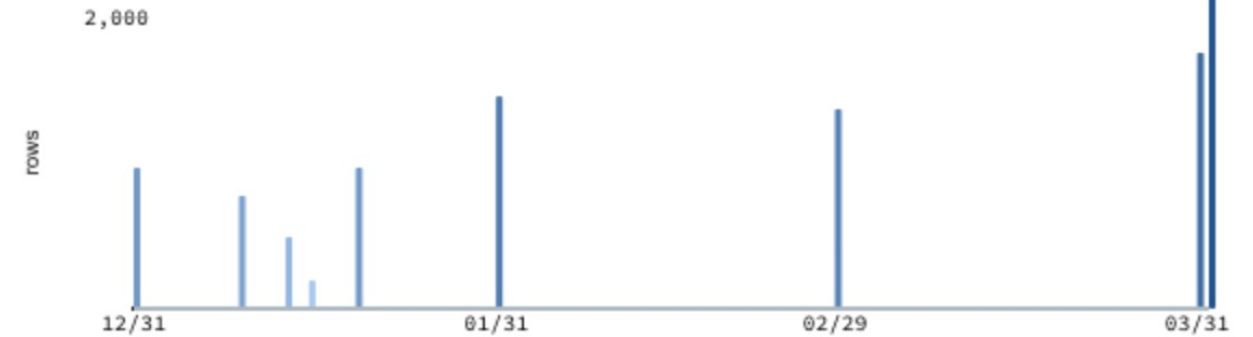
# Lineage Example across jobs



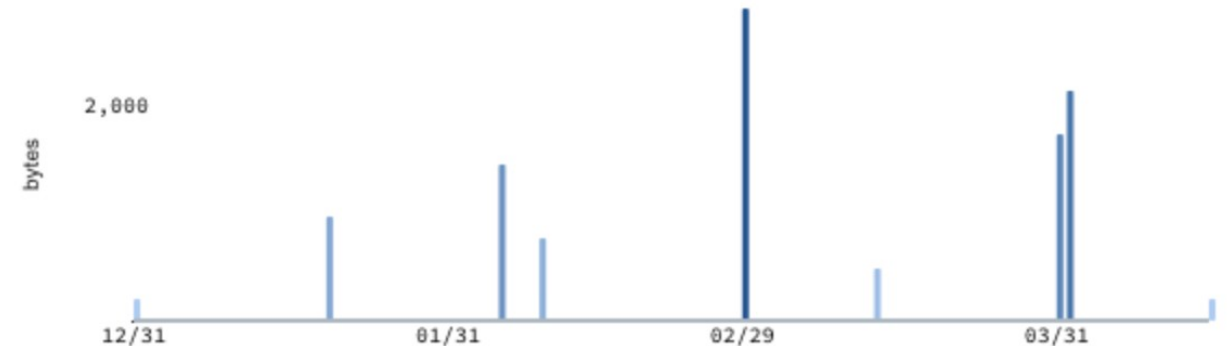
# Example of OpenLineage metadata usage:

## Data volume evolution

Row count



Byte count



# Join the conversation

## OpenLineage:

**Github:** [github.com/OpenLineage](https://github.com/OpenLineage)★

**Slack:** [OpenLineage.slack.com](https://openlineage.slack.com)

**Twitter:** [@OpenLineage](https://twitter.com/OpenLineage)🐦

**Email:** [groups.google.com/g/openlineage](https://groups.google.com/g/openlineage)

## Marquez:

**Github:** [github.com/MarquezProject/marquez](https://github.com/MarquezProject/marquez)★

**Slack:** [MarquezProject.slack.com](https://marquezproject.slack.com)

**Twitter:** [@MarquezProject](https://twitter.com/MarquezProject)🐦