



# Simplifying Data Lineage: How OpenLineage Empowers Airflow and Beyond

3.0

Julien Le Dem & Harel Shein, Datadog

# Who we are

Datadog is a world-class data platform ingesting more than a 100 trillion events a day, providing real-time insights.



**Harel Shein**  
Senior Engineering Manager  
OpenLineage TSC



**Julien Le Dem**  
Principal Engineer  
co-creator of Parquet, Arrow  
and OpenLineage



# Imagine we work at Canoe

The screenshot shows the homepage of the Canoe travel website. At the top, there's a navigation bar with the word "Canoe" on the left and links for "Hotels", "Flights", "Car Rentals", "Things to Do", and "Sign in". Below the navigation is a large banner with the text "Find the best deals on hotels and flights" overlaid on a background image of a coastal town with orange-roofed buildings and hills. A search form is centered in the banner, featuring fields for "Destination" and "Check in" with a blue "Search" button. Below the banner, the page title "Featured Hotels" is displayed between two yellow sun icons. There are four hotel cards visible: "Oceanview Inn" in Miami Beach, FL, "Grand Hotel" in Paris, France, "Mountain Lodge" in Aspen, Colorado, and "Cityscape Hotel" in New York, NY. Each card includes a small photo of the hotel, the name, location, a five-star rating icon, and the price per night.

Canoe

Hotels Flights Car Rentals Things to Do Sign in

Find the best deals on hotels and flights

Destination Check in Search

Featured Hotels

Oceanview Inn  
Miami Beach, FL  
★★★★★  
\$210 /night

Grand Hotel  
Paris, France  
★★★★★  
\$350 /night

Mountain Lodge  
Aspen, Colorado  
★★★★★  
\$250 /night

Cityscape Hotel  
New York, NY  
★★★★★  
\$450 /night

# But something is wrong...

The screenshot shows the homepage of the Canoe travel website. At the top, there is a navigation bar with the word "Canoe" on the left and links for "Hotels", "Flights", "Car Rentals", "Things to Do", and "Sign in". Below the navigation is a large banner with the text "Find the best deals on hotels and flights" overlaid on a background image of a coastal town with orange-roofed buildings and hills. A search form is centered in the banner, featuring fields for "Destination" and "Check in" and a blue "Search" button. Below the banner, the page title "Featured Hotels" is displayed, followed by four hotel cards:

- Oceanview Inn**  
Miami Beach, FL  
★★★★  
\$150 /night
- Grand Hotel**  
Paris, France  
★★★★  
\$300 /night
- Mountain Lodge**  
Aspen, Colorado  
★★★★
- Cityscape Hotel**  
New York, NY  
★★★★

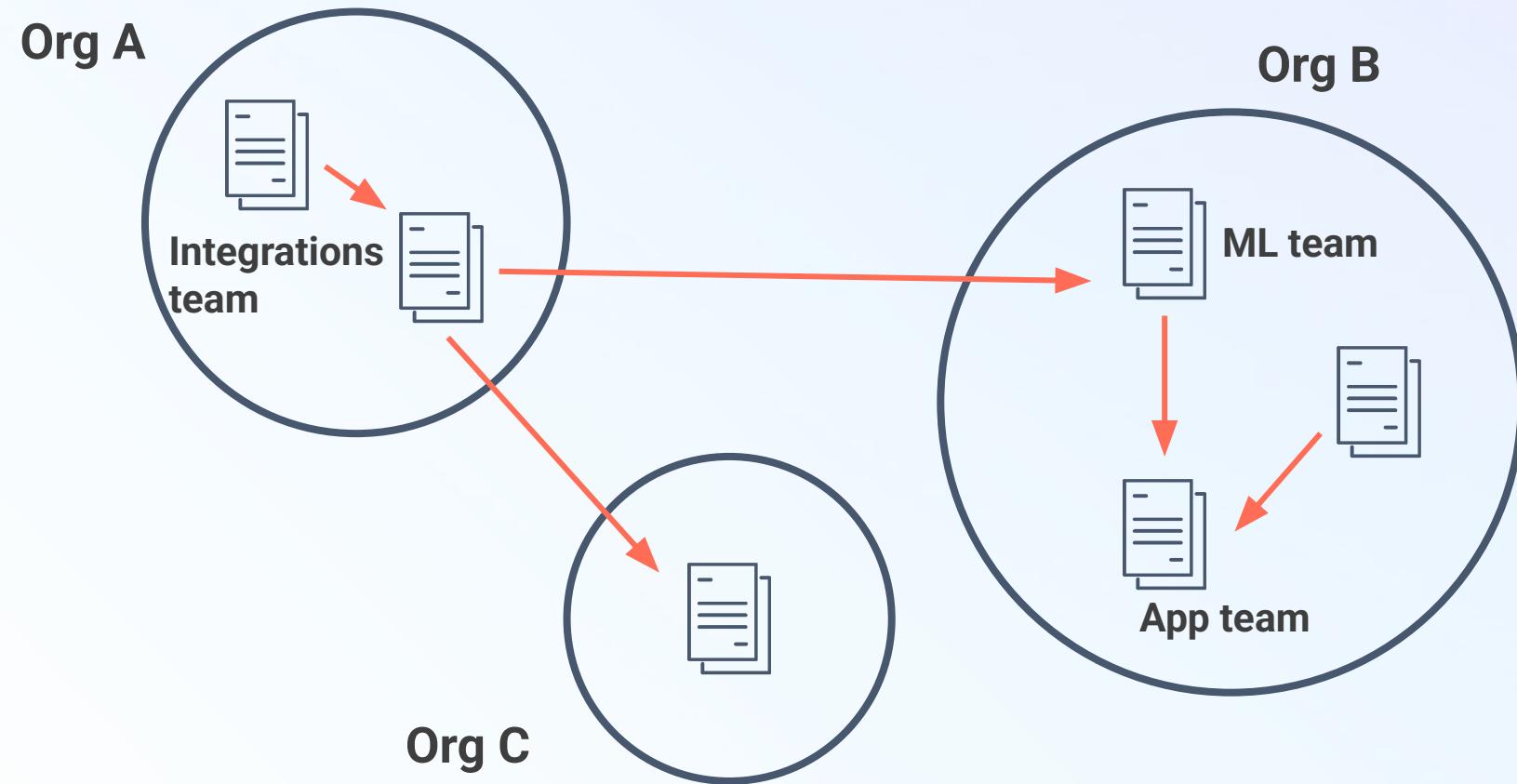
# But something is wrong...

The screenshot shows a travel search interface with the word "Canoe" in the top left corner. The top navigation bar includes links for Hotels, Flights, Car Rentals, Things to Do, and Sign in. A large banner in the center says "Find the best deals on hotels and flights" over a background image of a coastal town. Below the banner is a search form with fields for Destination and Check in, and a blue Search button. The main content area is titled "Featured Hotels" and lists four options:

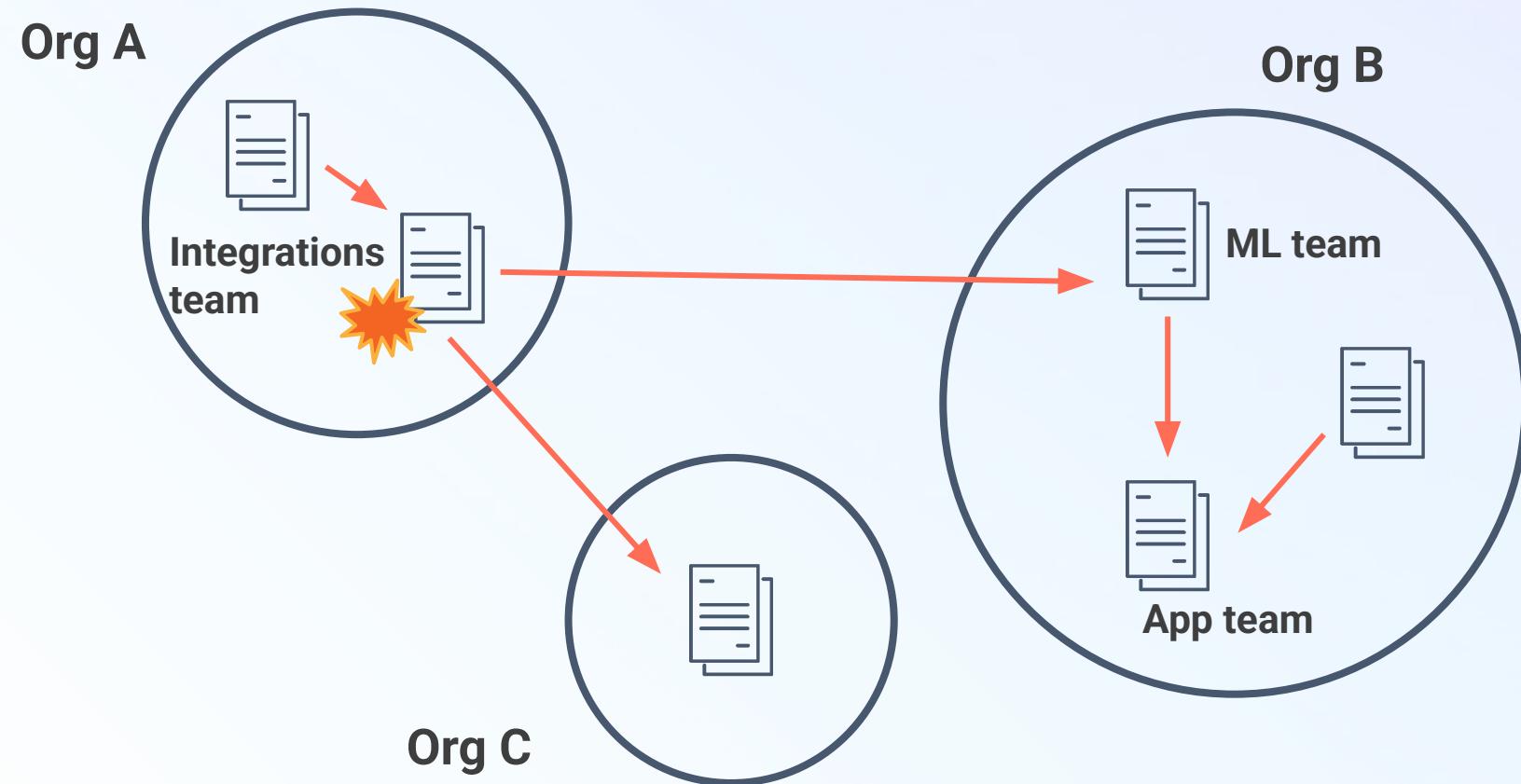
Hotel Name	Location	Rating	Price
Oceanview Inn	Miami Beach, FL	★★★★	\$150 /night
Grand Hotel	Paris, France	★★★★	\$300 /night
Mountain Lodge	Aspen, Colorado	★★★★	
Cityscape Hotel	New York, NY	★★★★	

Two black bat icons are positioned on the left and right sides of the page, appearing to fly around the content.

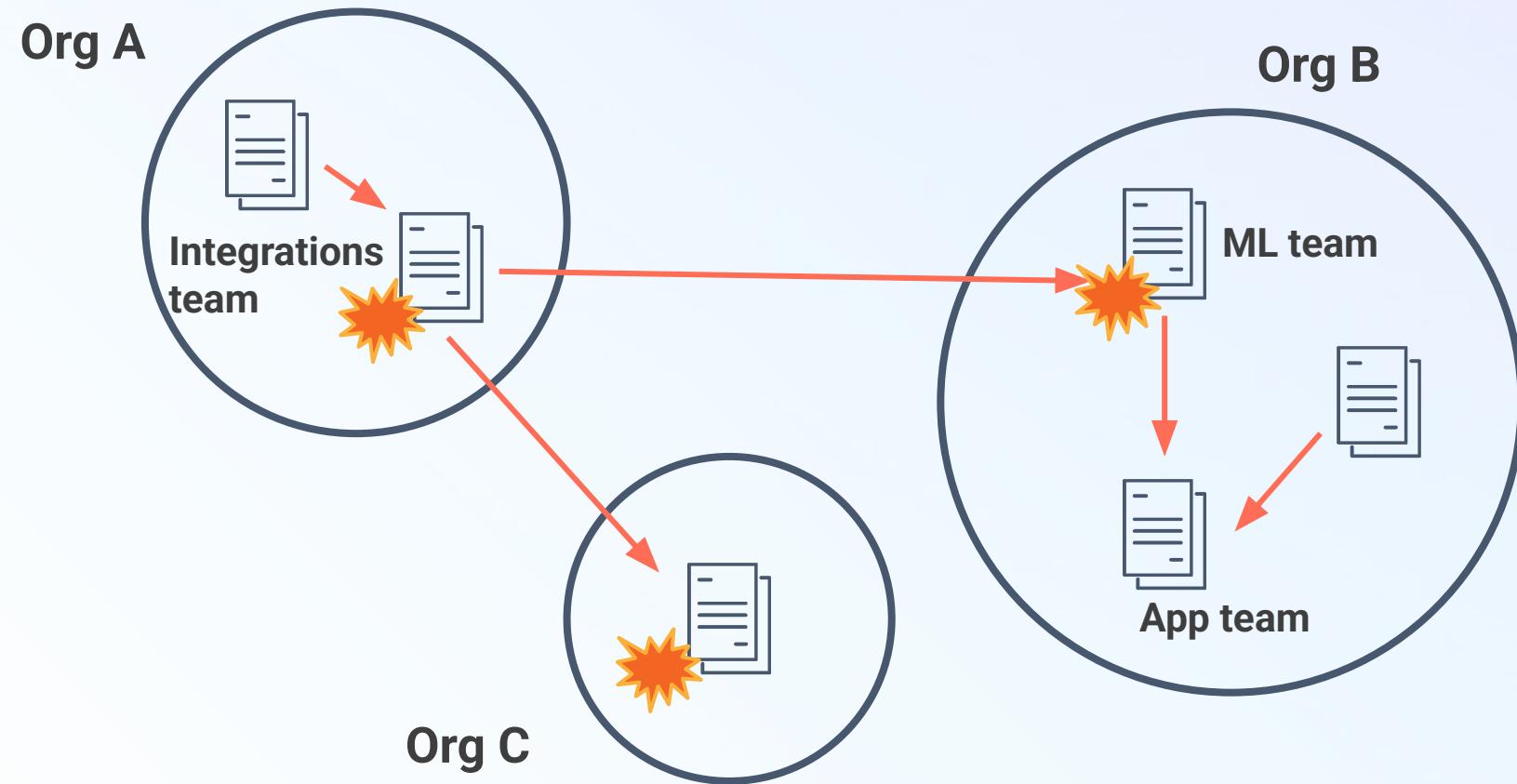
# Data ecosystem at Canoe



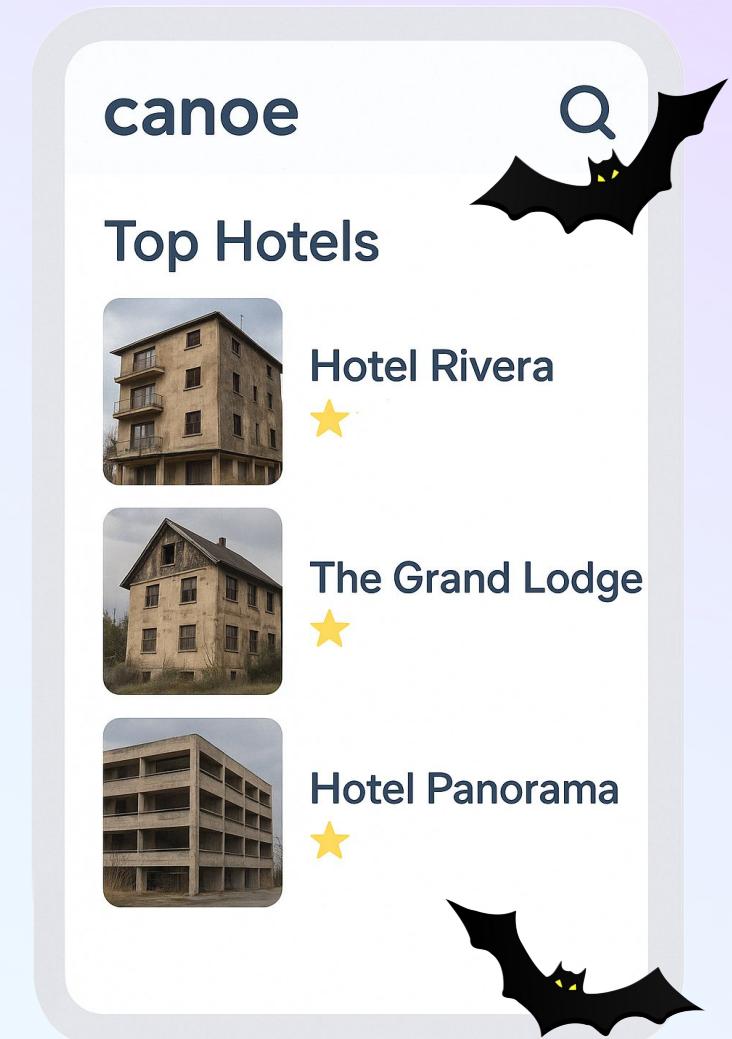
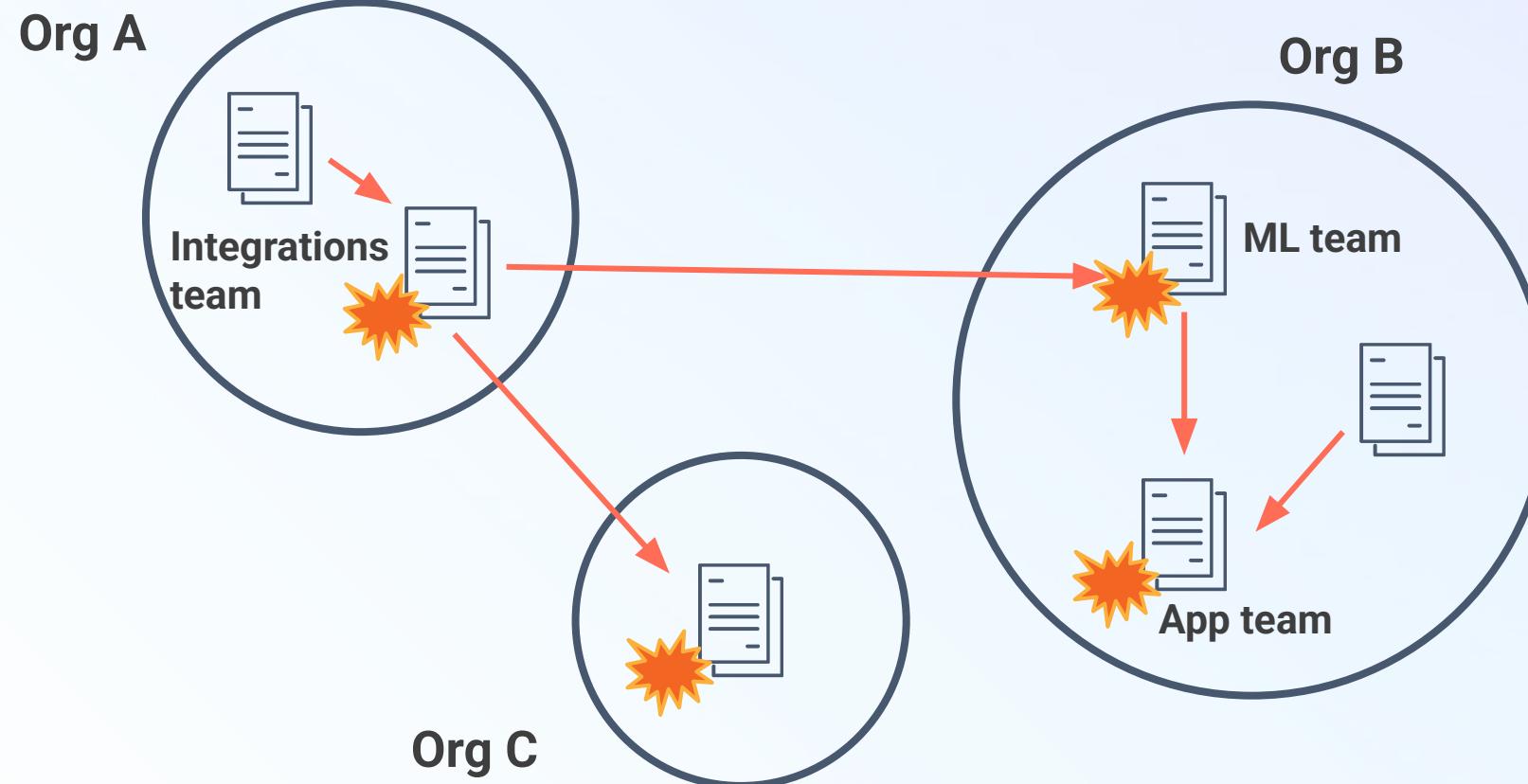
# Data ecosystem at Canoe



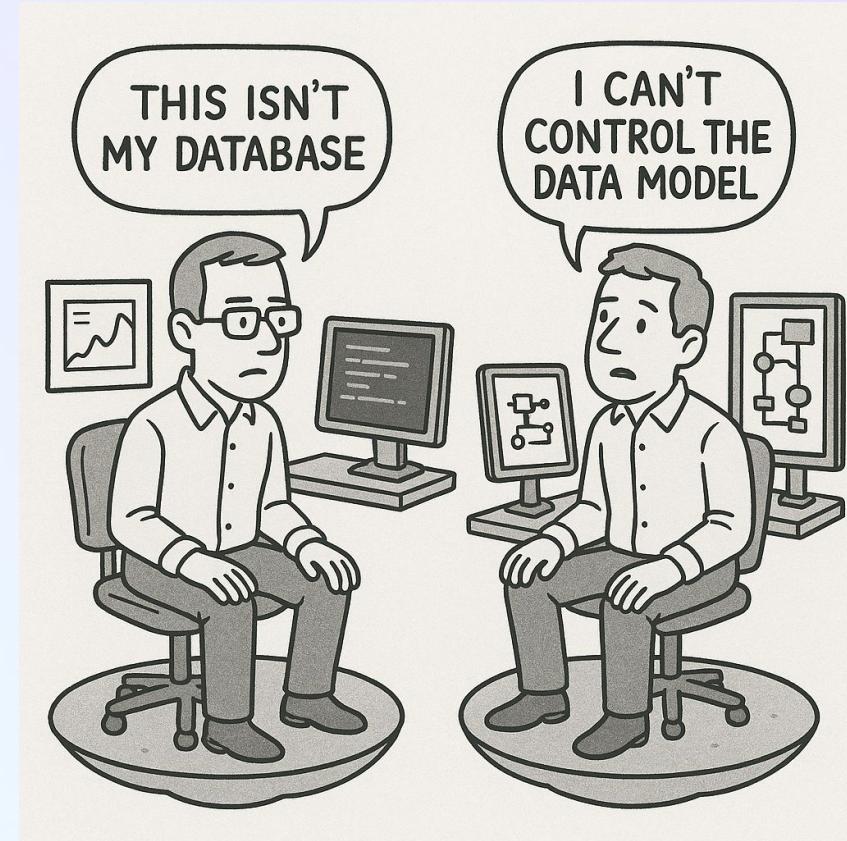
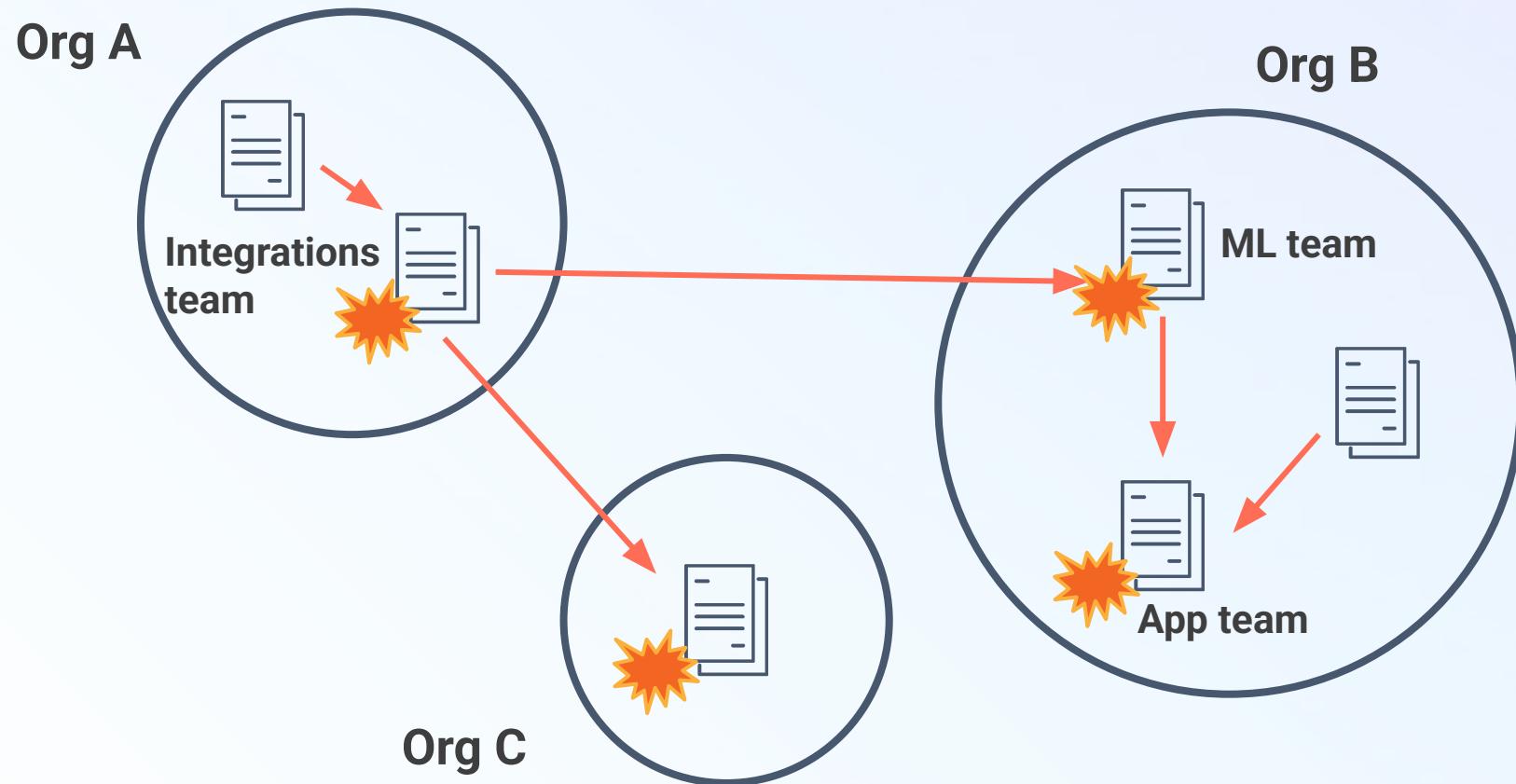
# Data ecosystem at Canoe



# Data ecosystem at Canoe



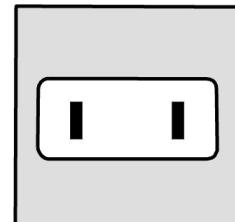
# Data ecosystem at Canoe



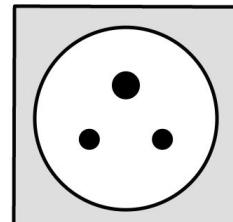
# And the ecosystem is fragmented



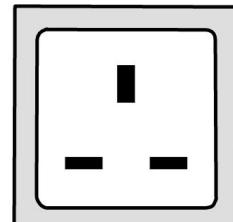
# We need everyone to speak the same language



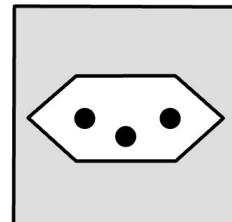
A



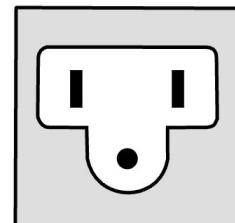
D



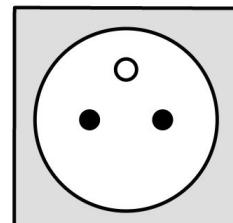
G



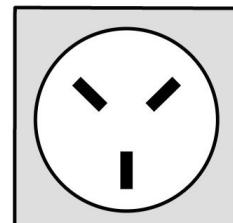
J



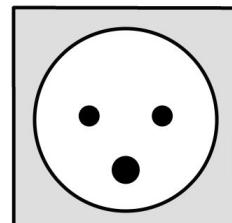
B



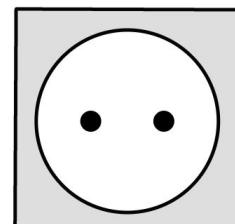
E



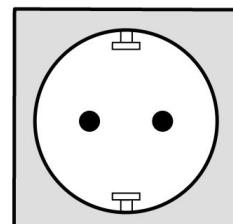
H



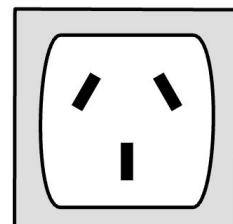
K



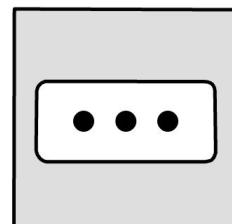
C



F



I

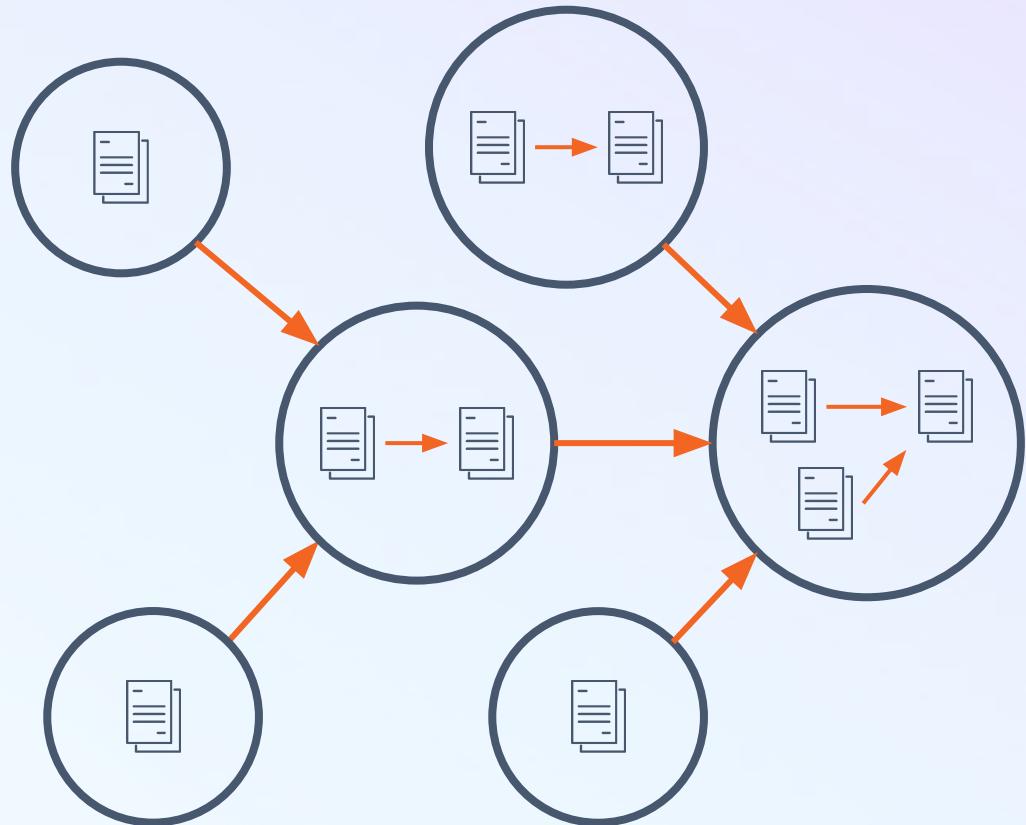


L

# Need a common language

Metadata about data pipelines:

- where data is from
- how it changes
- where it goes



# HOW STANDARDS PROLIFERATE:

(SEE: A/C CHARGERS, CHARACTER ENCODINGS, INSTANT MESSAGING, ETC)

SITUATION:  
THERE ARE  
14 COMPETING  
STANDARDS.

14?! RIDICULOUS!  
WE NEED TO DEVELOP  
ONE UNIVERSAL STANDARD  
THAT COVERS EVERYONE'S  
USE CASES.

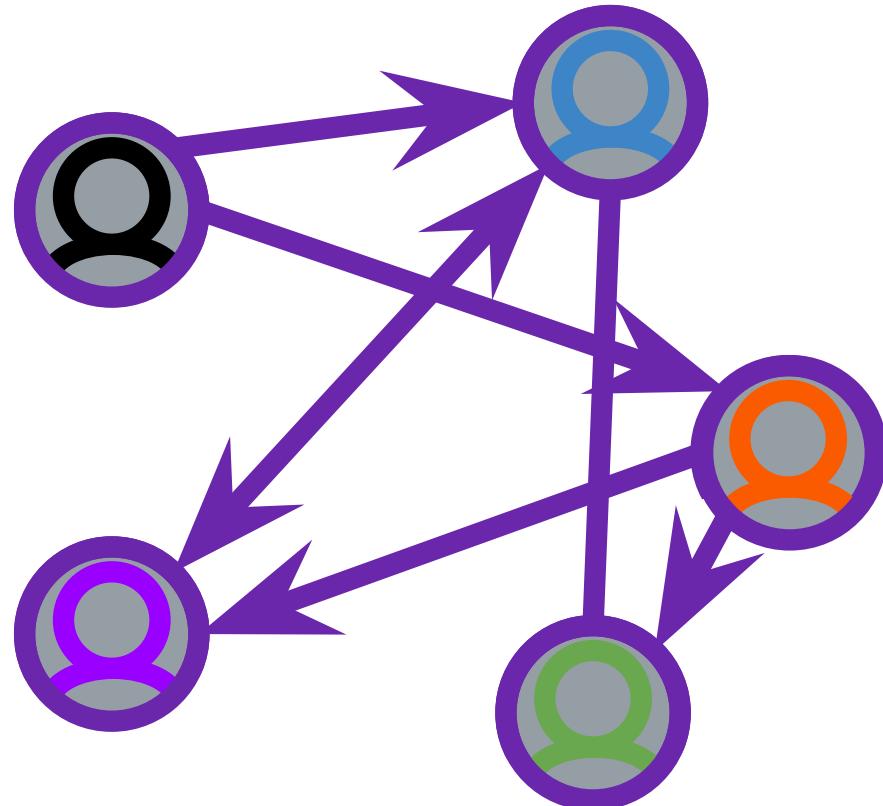


YEAH!

SOON:

SITUATION:  
THERE ARE  
15 COMPETING  
STANDARDS.

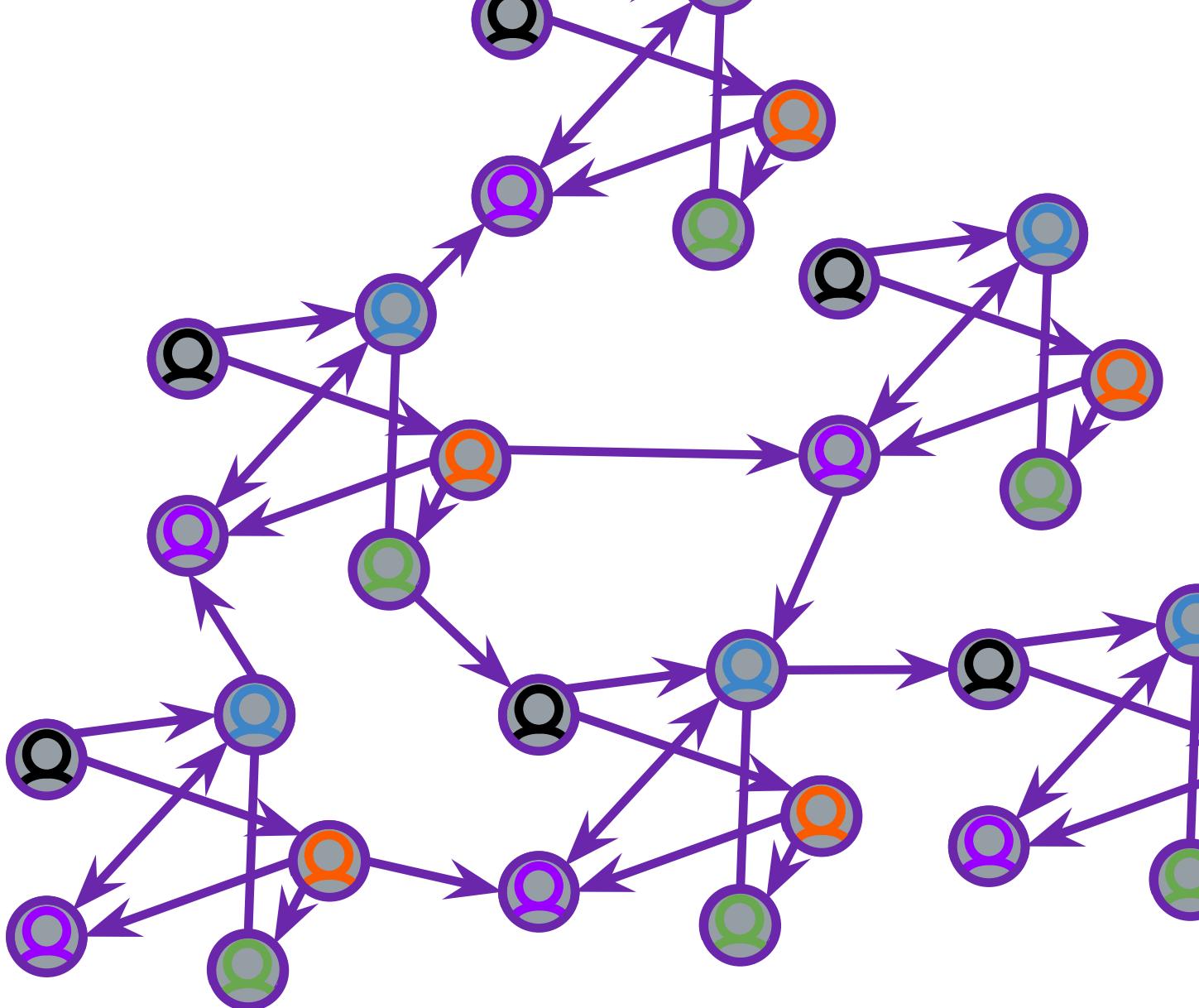
# How do standards get started?



# We Make Stone Soup



Align incentives  
to build a  
network effect

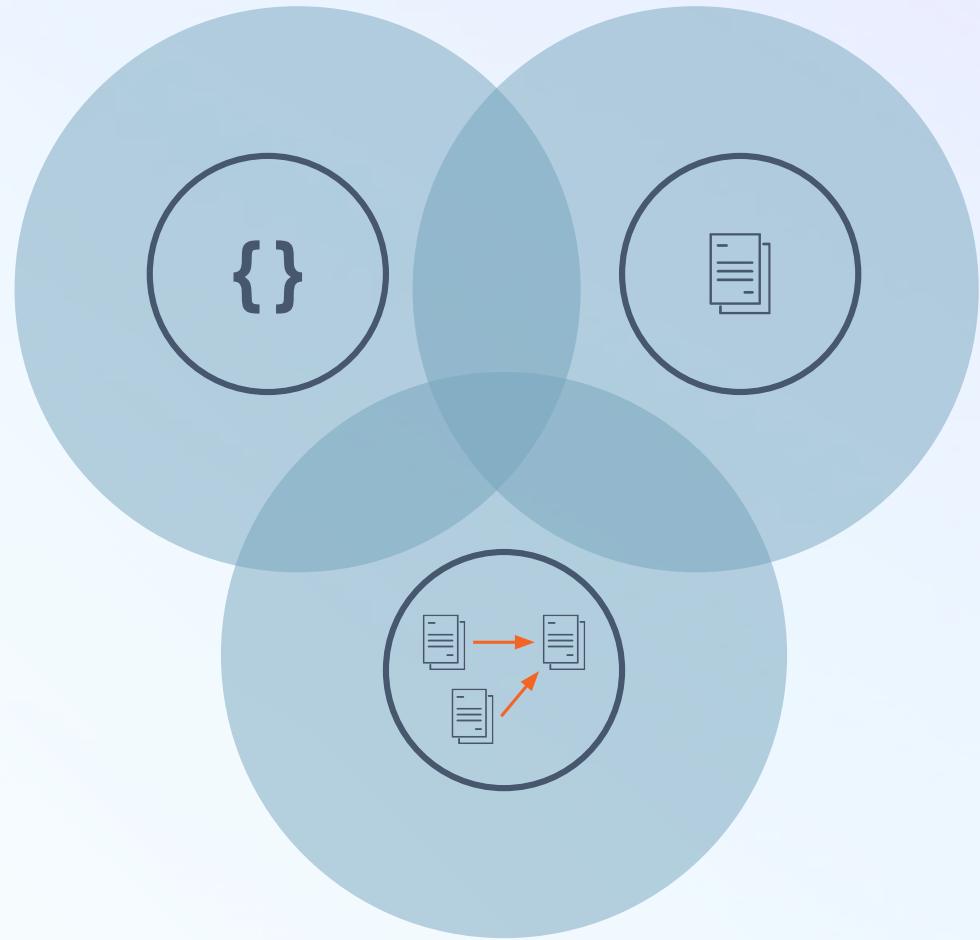


# The snowball effect

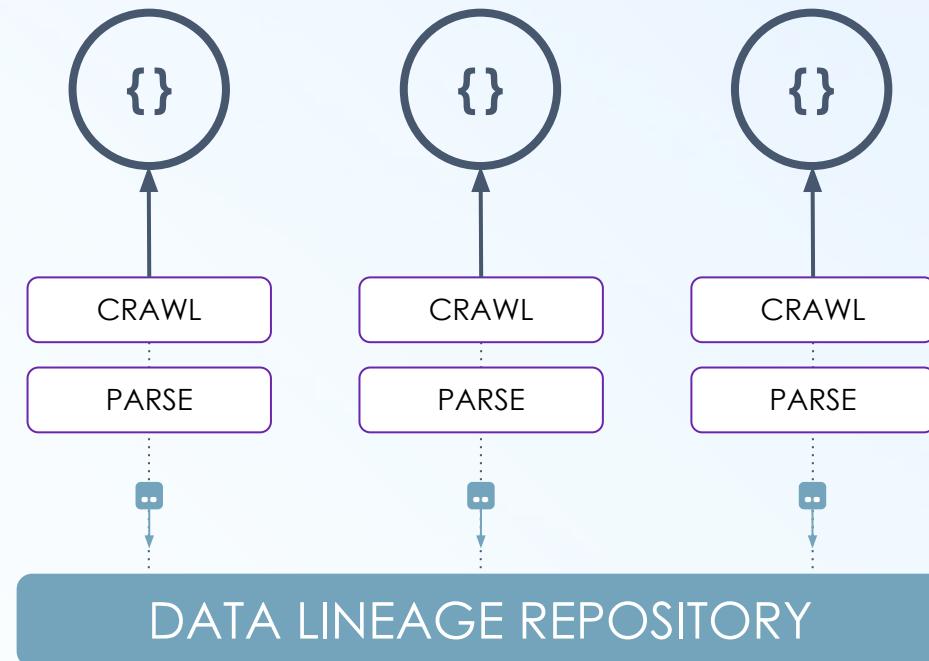


so... OpenLineage

# How to observe the lineage?



# Analyze source code

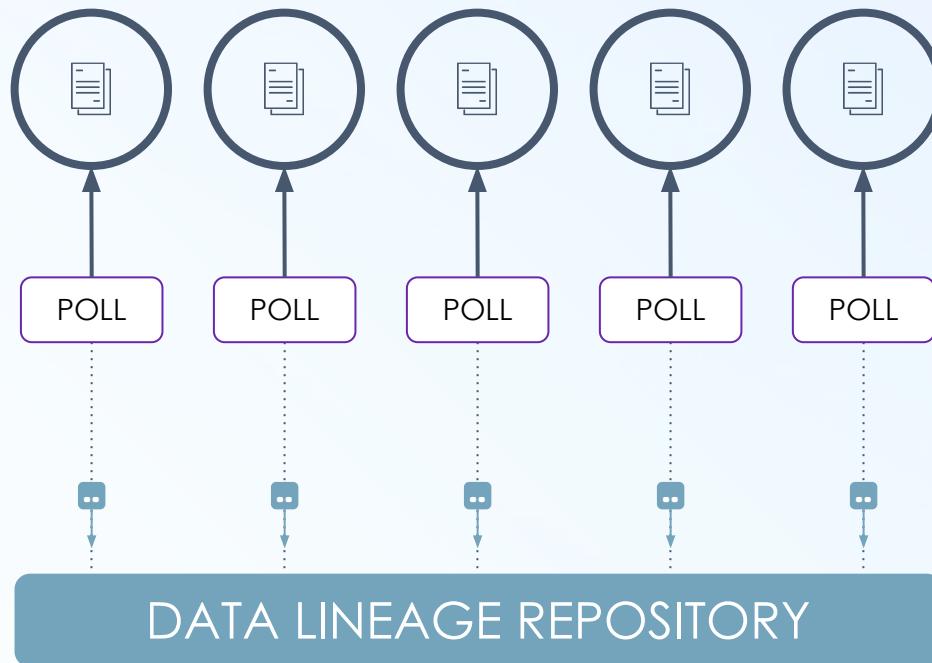


Integrate with source code repositories

Look for queries and parse them for lineage

Report to a lineage metadata repository

# Process query or activity logs

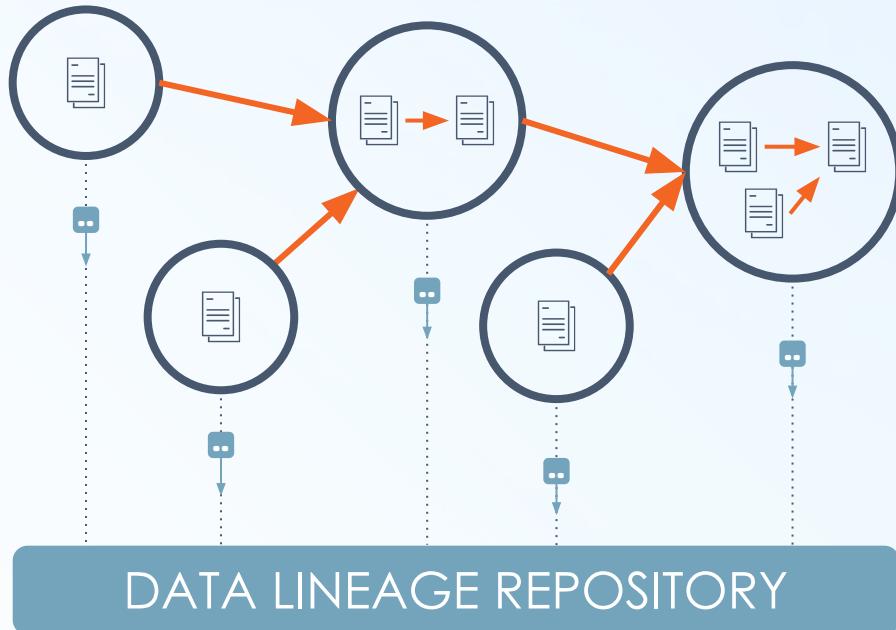


Integrate with data stores and warehouses

Regularly process query logs to trace lineage

Report to a lineage metadata repository

# Observe the pipeline



- Integrate with data orchestration systems
- As jobs run, observe the way they affect data
- Report to a lineage metadata repository

# Why Collect Data at Runtime?

The best time to collect metadata



You can try to infer data after the fact...



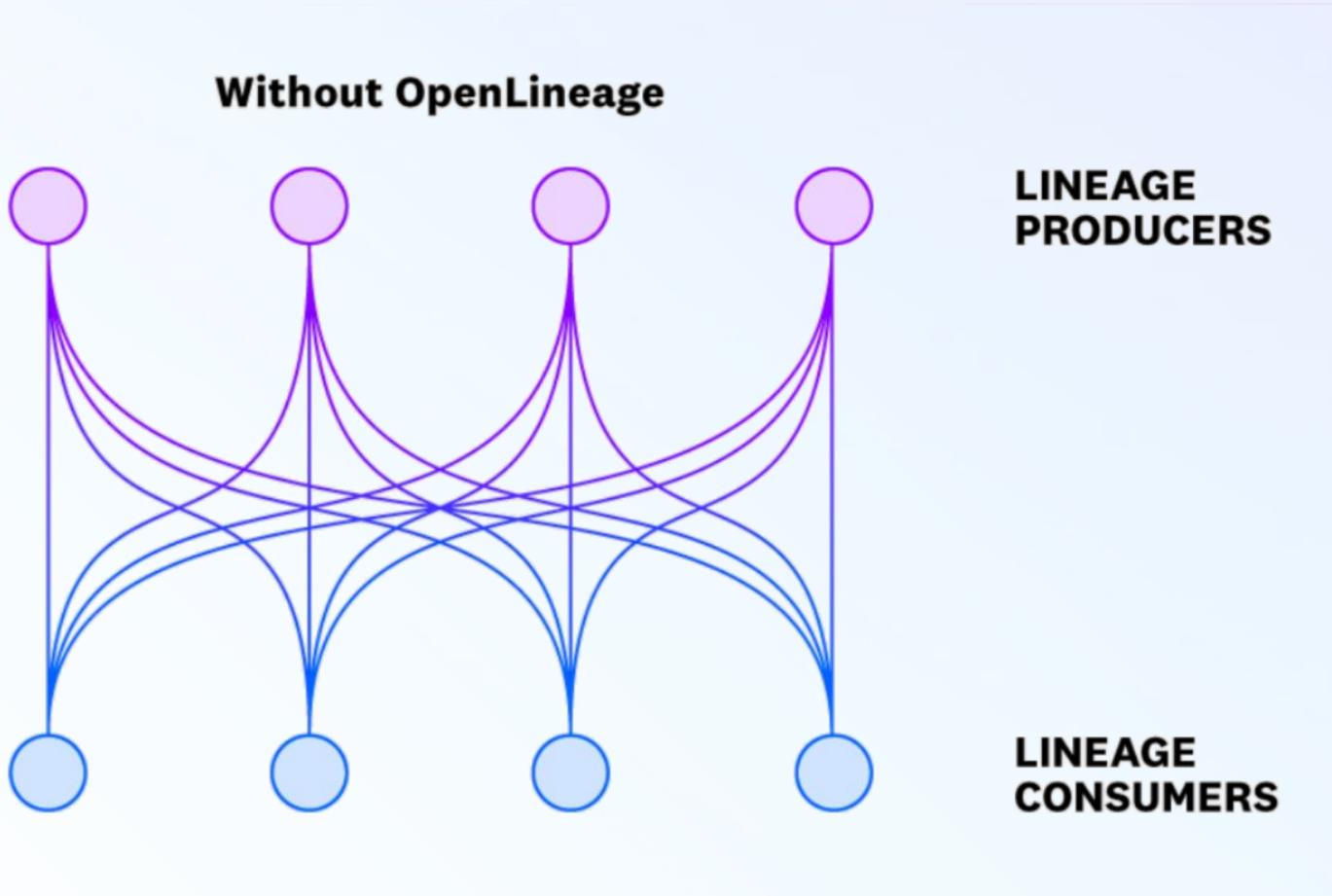
...or you can capture it when originally created.

# What is OpenLineage?

- Spec: open, vendor neutral
- Libraries: for common languages
- Integrations with data tools

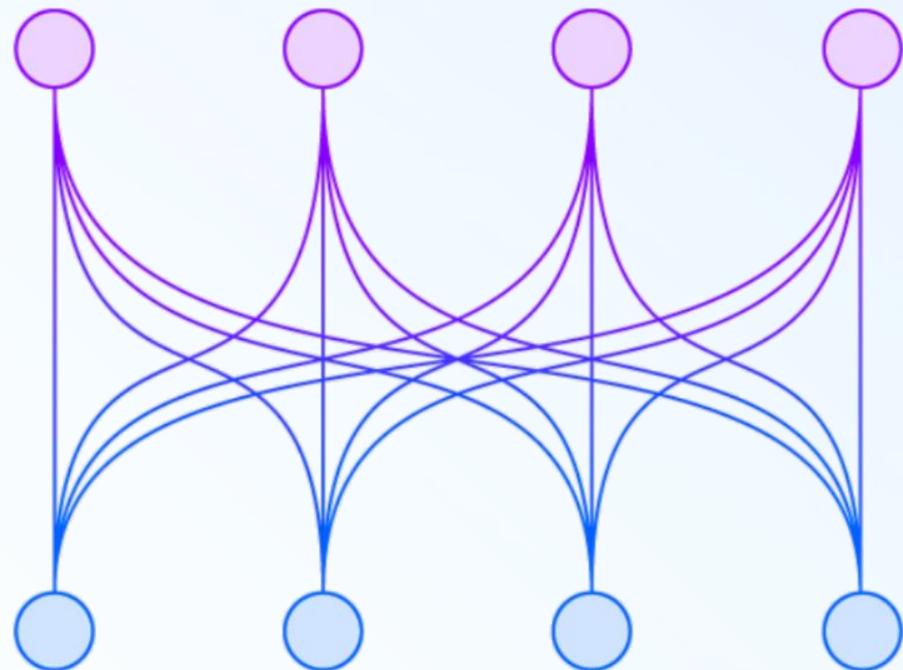


# Open Standards Matter

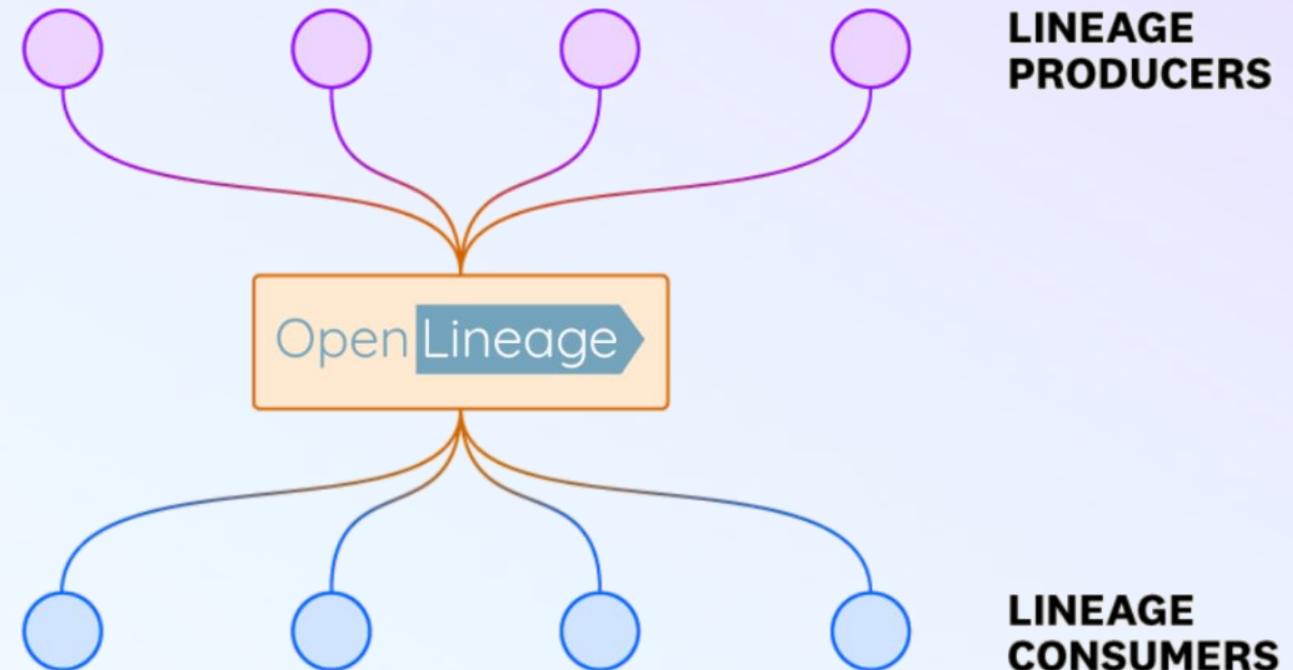


# Simplified

**Without OpenLineage**



**With OpenLineage**

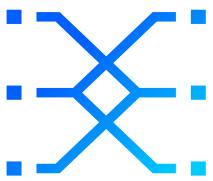


**LINEAGE  
PRODUCERS**

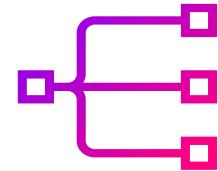
**LINEAGE  
CONSUMERS**

# OpenLineage at Datadog

# Datadog Data Observability

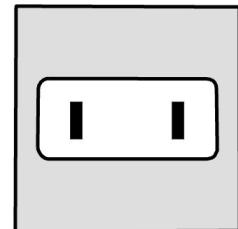


Data Jobs Monitoring

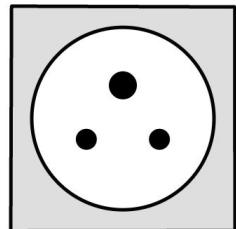


Data Quality  
Monitoring

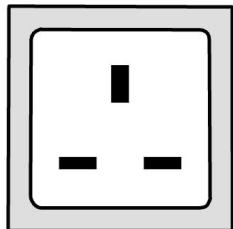
# OpenLineage benefits



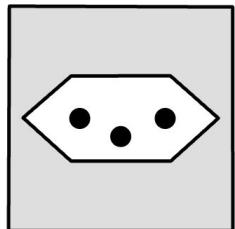
A



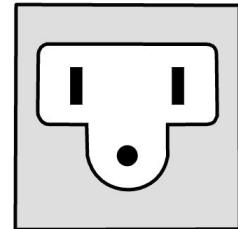
D



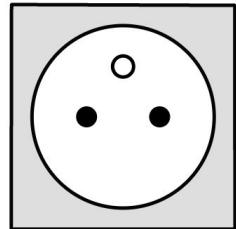
G



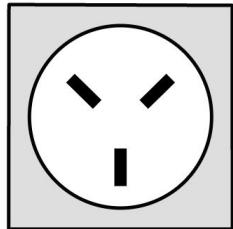
J



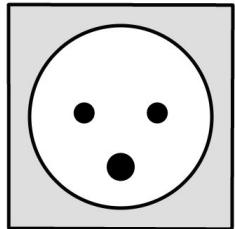
B



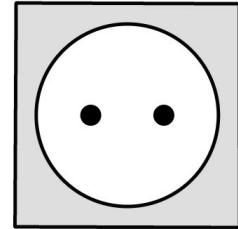
E



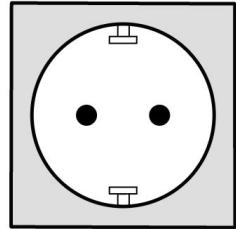
H



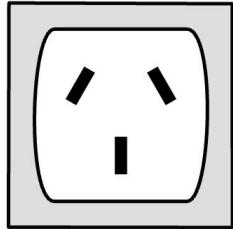
K



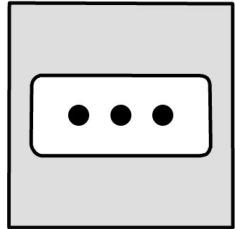
C



F



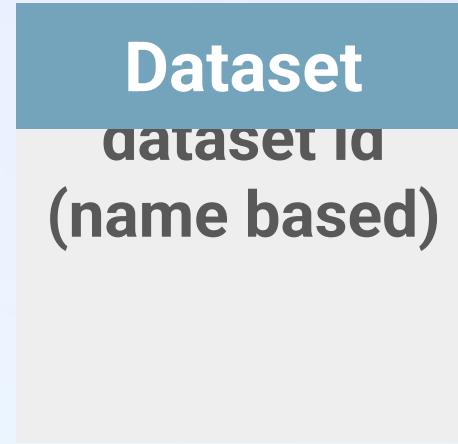
I



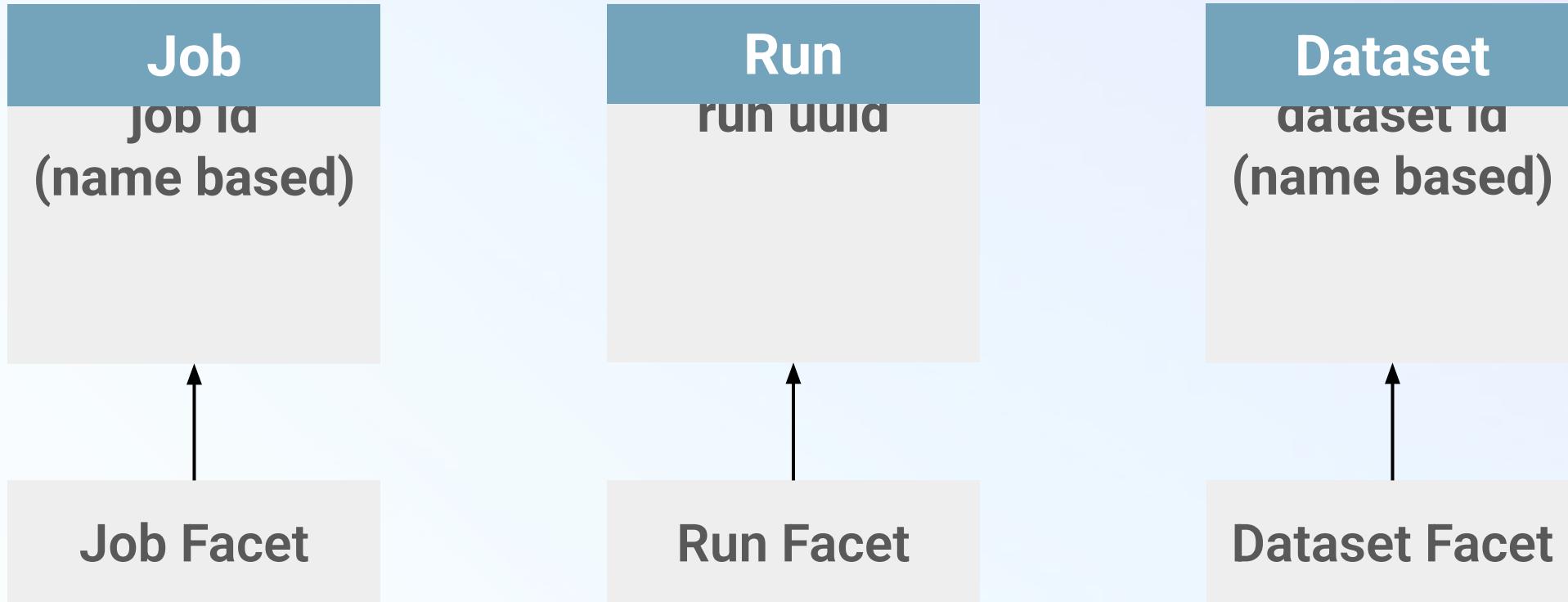
L

# OpenLineage Spec deep dive

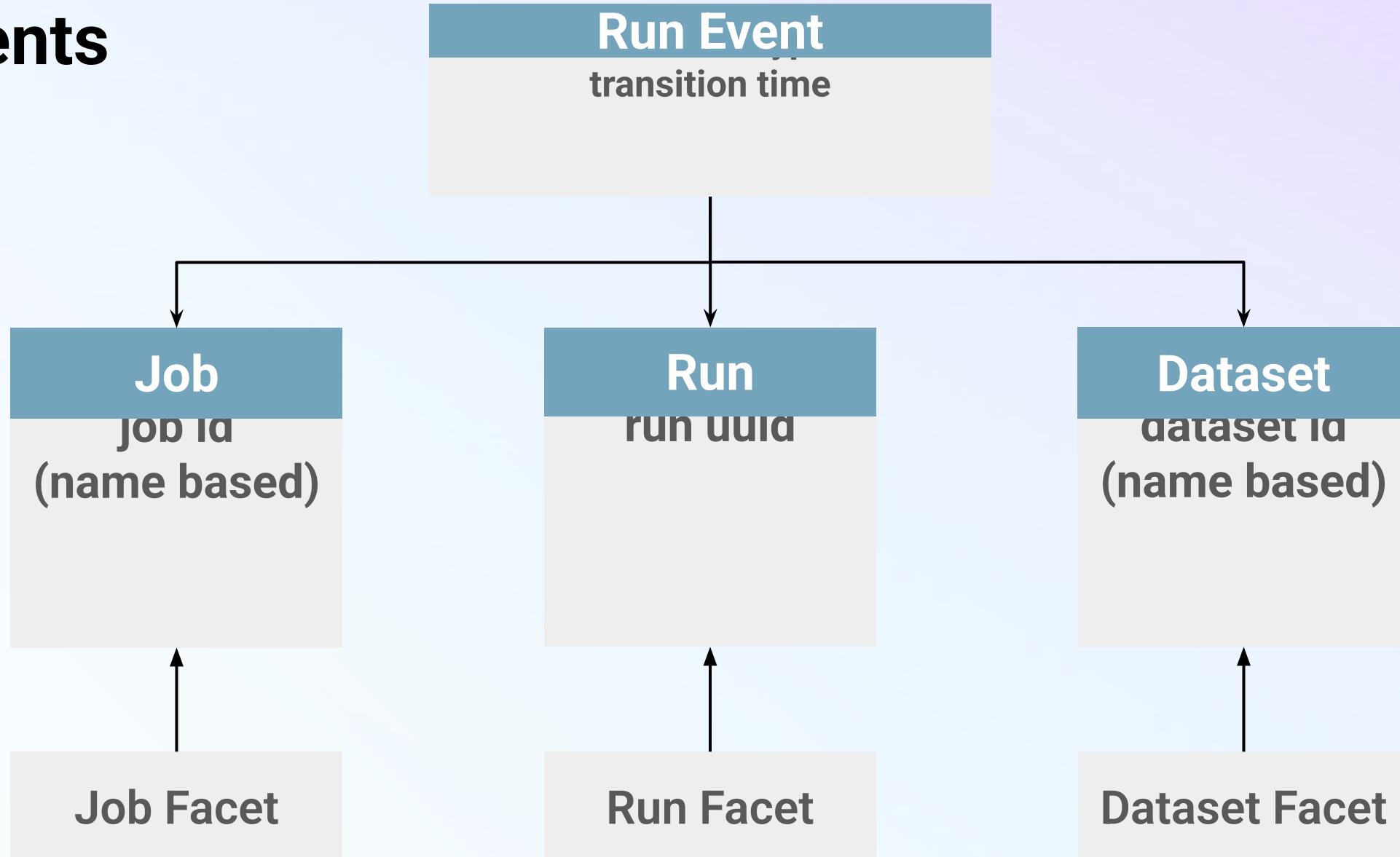
# Core Concepts: Job, Run, Dataset



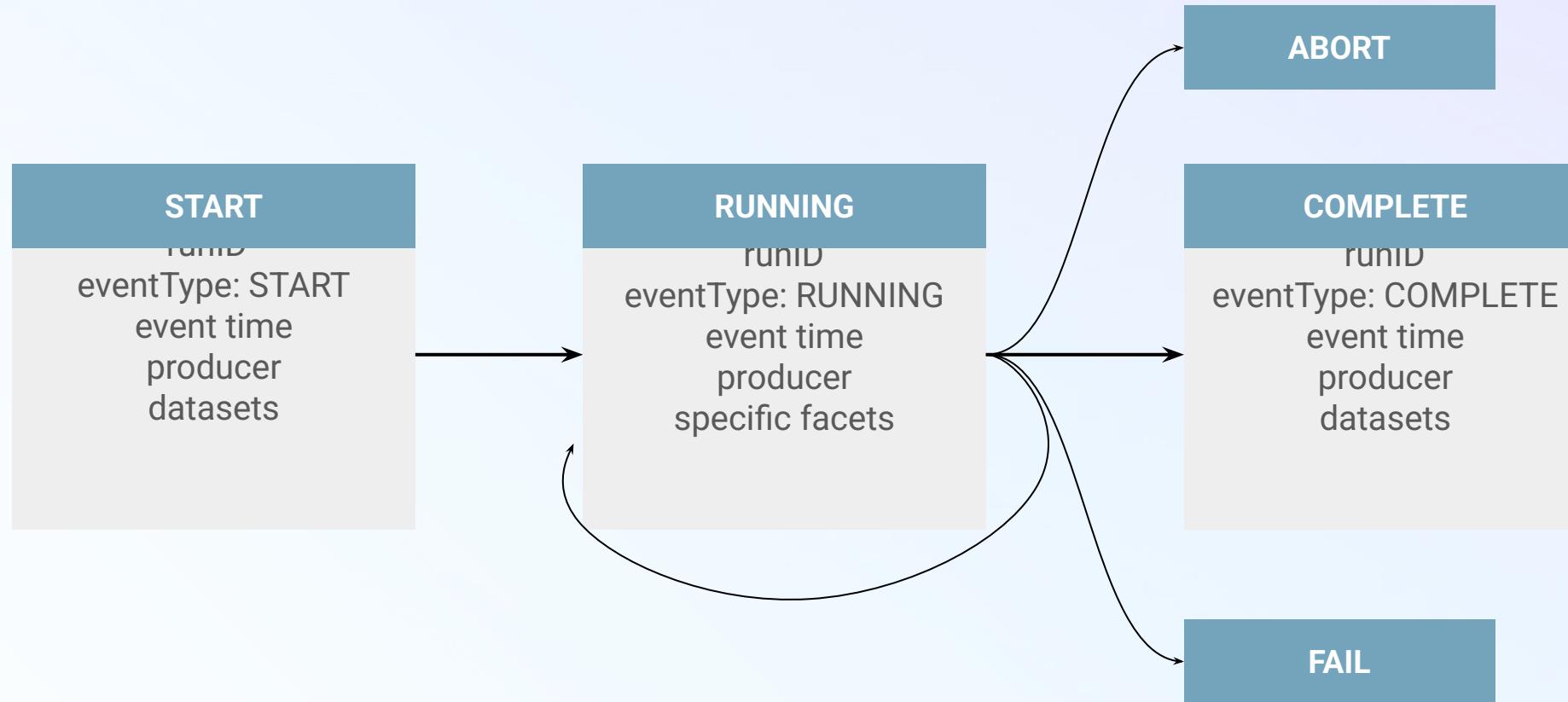
# Core Concepts: Facets



# Events



# Run Event Lifecycle



# Example of OpenLineage event

```
{  
  "eventType": "COMPLETE",  
  "eventTime": "2025-06-08T12:00:00.000Z",  
  "run": {  
    "runId": "123e4567-e89b-12d3-a456-426614174000",  
    "facets": {  
      "parent": {  
        "run": {  
          "runId": "111e4567-e89b-12d3-a456-426614174abc",  
          "namespace": "my_company_etl"  
        },  
        "job": {  
          "namespace": "my_company_etl",  
          "name": "parent_job"  
        }  
      }  
    }  
  },  
  "job": {  
    "namespace": "my_company_etl",  
    "name": "daily_sales_etl",  
    "facets": {  
      "documentation": {  
        "description": "ETL job to aggregate daily sales for dashboard reporting"  
      }  
    }  
  }  
},
```

```
"inputs": [  
  {  
    "namespace": "my_data_warehouse",  
    "name": "raw_sales_data",  
    "facets": {  
      "schema": {  
        "fields": [  
          {"name": "sale_id", "type": "string"},  
          {"name": "amount", "type": "float"},  
          {"name": "timestamp", "type": "timestamp"}  
        ]  
      }  
    }  
  }  
,  
  "outputs": [  
    {  
      "namespace": "my_data_warehouse",  
      "name": "aggregated_sales_data",  
      "facets": {  
        "outputStatistics": {  
          "rowCount": 2500,  
          "size": 102400,  
          "columnStats": {  
            "amount": {  
              "max": 999.99,  
              "min": 1.23,  
              "nullCount": 0  
            }  
          }  
        }  
      }  
    }  
,  
    "producer": "https://mycompany.com/lineage/collector"  
  }  
]
```

# Example of OpenLineage event

```
{  
  "eventType": "COMPLETE",  
  "eventTime": "2025-06-08T12:00:00.000Z",  
  "run": {  
    "runId": "123e4567-e89b-12d3-a456-426614174000",  
    "facets": {  
      "parent": {  
        "run": {  
          "runId": "111e4567-e89b-12d3-a456-426614174abc",  
          "namespace": "my_company_etl"  
        },  
        "job": {  
          "namespace": "my_company_etl",  
          "name": "parent_job"  
        }  
      }  
    }  
  }  
}  
  
"job": {  
  "namespace": "my_company_etl",  
  "name": "daily_sales_etl",  
  "facets": {  
    "documentation": {  
      "description": "ETL job to aggregate daily sales for dashboard reporting"  
    }  
  }  
},
```

JOB

```
"inputs": [  
  {  
    "namespace": "my_data_warehouse",  
    "name": "raw_sales_data",  
    "facets": {  
      "schema": {  
        "fields": [  
          {"name": "sale_id", "type": "string"},  
          {"name": "amount", "type": "float"},  
          {"name": "timestamp", "type": "timestamp"}  
        ]  
      }  
    }  
  }  
,  
  "outputs": [  
    {  
      "namespace": "my_data_warehouse",  
      "name": "aggregated_sales_data",  
      "facets": {  
        "outputStatistics": {  
          "rowCount": 2500,  
          "size": 102400,  
          "columnStats": {  
            "amount": {  
              "max": 999.99,  
              "min": 1.23,  
              "nullCount": 0  
            }  
          }  
        }  
      }  
    }  
,  
    "producer": "https://mycompany.com/lineage/collector"  
  ]
```

## product\_revenue\_report

Env All Status All

+ Filter

## Summary

## Summary

## Job Runs

## Tasks

ALL MONITORS 2 ALERT 3 OK

Latest run completed 4m 24s ago

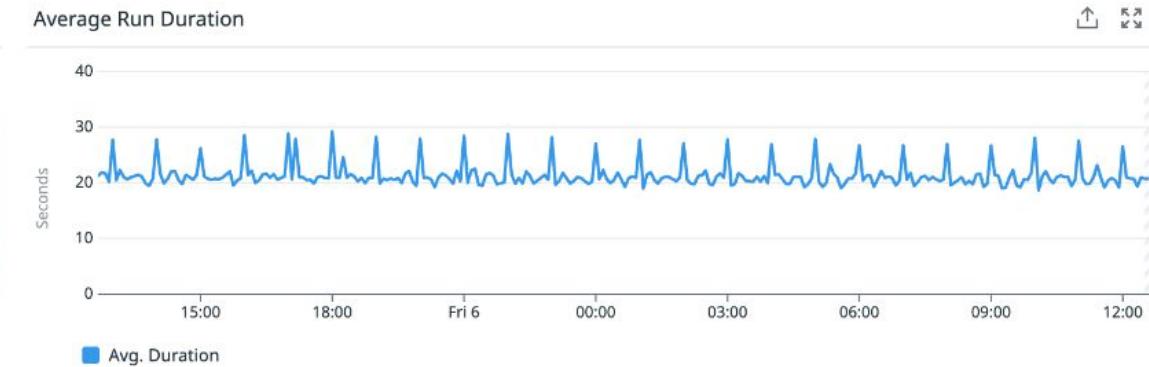
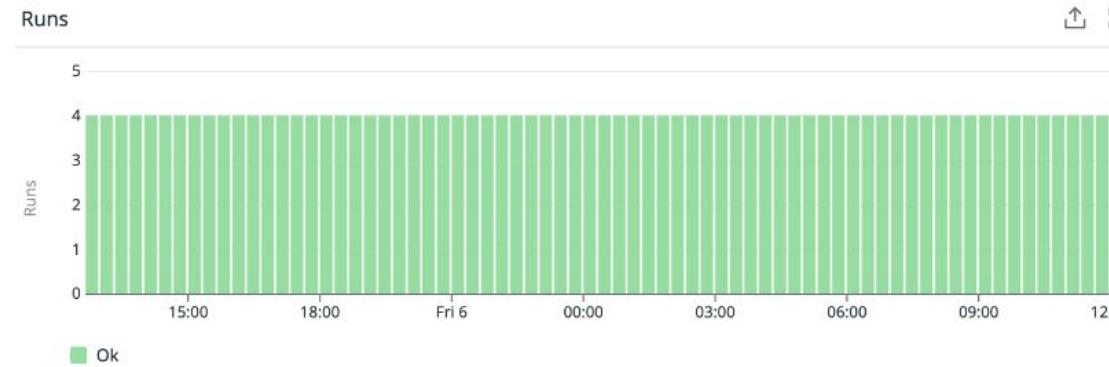
288 Total Runs

0% Failed

Average run duration is 21.3s

4%

Compared to previous period



## Runs

Showing 1-15 of 288 total runs

[View in Trace Explorer](#) 

↓ STARTED	STATUS	ENV	DURATION
Jun 7, 12:35:00 pm	OK	demo-env	20.8s
Jun 7, 12:30:01 pm	OK	demo-env	20.7s
Jun 7, 12:25:00 pm	OK	demo-env	21.0s

# Example of OpenLineage event

```
{  
  "eventType": "COMPLETE",  
  "eventTime": "2025-06-08T12:00:00.000Z",  
  "run": {  
    "runId": "123e4567-e89b-12d3-a456-426614174000",  
    "facets": {  
      "parent": {  
        "run": {  
          "runId": "111e4567-e89b-12d3-a456-426614174abc",  
          "namespace": "my_company_etl"  
        },  
        "job": {  
          "namespace": "my_company_etl",  
          "name": "parent_job"  
        }  
      }  
    }  
  },  
  "job": {  
    "namespace": "my_company_etl",  
    "name": "daily_sales_etl",  
    "facets": {  
      "documentation": {  
        "description": "ETL job to aggregate daily sales for dashboard reporting"  
      }  
    }  
  },  
}
```

RUN

```
"inputs": [  
  {  
    "namespace": "my_data_warehouse",  
    "name": "raw_sales_data",  
    "facets": {  
      "schema": {  
        "fields": [  
          {"name": "sale_id", "type": "string"},  
          {"name": "amount", "type": "float"},  
          {"name": "timestamp", "type": "timestamp"}  
        ]  
      }  
    }  
  }  
,  
  "outputs": [  
    {  
      "namespace": "my_data_warehouse",  
      "name": "aggregated_sales_data",  
      "facets": {  
        "outputStatistics": {  
          "rowCount": 2500,  
          "size": 102400,  
          "columnStats": {  
            "amount": {  
              "max": 999.99,  
              "min": 1.23,  
              "nullCount": 0  
            }  
          }  
        }  
      }  
    }  
,  
    "producer": "https://mycompany.com/lineage/collector"  
  }  
]
```

# visualization\_dag ⚡ RUNNING

Env All Status All

Summary

Job Runs

Tasks

Showing 1-15 of 1.91k

↓ STARTED

Oct 8, 1:25:01 pm

Oct 8, 1:20:01 pm

Oct 8, 1:15:01 pm

Oct 8, 1:10:00 pm

Oct 8, 1:05:00 pm

Oct 8, 1:00:01 pm

Oct 8, 12:55:00 pm

Oct 8, 12:50:00 pm

Oct 8, 12:45:01 pm

Oct 8, 12:40:01 pm

Oct 8, 12:35:01 pm

Oct 8, 12:30:00 pm

Oct 8, 12:25:01 pm

Oct 8, 12:20:01 pm

Oct 8, 12:15:00 pm

visualization\_dag > visualization\_dag > trace\_id 5270072548048800615

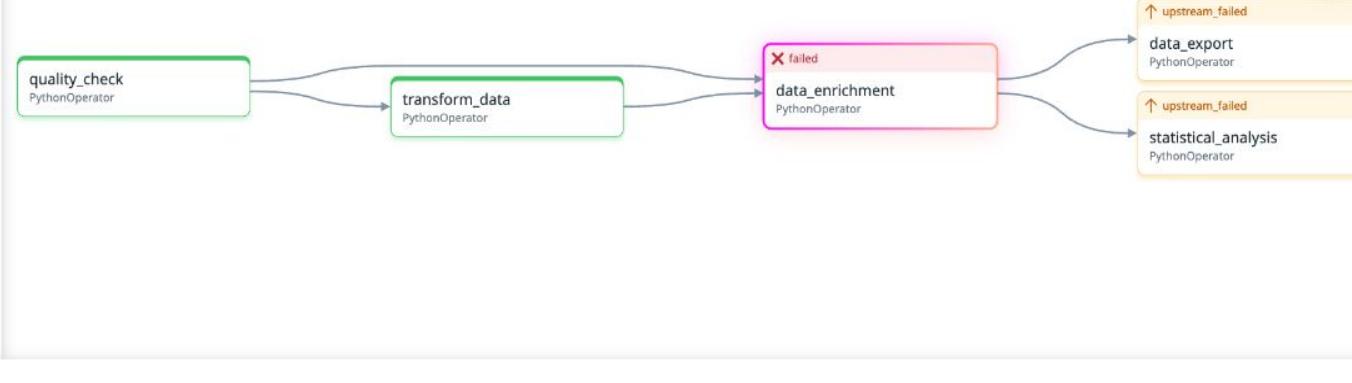
p87 16m 18s | Oct 08 11:30:00.544 (33m ago) | Airflow UI

Trace: Tasks Flame Graph Waterfall Span List 10 JSON

TASK ID	↓ STATUS	↓ <sup>2</sup> DOWNSTRE...	DURATION
unreliable_me...	FAILED	→ 7	6.50s
data_enrichm...	FAILED	→ 4	2.00μs
statistical Ана...	UPSTREAM_F...	→ 2	0ns
generate_repor...	UPSTREAM_F...	→ 2	0ns
data_export	UPSTREAM_F...	→ 2	0ns
collect_metrics	UPSTREAM_F...	→ 2	0ns
backup data	UPSTREAM_F...	→ 2	0ns

visualization\_dag airflow.task data\_enrichment

← View Full DAG Map | Showing direct dependencies



🕒 8.66s 0.88% total exec time

Span: Overview Errors 5 Logs 40

visualization\_dag > data\_enrichment

Pretty Raw

Fix With Bits

Enrichment failure occurred based on run minute timing

Traceback (most recent call last):

Show 6 third-party frames

> File "smoke\_testing/visualization\_dag.py", line 145, in enrichment\_task

Exception: Enrichment failure occurred based on run minute timing

visualization\_dag > data\_enrichment

Collapse

# Example of OpenLineage event

```
{
  "eventType": "COMPLETE",
  "eventTime": "2025-06-08T12:00:00.000Z",
  "run": {
    "runId": "123e4567-e89b-12d3-a456-426614174000",
    "facets": {
      "parent": {
        "run": {
          "runId": "111e4567-e89b-12d3-a456-426614174abc",
          "namespace": "my_company_etl"
        },
        "job": {
          "namespace": "my_company_etl",
          "name": "parent_job"
        }
      }
    }
  },
  "job": {
    "namespace": "my_company_etl",
    "name": "daily_sales_etl",
    "facets": {
      "documentation": {
        "description": "ETL job to aggregate daily sales for dashboard reporting"
      }
    }
  }
},
```

```
"inputs": [
  {
    "namespace": "my_data_warehouse",
    "name": "raw_sales_data",
    "facets": {
      "schema": {
        "fields": [
          {"name": "sale_id", "type": "string"},
          {"name": "amount", "type": "float"},
          {"name": "timestamp", "type": "timestamp"}
        ]
      }
    }
  }
],
"outputs": [
  {
    "namespace": "my_data_warehouse",
    "name": "aggregated_sales_data",
    "facets": {
      "outputStatistics": {
        "rowCount": 2500,
        "size": 102400,
        "columnStats": {
          "amount": {
            "max": 999.99,
            "min": 1.23,
            "nullCount": 0
          }
        }
      }
    }
  }
],
"producer": "https://mycompany.com/lineage/collector"
```

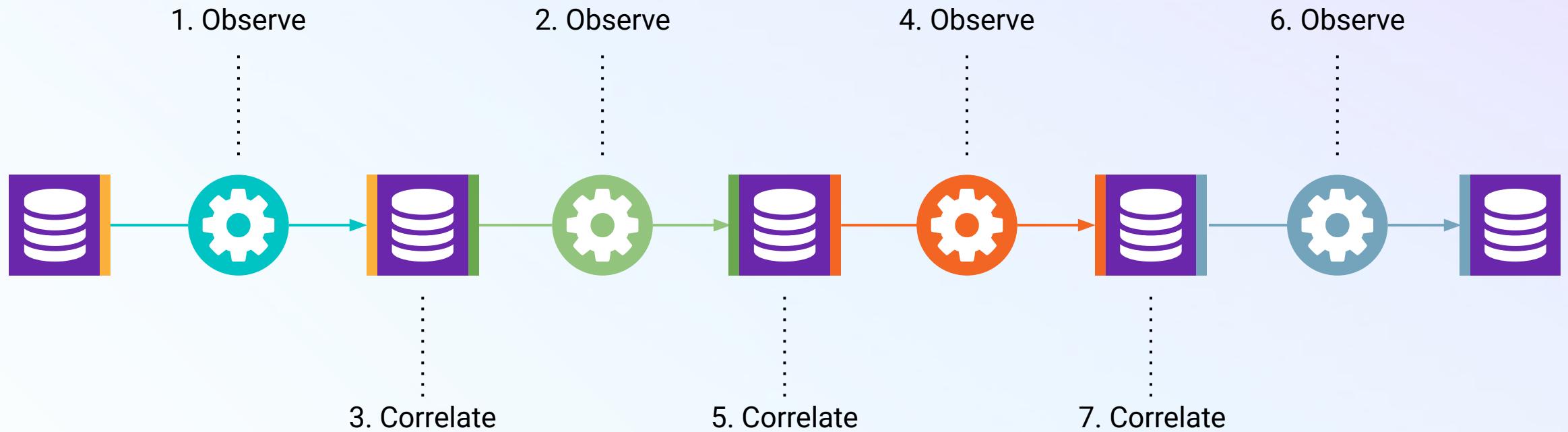
# Example of OpenLineage event

```
{  
  "eventType": "COMPLETE",  
  "eventTime": "2025-06-08T12:00:00.000Z",  
  "run": {  
    "runId": "123e4567-e89b-12d3-a456-426614174000",  
    "facets": {  
      "parent": {  
        "run": {  
          "runId": "111e4567-e89b-12d3-a456-426614174abc",  
          "namespace": "my_company_etl"  
        },  
        "job": {  
          "namespace": "my_company_etl",  
          "name": "parent_job"  
        }  
      }  
    }  
  },  
  "job": {  
    "namespace": "my_company_etl",  
    "name": "daily_sales_etl",  
    "facets": {  
      "documentation": {  
        "description": "ETL job to aggregate daily sales for dashboard reporting"  
      }  
    }  
  }  
},
```

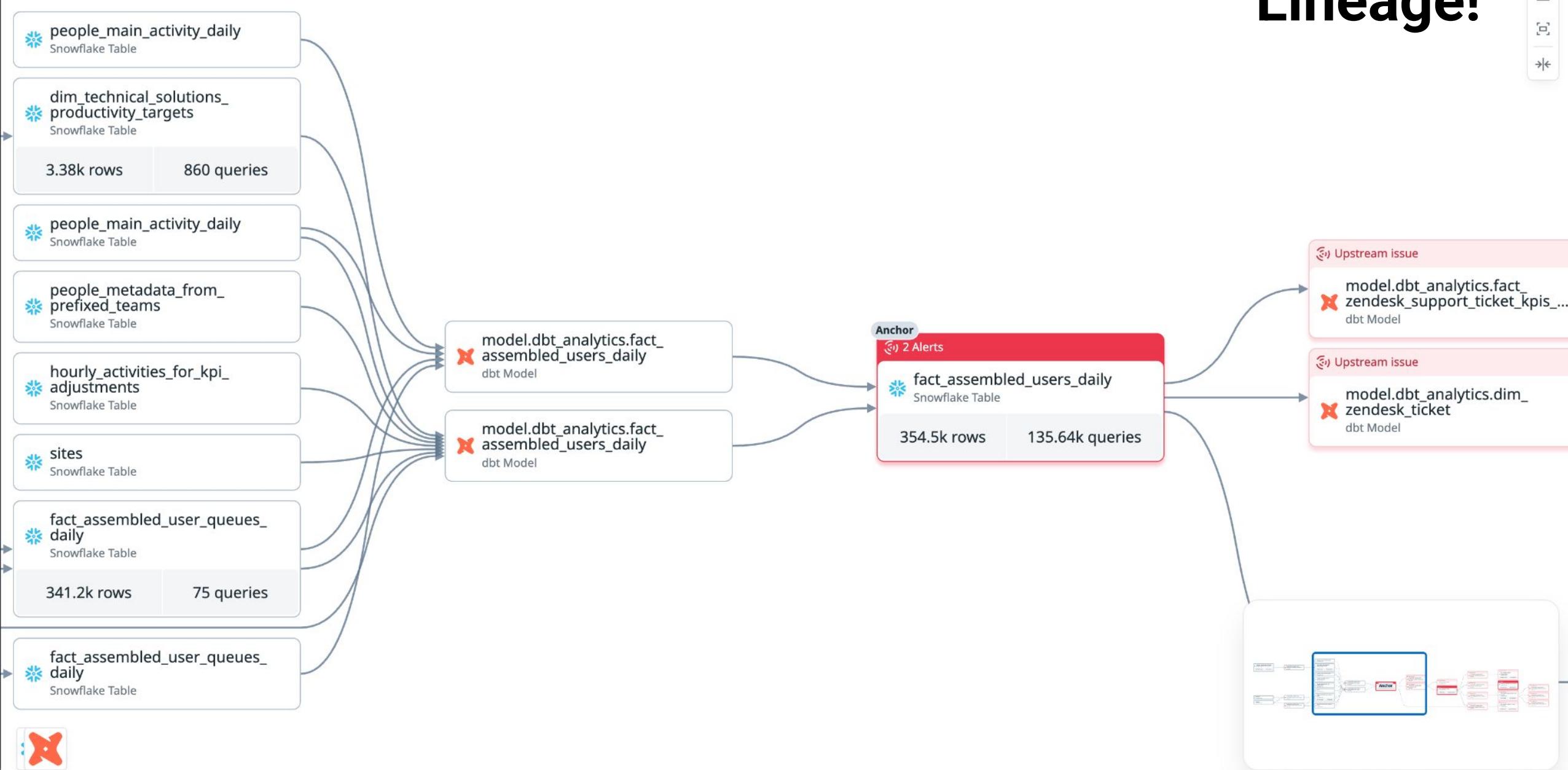
"inputs": [  
 {  
 "namespace": "my\_data\_warehouse",  
 "name": "raw\_sales\_data",  
 "facets": {  
 "schema": {  
 "fields": [  
 {"name": "sale\_id", "type": "string"},  
 {"name": "amount", "type": "float"},  
 {"name": "timestamp", "type": "timestamp"}  
 ]  
 }  
 }  
 }  
,  
 "outputs": [  
 {  
 "namespace": "my\_data\_warehouse",  
 "name": "aggregated\_sales\_data",  
 "facets": {  
 "outputStatistics": {  
 "rowCount": 2500,  
 "size": 102400,  
 "columnStats": {  
 "amount": {  
 "max": 999.99,  
 "min": 1.23,  
 "nullCount": 0  
 }  
 }  
 }  
 }  
 }  
,  
 "producer": "<https://mycompany.com/lineage/collector>"  
 }  
]

**OUTPUTS**

# Lineage is built on correlations



# Lineage!



# Example of OpenLineage event

```
{  
  "eventType": "COMPLETE",  
  "eventTime": "2025-06-08T12:00:00.000Z",  
  "run": {  
    "runId": "123e4567-e89b-12d3-a456-426614174000",  
    "facets": {  
      "parent": {  
        "run": {  
          "runId": "111e4567-e89b-12d3-a456-426614174abc",  
          "namespace": "my_company_etl"  
        },  
        "job": {  
          "namespace": "my_company_etl",  
          "name": "parent_job"  
        }  
      }  
    }  
  },  
  "job": {  
    "namespace": "my_company_etl",  
    "name": "daily_sales_etl",  
    "facets": {  
      "documentation": {  
        "description": "ETL job to aggregate daily sales for dashboard reporting"  
      }  
    }  
  },  
}
```

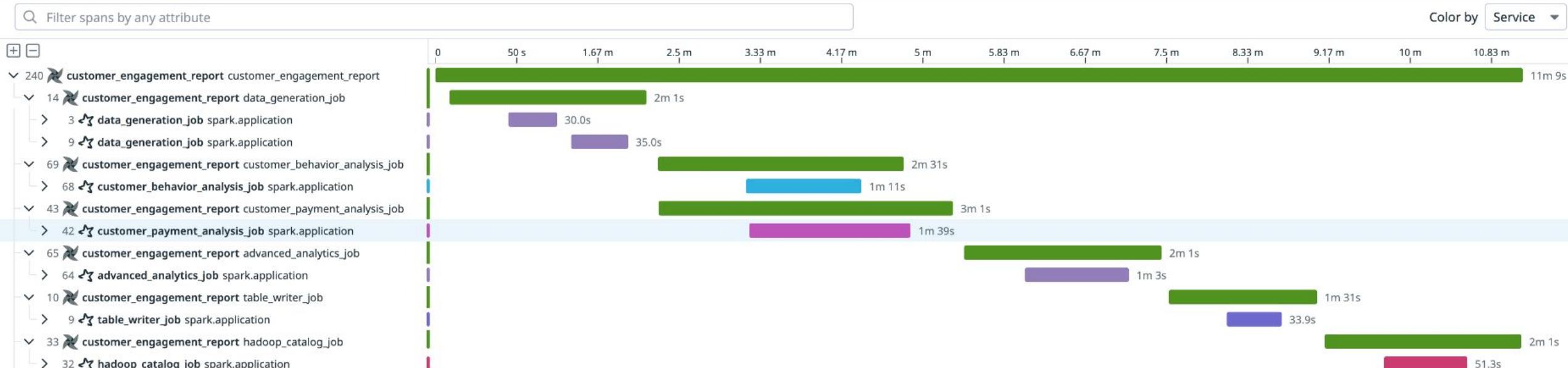
PARENT

```
"inputs": [  
  {  
    "namespace": "my_data_warehouse",  
    "name": "raw_sales_data",  
    "facets": {  
      "schema": {  
        "fields": [  
          {"name": "sale_id", "type": "string"},  
          {"name": "amount", "type": "float"},  
          {"name": "timestamp", "type": "timestamp"}  
        ]  
      }  
    }  
  }  
,  
  "outputs": [  
    {  
      "namespace": "my_data_warehouse",  
      "name": "aggregated_sales_data",  
      "facets": {  
        "outputStatistics": {  
          "rowCount": 2500,  
          "size": 102400,  
          "columnStats": {  
            "amount": {  
              "max": 999.99,  
              "min": 1.23,  
              "nullCount": 0  
            }  
          }  
        }  
      }  
    }  
,  
    "producer": "https://mycompany.com/lineage/collector"  
  }  
]
```

⌚ 11m 9s | Jun 07 15:00:00.866 (44m ago)

Trace: Flame Graph Waterfall Span List 241

Filter spans by any attribute

 customer\_payment\_analysis\_job spark.application spark.application ⋮

⌚ 1m 39s 14.8% total exec time

Span: Overview Infrastructure Logs 50+ User Outputs 50+ Configuration Dev Agent

## Executors CPU Breakdown



## Input / Output

## INPUT

Size 7.72 KiB

Records 2.21k

## OUTPUT

Size 7.72 KiB

Records 105

## SHUFFLE

Size 96.4 KiB

Records 815

## Memory usage

⌚ 11m 9s | Jun 07 15:00:00.866 (44m ago)

Trace: Flame Graph Waterfall Span List 241

Filter spans by any attribute

Color by Service Hide Legend

# OpenLineage parent id



customer\_payment\_analysis\_job spark.application spark.application ::

⌚ 1m 39s 14.8% total exec time

Span: Overview Infrastructure Logs 50+ User Outputs 50+ Configuration Dev Agent

- spark\_openlineage\_parentJobName
- spark\_openlineage\_parentJobNamespace
- spark\_openlineage\_parentRunId
- spark\_openlineage\_rootParentJobName
- spark\_openlineage\_rootParentJobNamespace
- spark\_openlineage\_rootParentRunId
- spark\_org\_apache\_hadoop\_yarn\_server\_webproxy\_amfilter\_AmIpFilter\_param\_PROXY\_HOSTS
- spark\_org\_apache\_hadoop\_yarn\_server\_webproxy\_amfilter\_AmIpFilter\_param\_PROXY\_URI\_BASES
- spark\_resourceManager\_cleanupExpiredHost
- spark\_shuffle\_service\_enabled
- spark\_sql\_catalog\_prod\_catalog

- customer\_engagement\_report.customer\_payment\_analysis\_job
- demo-env
- 01974b8c-ad00-7c2f-a278-260f9096e0b1
- customer\_engagement\_report
- demo-env
- 01974b8c-ad00-7c02-a664-6552cc96d2ce
- ip-192-168-17-253.ec2.internal
- [http://ip-192-168-17-253.ec2.internal:20888/proxy/application\\_1747825417854\\_2386](http://ip-192-168-17-253.ec2.internal:20888/proxy/application_1747825417854_2386)
- true
- true
- org.apache.iceberg.spark.SparkCatalog

# OMG the possibilities are endless

- Dependency tracing
- Root cause identification
- Issue prioritization
- Impact mapping
- Precision backfills
- Anomaly detection
- Change management
- Historical analysis
- Compliance



# Problem fixed!

The screenshot shows the homepage of the Canoe travel website. At the top, there is a navigation bar with the word "Canoe" on the left and links for "Hotels", "Flights", "Car Rentals", "Things to Do", and "Sign in". Below the navigation is a large banner with the text "Find the best deals on hotels and flights" overlaid on a background image of a coastal town with orange-roofed buildings and hills. A search form is centered in the banner, featuring fields for "Destination" and "Check in" and a blue "Search" button. Below the banner, the page title "Featured Hotels" is displayed, flanked by two yellow sun icons. There are four hotel cards visible: "Oceanview Inn" in Miami Beach, FL, "Grand Hotel" in Paris, France, "Mountain Lodge" in Aspen, Colorado, and "Cityscape Hotel" in New York, NY. Each card includes a small thumbnail image of the hotel, the name, location, a five-star rating icon, and the price per night.

Canoe

Hotels Flights Car Rentals Things to Do Sign in

Find the best deals on hotels and flights

Destination Check in Search

Featured Hotels

Oceanview Inn  
Miami Beach, FL  
★★★★★  
\$210 /night

Grand Hotel  
Paris, France  
★★★★★  
\$350 /night

Mountain Lodge  
Aspen, Colorado  
★★★★★  
\$250 /night

Cityscape Hotel  
New York, NY  
★★★★★  
\$450 /night

# How Can OpenLineage Benefit You?

- Identify your key datasets and jobs
- Instrument with OpenLineage: Airflow, Spark, dbt, Flink, etc.
- Profit!
  - ... or at least have trust in your data!

# Thank you!

OpenLineage

[github.com/OpenLineage](https://github.com/OpenLineage)

[OpenLineage.io](https://OpenLineage.io)

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

[docs.datadoghq.com/data\\_observability](https://docs.datadoghq.com/data_observability)