# 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**
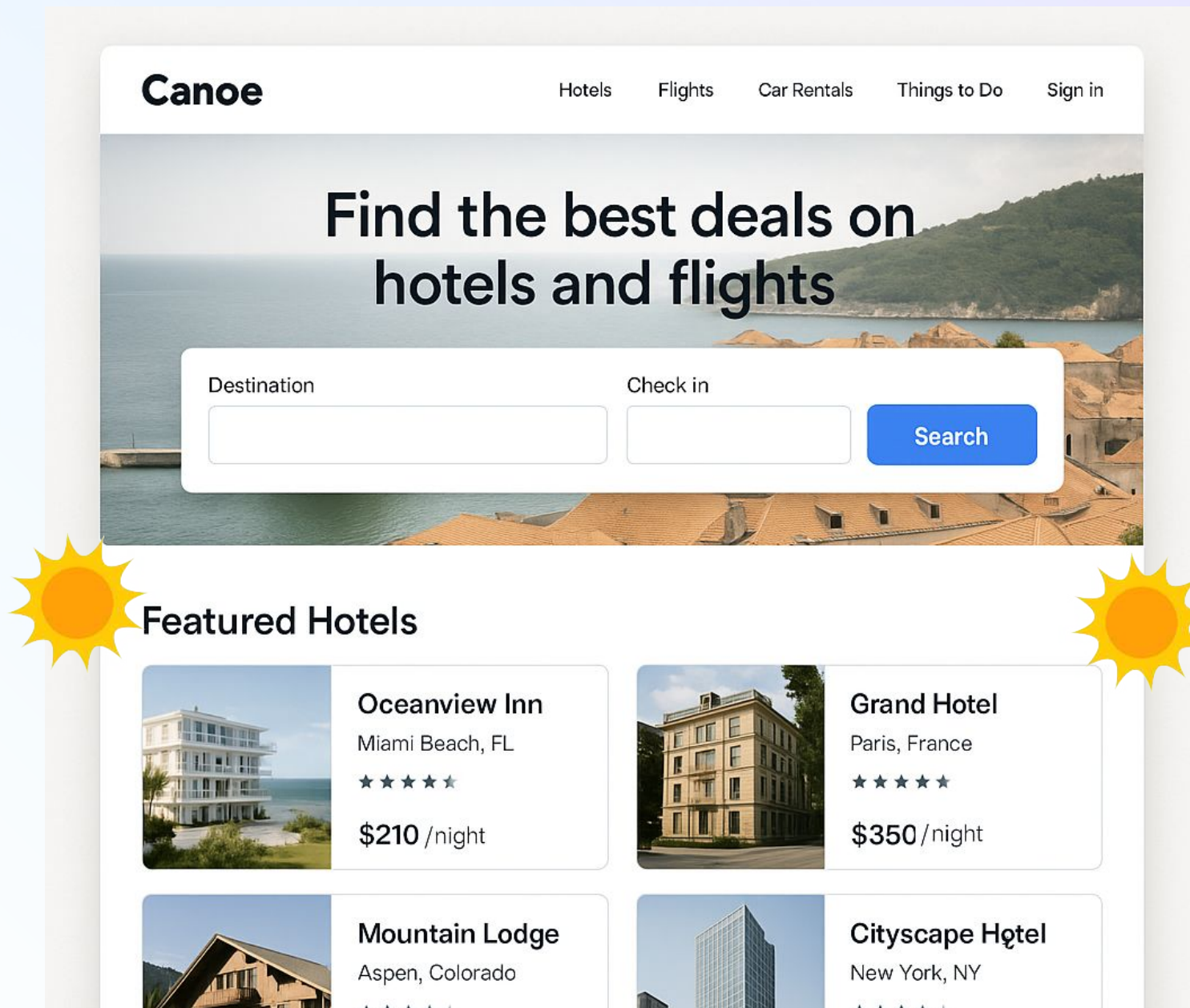Senior Engineering Manager
OpenLineage TSC

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

DATADOG

# Imagine we work at Canoe

# But something is wrong...

# But something is wrong...

# Data ecosystem at Canoe

# Data ecosystem at Canoe

# Data ecosystem at Canoe

# Data ecosystem at Canoe



Org A

Integrations team

Org B

ML team

App team

Org C

canoe

Top Hotels

Hotel Rivera

The Grand Lodge

Hotel Panorama

# Data ecosystem at Canoe

# And the ecosystem is fragmented

# We need everyone to speak the same language

# Need a common language

Metadata about data pipelines:

- where data is from

- how it changes

- where it goes

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

{} {} {}

CRAWL    CRAWL    CRAWL

PARSE    PARSE    PARSE

DATA LINEAGE REPOSITORY

Integrate with source code repositories

Look for queries and parse them for lineage

Report to a lineage metadata repository

# Process query or activity logs

POLL · POLL · POLL · POLL · POLL

DATA LINEAGE REPOSITORY

- 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

DATA LINEAGE 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



Without OpenLineage

LINEAGE PRODUCERS

LINEAGE CONSUMERS

# Simplified



Without OpenLineage

With OpenLineage

LINEAGE PRODUCERS

OpenLineage

LINEAGE CONSUMERS

# OpenLineage at Datadog

# Datadog Data Observability

Data Jobs Monitoring

Data Quality Monitoring

# OpenLineage benefits

# OpenLineage Spec deep dive

# Core Concepts: Job, Run, Dataset

| Job |
|---|
| **job id**<br>**(name based)** |

| Run |
|---|
| **run uuid** |

| Dataset |
|---|
| **dataset id**<br>**(name based)** |

# Core Concepts: Facets

| Job |
|-----|
| **Job Id**<br>**(name based)** |

| Run |
|-----|
| **run uuid** |

| Dataset |
|---------|
| **dataset Id**<br>**(name based)** |

↑

↑

↑

| Job Facet |
|-----------|

| Run Facet |
|-----------|

| Dataset Facet |
|---------------|

# Events

```
                    ┌─────────────────┐
                    │   Run Event     │
                    ├─────────────────┤
                    │ transition time │
                    │                 │
                    └────────┬────────┘
                             │
       ┌─────────────────────┼─────────────────────┐
       ▼                     ▼                     ▼
┌─────────────┐      ┌─────────────┐       ┌─────────────┐
│     Job     │      │     Run     │       │   Dataset   │
├─────────────┤      ├─────────────┤       ├─────────────┤
│   Job id    │      │   run uuid  │       │  dataset id │
│(name based) │      │             │       │(name based) │
│             │      │             │       │             │
└──────▲──────┘      └──────▲──────┘       └──────▲──────┘
       │                    │                     │
┌─────────────┐      ┌─────────────┐       ┌─────────────┐
│  Job Facet  │      │  Run Facet  │       │Dataset Facet│
└─────────────┘      └─────────────┘       └─────────────┘
```

# Run Event Lifecycle



START
runID
eventType: START
event time
producer
datasets

RUNNING
runID
eventType: RUNNING
event time
producer
specific facets

ABORT

COMPLETE
runID
eventType: COMPLETE
event time
producer
datasets

FAIL

# 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": {                                    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"
}
```

Integrations

UTC-07:00
1d   Past 1 Day

# product_revenue_report

Env
All

Status
All

+ Filter

### Summary

### Job Runs

### Tasks

## ⊕ Summary

⟳ **ALL MONITORS**   **2 ALERT**   **3 OK**   + Add

Latest run completed **4m 24s ago**

▤ **288 Total Runs**

0% Failed

Average run duration is **21.3s**

**4% ↘**

Compared to previous period

### Runs



Ok

### Average Run Duration



Avg. Duration

## ▤ 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

```json
{
  "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"
      }
    }
  },
```
```json
  "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"
}
```

**RUN**

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 |

Collapse

---

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

Showing **1–18** of **18 tasks** 🔍 Search

← View Full DAG Map | Showing direct dependencies

| TASK ID | STATUS | DOWNSTRE... | DURATION |
|---|---|---|---|
| 🐍 unreliable_me... | **FAILED** | → 7 | 6.50s |
| 🐍 data_enrichm... | **FAILED** | → 4 | 2.00μs |
| 🐍 statistical_ana... | **UPSTREAM_F...** | → 2 | 0ns |
| 🐍 generate_rep... | **UPSTREAM_F...** | → 2 | 0ns |
| 🐍 data_export | **UPSTREAM_F...** | → 2 | 0ns |
| 🐍 collect_metrics | **UPSTREAM_F...** | → 2 | 0ns |
| 🐍 backup_data | **UPSTREAM_F...** | → 2 | 0ns |

quality_check
PythonOperator

transform_data
PythonOperator

✕ failed
data_enrichment
PythonOperator

↑ upstream_failed
data_export
PythonOperator

↑ upstream_failed
statistical_analysis
PythonOperator

---

visualization_dag *airflow.task* 🛑 data_enrichment ⋮

⏱ 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

# Example of OpenLineage event

```json
{
  "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"
      }
    }
  },
```

```json
  "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"
}
```

**INPUTS**

# Example of OpenLineage event

```json
{
  "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"
      }
    }
  },
```

```json
  "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"
}
```
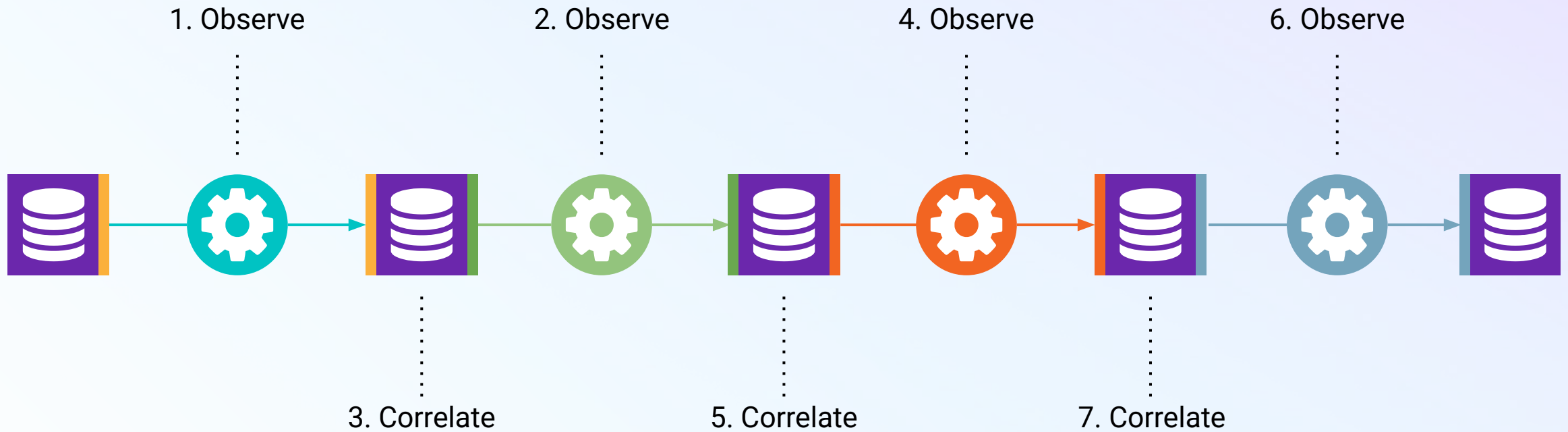
**OUTPUTS**

# Lineage is built on correlations

1. Observe

2. Observe

4. Observe

6. Observe

3. Correlate

5. Correlate

7. Correlate

DATADOG

# Example of OpenLineage event

```json
{
  "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**

```json
  "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"
}
```

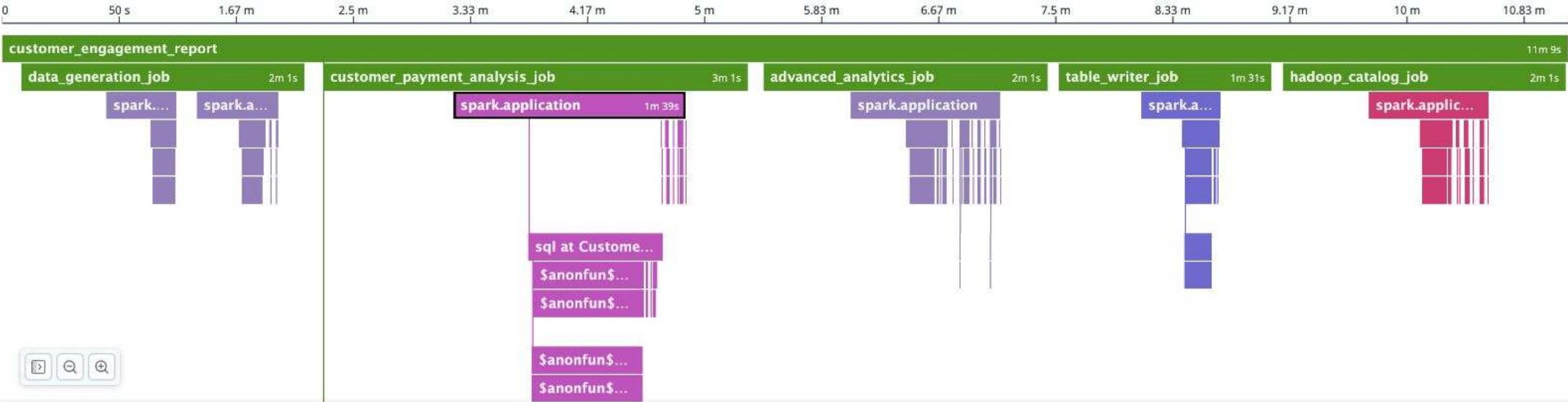# OpenLineage parent id

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 ▸

| 0 | 50 s | 1.67 m | 2.5 m | 3.33 m | 4.17 m | 5 m | 5.83 m | 6.67 m | 7.5 m | 8.33 m | 9.17 m | 10 m | 10.83 m |

% Exec Time ▾

| customer_engagement_report | | | | 11m 9s |

- customer_enga...  54.6%
- data_generatio...  9.72%
- advanced_anal...  9.56%
- customer_pay...  8.61%
- hadoop_catalo...  7.66%
- table_writer_job  5.07%
- customer_beha...  4.77%

data_generation_job 2m 1s   customer_payment_analysis_job 3m 1s   advanced_analytics_job 2m 1s   table_writer_job 1m 31s   hadoop_catalog_job 2m 1s

spark...  spark.a...   spark.application 1m 39s   spark.application   spark.a...   spark.applic...

sql at Custome...

$anonfun$...
$anonfun$...

$anonfun$...

$anonfun$...

---

🎷 customer_payment_analysis_job  *spark.application*  spark.application  ⋮          ⏱ 1m 39s    14.8% total exec time

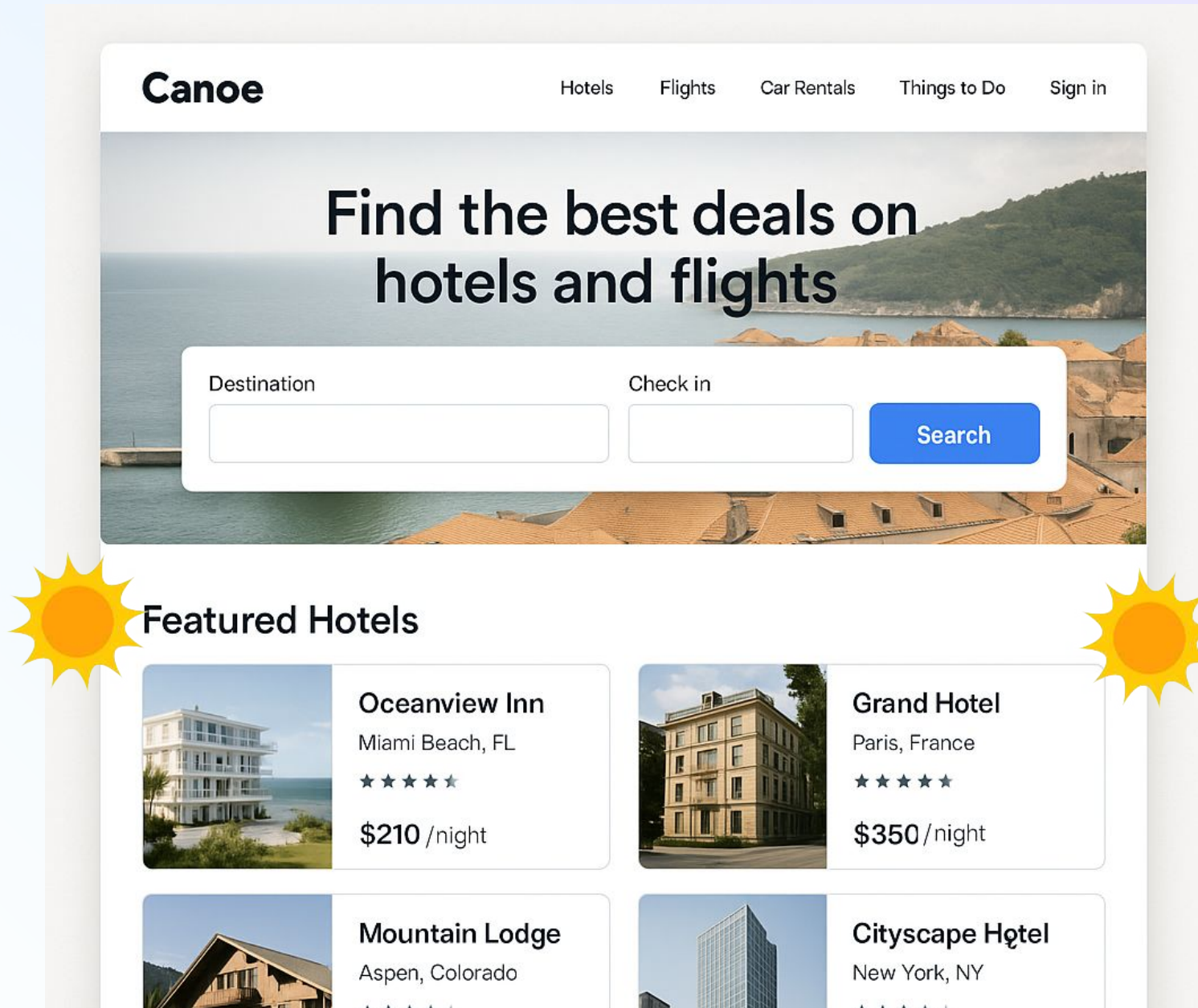Span:  Overview   Infrastructure   Logs 50+   User Outputs 50+   Configuration   Dev Agent

| spark_openlineage_parentJobName | customer_engagement_report.customer_payment_analysis_job |
| spark_openlineage_parentJobNamespace | demo-env |
| spark_openlineage_parentRunId | 01974b8c-ad00-7c2f-a278-260f9096e0b1 |
| spark_openlineage_rootParentJobName | customer_engagement_report |
| spark_openlineage_rootParentJobNamespace | demo-env |
| spark_openlineage_rootParentRunId | 01974b8c-ad00-7c02-a664-6552cc96d2ce |
| spark_org_apache_hadoop_yarn_server_webproxy_amfilter_AmIpFilter_param_PROXY_HOSTS | ip-192-168-17-253.ec2.internal |
| spark_org_apache_hadoop_yarn_server_webproxy_amfilter_AmIpFilter_param_PROXY_URI_BASES | http://ip-192-168-17-253.ec2.internal:20888/proxy/application_1747825417854_2386 |
| spark_resourceManager_cleanupExpiredHost | true |
| spark_shuffle_service_enabled | true |
| spark_sql_catalog_prod_catalog | 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!

# 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

OpenLineage.io

@OpenLineage

docs.datadoghq.com/data_observability