# Data Observability
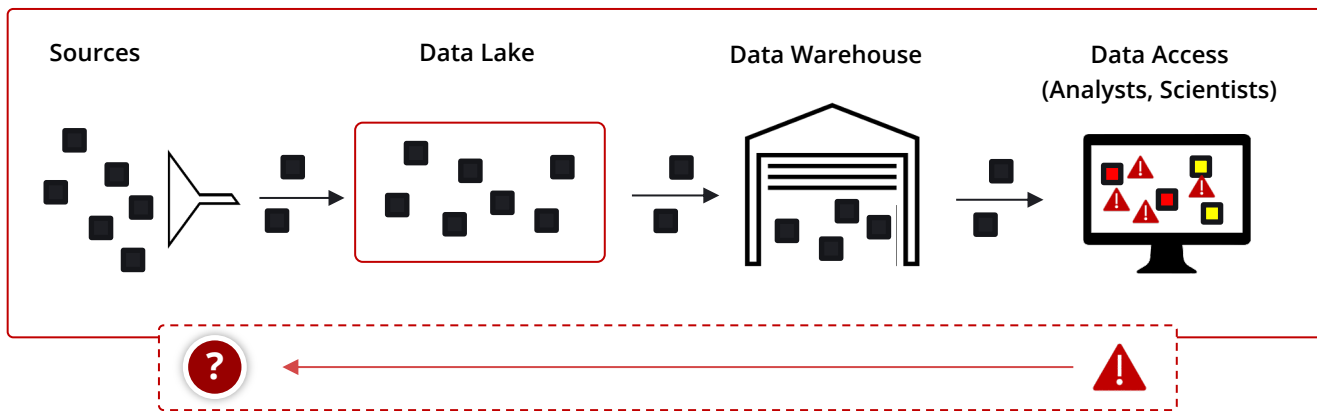
## Technical Overview

# What is Data Observability?

# Problem Statement

**Data engineers are reactive to data issues**

Sources          Data Lake          Data Warehouse          Data Access
(Analysts, Scientists)

Many data quality issues are **overlooked**

Platform only learns about issues when **reported by data consumers**

After issues are reported, they **are not resolved quickly**
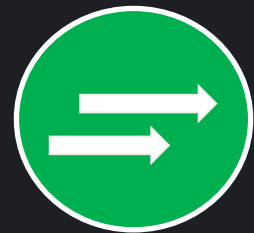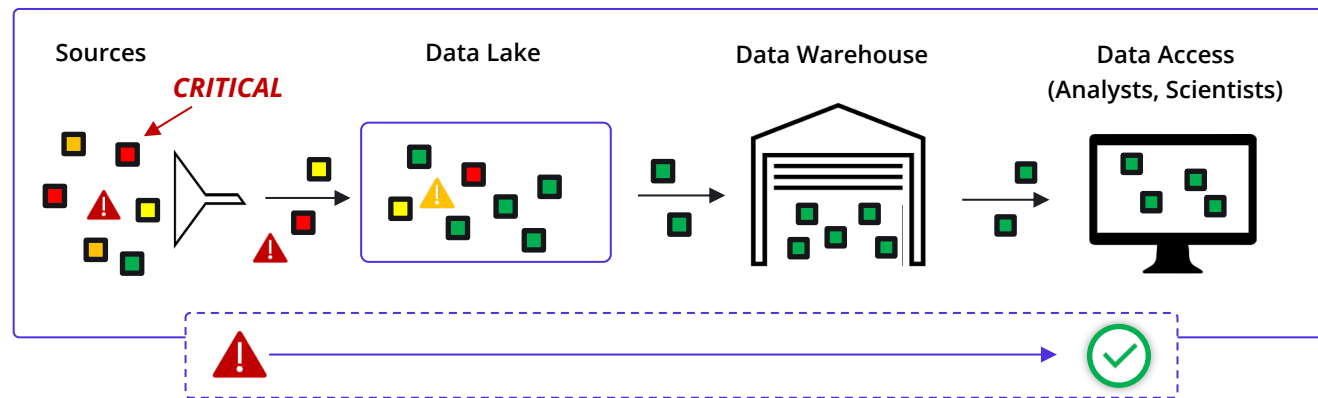
Databand

## The root causes

Fragmented toolchain

Volume of data

Flooding of noise

# Solution

**Proactive Data Observability:** *Shift left and solve problems at the source*



| Improve **MTTD** | Improve **MTTR** | Improve data product **quality** |
|---|---|---|
| *Discover issues in real time, early as ingest* | *Identify the cause of issues instantly* | *Enhance trust and consumer satisfaction* |

Our solution focuses on observing **data in motion**

# Solution

**Proactive Data Observability:** *Shift left and solve problems at the source*



Databand focuses on observing **data in motion**

Observe data pipeline **process quality**
- Status
- Performance
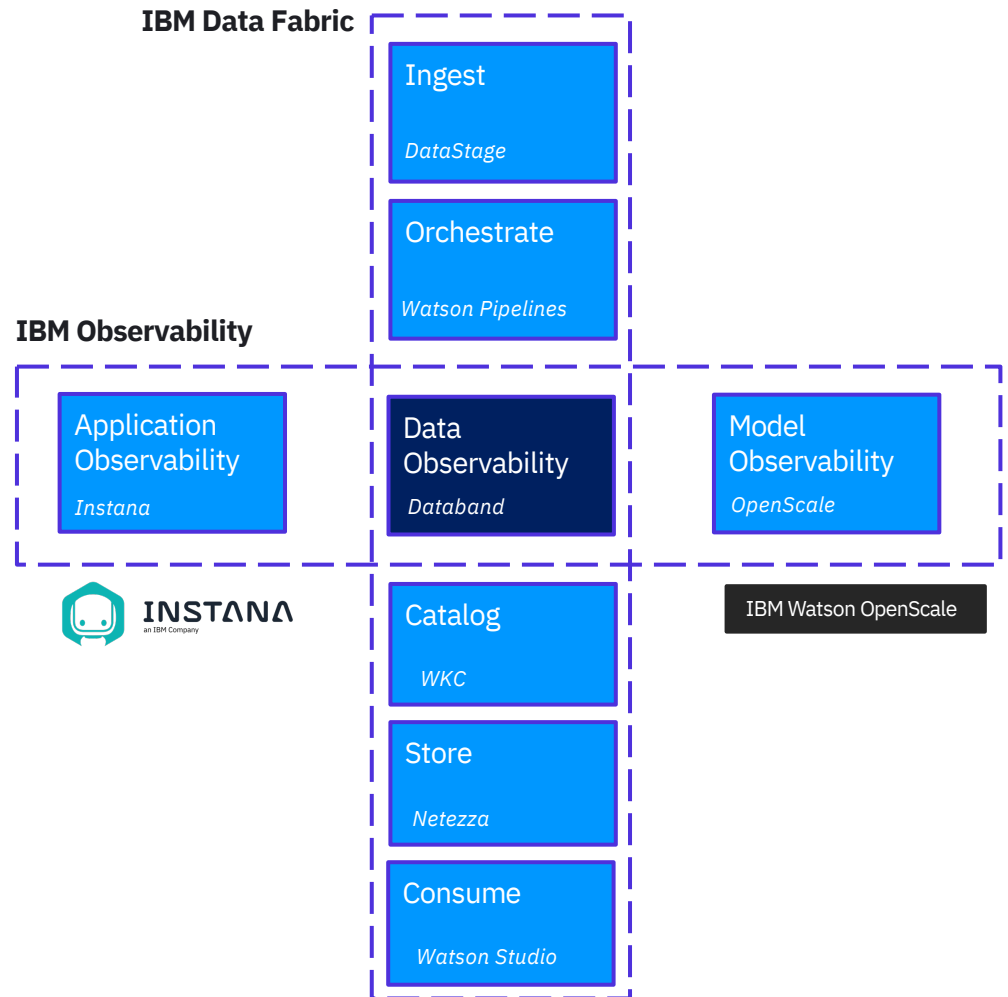- Latency

Observe **data quality and reliability**
- *Schema changes*
- *Data shape*
- *Data freshness*

Supported **data pipelines** and **workflow managers**
- *Airflow*: Python, Spark, dbt, SQL and other operators

The Databand bridges **two** of IBM's strategic directions, **IBM Data Fabric** and **IBM Observability**.

Bridging these strategies unlocks powerful new use cases for customers and growth for IBM.

**IBM Data Fabric**

**Ingest**
*DataStage*

**Orchestrate**
*Watson Pipelines*

**IBM Observability**

**Application Observability**
*Instana*

**Data Observability**
*Databand*

**Model Observability**
*OpenScale*

INSTANA
an IBM Company

**Catalog**
*WKC*

IBM Watson OpenScale

**Store**
*Netezza*

**Consume**
*Watson Studio*

# Databand integration

| Pipeline implementation and deployment | Databand integration | Additional details |
|---|---|---|
| Spark | No code, optional SDK | Configuration in a Spark cluster – Databand provides a listener |
| Airflow (all pipelines) | No code for pipeline status, SDK for dataset monitoring | |
| Python, PySpark, Java, Scala without an workflow engine | SDK for pipeline status and dataset monitoring | |
| dbt | No code, optional SDK | Syncer for dbt Cloud for a no-code integration. Python SDK can be used to retrieve information about a specific job run |
| DataStage * (Q4) | No code | Supported for *DataStage Next Gen* |

# Data Pipelines and Workflow managers

# Data pipelines

- *Data pipeline* is a generic term that describes the process of moving data between data sources

  - While in most cases data pipelines performing ETL tasks, a data pipelines can move data without transformations

- Data pipelines can be implemented in a variety of programing languages, technologies, and tools

  1. Languages: *Python, Java, Scala, SQL*
  2. Technologies: *Spark*
  3. Tools: *dbt, DataStage, Azure Data Factory, AWS Data Pipeline, and others*

# Workflow engines

- *Workflow* engines are used to orchestrate execution of tasks

  - Many types of workflows are supported by generic *operators*, not just ETL

    - Examples of workflows: ETL, MLOps, DevOps

- Examples of workflow engines

  - *Airflow, Azkaban, MLFlow, Kubeflow, Luigi,* and others
  - *Airflow* is one of the first and one of the most generic workflow engines



- *Airflow* is one of the first (open source) and one of the most generic workflow engines
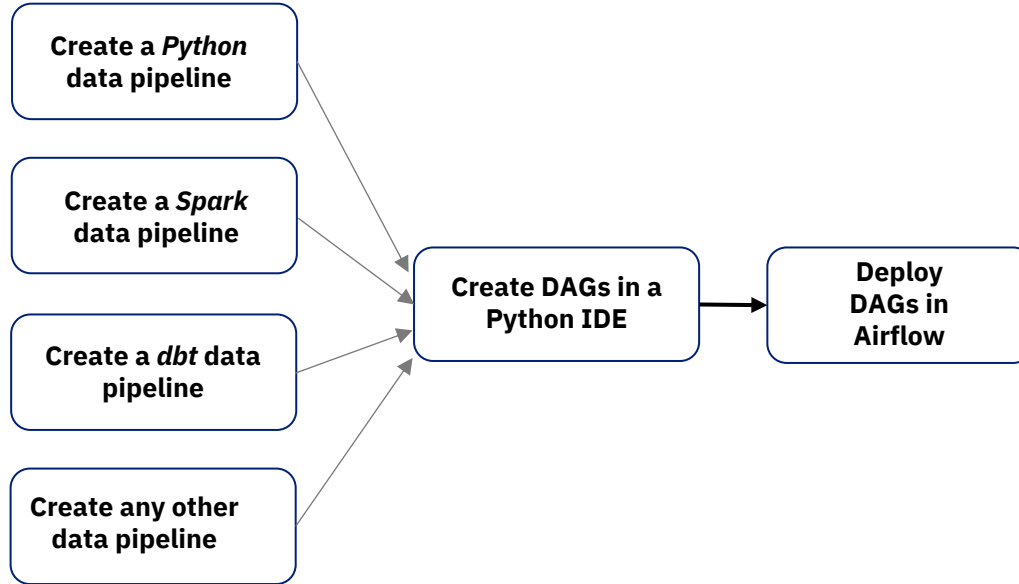
# Apache Airflow

- **Important concepts**

  - *DAGs (*Directed Acyclic Graph)
    - A collection of the tasks in a job with relationships and dependencies
    - A DAG is defined in a Python script

  - *Operators*
    - Pre-built functions for frequently used tasks: Python, bash, SQL, and others

  - *Admin console*
    - Understand the features of the admin console



*For Databand integration, all DAGS that we are discussing contain are **data pipelines** (and not other types of pipelines)*

# Development and deployment of data pipelines

# Development and deployment of data pipelines

**Create a *Python* data pipeline**

**Create a *Spark* data pipeline**

**Create a *dbt* data pipeline**

**Create any other data pipeline**

**Create DAGs in a Python IDE**

**Deploy DAGs in Airflow**

*Databand*

# Development and deployment of data pipelines

**Create a *Python* data pipeline**

**Create a *(PySpark* data pipeline**

**Create a *dbt* data pipeline (model)**

**Create Java or Scala data pipeline**

**Add Databand API**

*Databand*

# Development and deployment of data pipelines

**1**

**Create a *Python* data pipeline**

**Create a *Spark* data pipeline**

**Create a *dbt* data pipeline**

**Create any other data pipeline**
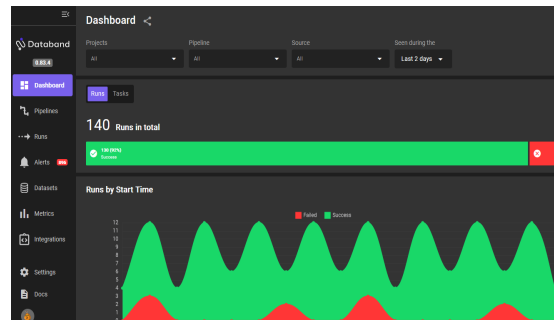
*Develop*

**2**

**Create DAGs in a Python IDE** → **Deploy DAGs in Airflow**

**and/or**

**Add Databand API**

*Develop*

**3**

*Databand*



*Monitor*

# The Databand.ai solution



**1**
**Automatically collect metadata**
From key solutions in the modern data stack.

**2**
**Build historical baseline**
Based on common data pipeline behavior.

**4**
**Resolve through automation**
Create smart workflows to remediate data quality issues and keep SLAs on track.

**3**
**Alert on anomalies and rules**
Based on deviations or breaches.

Collect

Profile

Alert

Resolve

# Dashboard

| Projects | Pipeline | Source | Seen during the |
|---|---|---|---|
| All | All | All | Last 2 days |

## 98 Runs in total

90 (91%)
Success

90 (91%) Success

## 296 Task Runs

268 (90%)
Success

### Runs by Start Time

■ Failed ■ Success

9
8
7
6
5
4
3
2
1
0

Mar 20 02:00 PM
Mar 20 05:00 PM
Mar 20 08:00 PM
Mar 20 11:00 PM
Mar 21 02:00 AM
Mar 21 05:00 AM
Mar 21 08:00 AM
Mar 21 11:00 AM
Mar 21 02:00 PM
Mar 21 05:00 PM
Mar 21 08:00 PM
Mar 21 11:00 PM
Mar 22 02:00 AM
Mar 22 05:00 AM
Mar 22 08:00 AM
Mar 22 11:00 AM
Mar 22 02:00 PM

### Tasks by Start Time

■ Failed ■ Success ■ Upstream ■ Skipped

28
26
24
22
20
18
16
14
12
10
8
6
4
2
0

Mar 20 02:00 PM
Mar 20 05:00 PM
Mar 20 08:00 PM
Mar 20 11:00 PM
Mar 21 02:00 AM
Mar 21 0...

## Reliability Dashboard

## Top Errors

## Metrics

Add New Metric

NYC 311 API_read_rows  (get_hourly_data)

service_311_get_data  >  Service 311  databand-internal-sa-demo-af

Count: 46  ↑ 0%

Mar 21, 2001 02:00:16 PM
max: 17031.25
min: 12467.5
0

40000
35000
30000
25000
20000
15000
10000
5000
0
-5000

4PM 5PM 6PM 7PM 8PM 9PM 10PM 11PM 12AM 1AM 2AM 3AM 4AM 5AM 6AM 7AM 8AM 9AM 10AM 11AM 12PM 1PM 2PM 3PM 4PM 5PM 6PM 7PM 8PM 9PM 10PM 11PM 12AM 1AM 2AM 3AM 4AM 5AM 6AM 7AM 8AM 9AM 10AM 11AM 12PM 1PM 2PM 3PM

## Last Active

Runs | Tasks

All runs

| Status | Run name | Pipeline | Source | Project | | |
|--------|----------|----------|--------|---------|---|---|
| ✓ | ⇢ trig__2022-03-22T18:00:00+00:00 > | service_311_closed_requests > | Airflow | Service 311 | | |
| ✓ | ⇢ scheduled__2022-03-22T18:00:00+00:00 > | service_311_get_data > | Airflow | Service 311 | 1 m, 43 s | ✓ 6 Today at 02:00:00 PM |

**Detect Data Anomalies**

0.67.2

Details  Metrics  Logs  Code  Run Info  Affected Datasets  Histograms  ⧉

AIRFLOW

## ⊘ Error Logs ⧉

```
Traceback (most recent call last):
  File "/opt/airflow/dags/repo/dbnd-demo-airflow/dags/demo/service_311/modules/get_data.py", li
    hourly_data = api.api_read_to_df(NYC_DATA_BASE_URL, params)
  File "/opt/airflow/dags/repo/dbnd-demo-airflow/dags/demo/service_311/op_functions/api_operati
    response.raise_for_status()
  File "/home/airflow/.local/lib/python3.7/site-packages/requests/models.py", line 941, in rais
    raise HTTPError(http_error_msg, response=self)
requests.exceptions.HTTPError: 404 Client Error: Not Found for url: https://data.cityofnewyork.
```
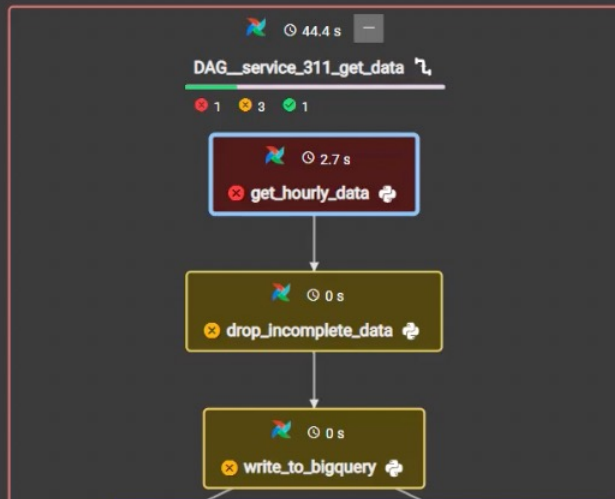
🦋 ⏱ 44.4 s  ⊟

DAG__service_311_get_data ⤵

⊗ 1   ⊗ 3   ✓ 1

🦋 ⏱ 2.7 s

⊗ get_hourly_data 🔁

🦋 ⏱ 0 s

⊗ drop_incomplete_data 🔁

🦋 ⏱ 0 s

⊗ write_to_bigquery 🔁

⊗ verif...

## Investigate Error Logs

### User Params

| templates_dict | None |
|---|---|
| op_args.0 | 'dbnd-demo-service311' |
| op_args.1 | 'data/databand-internal-sa-demo-af/' |
| op_args.2 | '2022-03-22' |
| op_args.3 | |
| function_name | get_hourly_data |

### Basic Params

| task_version | |
|---|---|

# Alerts

Add New Receiver 1    Add Alert    Search

🔔 Alerts 1787    🔔 11 Alerts Defined

| Projects | Pipeline | Source | Status | Severity | Triggered Time |
|---|---|---|---|---|---|
| All | All | All | Triggered ✕ | All | Last 7 days |

**119 Alerts**  Acknowledge  Resolve

| | | Severity | Description | Trigger Value | Origin | Time Triggered ↓ | Status |
|---|---|---|---|---|---|---|---|
| ☐ | ⋮ | ⚠ HIGH | Missing Dataset Operation in the Run… | 2 | 🗄 Multiple datasets 👁<br>⤷ service_311_closed_requests ›<br>⟶ trig__2022-03-22T17:00:00+00:00 › | Today at 02:02:22 PM | ● Triggered |
| ☐ | ⋮ | ⚠ HIGH | Missing Dataset Operation in the Run `service_311_get_data` | 5 | 🗄 Multiple datasets 👁<br>⤷ service_311_get_data ›<br>⟶ scheduled__2022-03-22T17:00:00+00:00 › | Today at 02:01:22 PM | ● Triggered |
| ☐ | ⋮ | 🔥 CRITICAL | Run Entered State: failed (Auto Alert)<br>Run Entered State: failed (Auto Alert) | failed | ⤷ service_311_get_data ›<br>⟶ scheduled__2022-03-22T17:00:00+00:00 › | Today at 02:01:17 PM | ● Triggered |
| ☐ | ⋮ | ⚠ HIGH | Missing Dataset Operation in the Run… | 1 | 🗄 BQ - Closed Requests ›<br>⤷ service_311_closed_requests ›<br>⟶ trig__2022-03-22T09:00:00+00:00 › | Today at 06:03:23 AM | ● Triggered |
| ☐ | ⋮ | 🔥 CRITICAL | Run Entered State: failed | failed | ⤷ service_311_closed_requests ›<br>⟶ trig__2022-03-22T09:00:00+00:00 › | | ● Triggered |
| ☐ | ⋮ | 🔥 CRITICAL | Run Entered State: failed (Auto Alert) | failed | ⤷ service_311_get_data ›<br>⟶ scheduled__2022-03-22T09:00:00+00:00 | | ● Triggered |
| ☐ | ⋮ | ⚠ HIGH | Run Entered State: failed | failed | ⤷ service_311_gcp_ingest_data › | Yesterday at 10:40:14 PM | ● Triggered |

Records per page: 25    1-25 of 119  |< ‹ › >|

# Data Schema

| Number of columns | 40 | Number of records | 14,909 | Size (Bytes) | 596,360 |

Dataset Name: GCS - Cleaned Hourly Data    Pipeline: service_311_get_data    Operation: Read

Type: Pandas.DataFrame    Run: scheduled__2022-03-22T07:00:00+00:00    Operation Status: Success

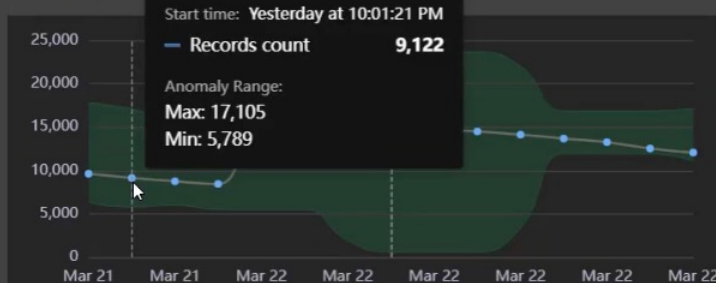| Column name | Field type |
| --- | --- |
| unique_key | int64 |

| Records count: | 14,909 |
| Mean: | 53,659,237.0048 |
| STD: | 296,299.4752 |
| Distinct values: | 14,909 |
| Null count: | 0 |
| Non-null count: | 14,909 |
| Null percentage: | 0% |
| Min: | 45,597,300 |
| Max: | 53,702,070 |
| 25%: | 53,689,335 |
| 50%: | 53,693,726 |
| 75%: | 53,697,938 |

Records count:

Run's history

Start time: Yesterday at 10:01:21 PM
— Records count    9,122
Anomaly Range:
Max: 17,105
Min: 5,789

25,000
20,000
15,000
10,000
5,000
0

Mar 21  Mar 21  Mar 22  Mar 22  Mar 22  Mar 22  Mar 22

| created_date | object |
| closed_date | object |
| agency | object |
| agency_name | object |

Details    AIRFLOW

Understand Schema Changes