
Beyond Testing: How to Build Circuit Breakers with Airflow

Prateek Chawla

MC

Agenda

- The Rise of Data Downtime
- The Data Reliability Workflow
- Introduction to Circuit Breaking
- Implementing Circuit Breakers with Airflow



Prateek Chawla, Founding Engineer



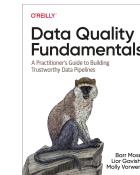
MONTE CARLO



\$236m in funding by backers of the world's best enterprise companies,
including Accel, GGV, Redpoint, ICONIQ, and IVP



Creator of the Data Observability category



Adopted

by 100s of customers across
all industries and sizes

affirm

CNN

FOX

GoodRx

HubSpot

jetBlue

PagerDuty

asics

PEPSICO

auth0

The **impact** of data downtime

~70

high severity events each year
per every 1k tables¹

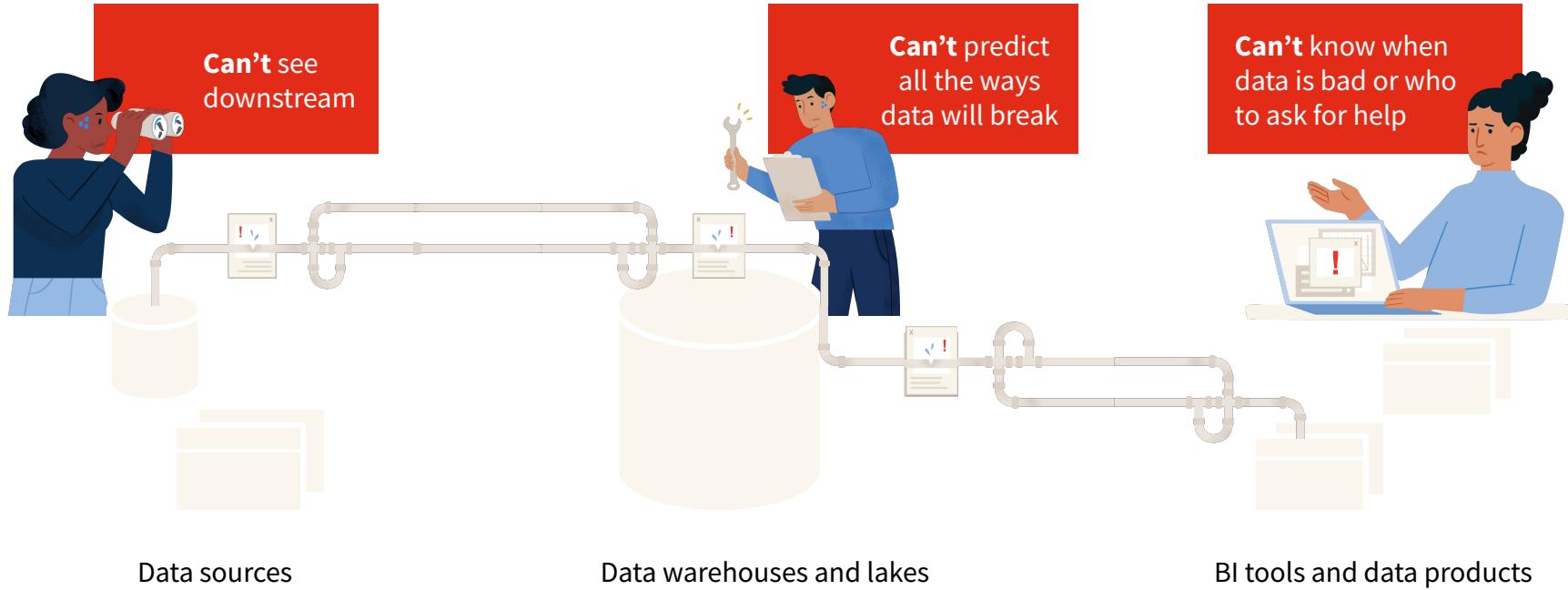
30-50%

data engineering time
spent on fire drills²

\$15M

avg. annual financial cost for
companies resulting from
poor data quality³

Detection, resolution, and schema management workflows are the issue



Good news: data downtime **looks similar** at all companies

- Is this data up-to-date?
- Why does this data size look off?
- Isn't this value suspiciously high?
- Why are there so many nulls?
- Why do we have duplicate IDs?
- What reports will I break with this schema update?
- Why are there 0s on tiles that usually show 100s?
- ...

Three Pillars of Observability Engineering



Cindy Sridharan
@copyconstruct

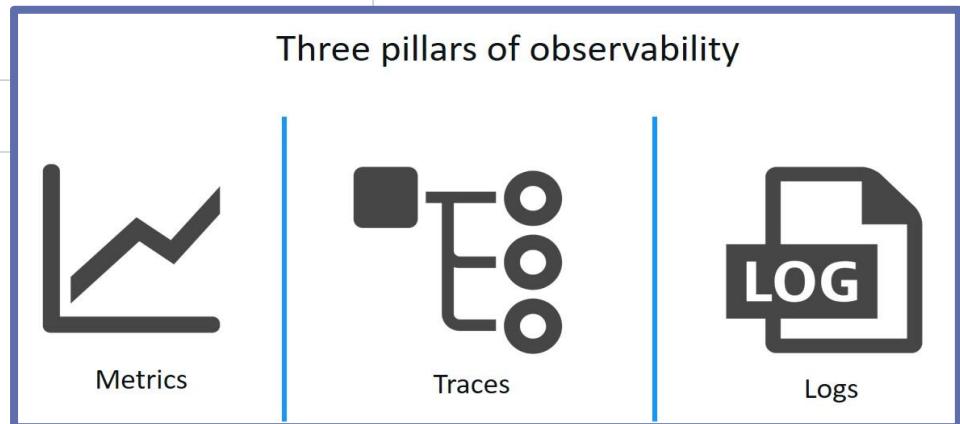


OH - "Observability - because devs don't like to do
"monitoring" we need to package it in new nomenclature to
make it palatable and trendy."

3:41 PM · Jul 28, 2017 from San Francisco, CA



54 23 people are Tweeting about this



The Solution? Detect. Resolve. Prevent.



Detect

- ML-powered anomaly detection
- Rule-based detection
- Targeted alerts to impacted owners & downstream users

Resolve

- Automated field-level lineage
- Impact radius assessment
- Code, data, and operational diagnostics

Prevent

- Auto-generated and on-demand insights
- Schema change notifications
- **Automated circuit breakers**



DATA OBSERVABILITY PILLARS

Freshness | Distribution | Volume | Schema | Lineage

The 5 Pillars of Data Observability

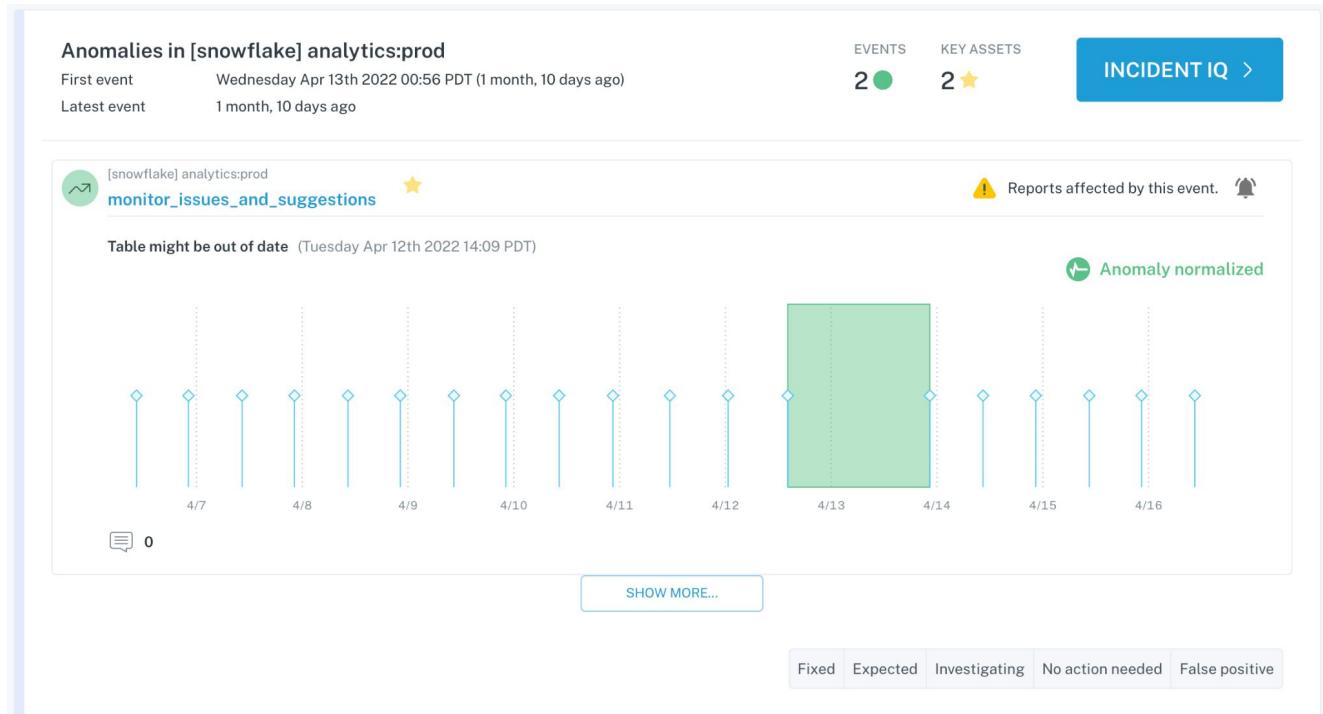
Freshness

Distribution

Volume

Schema

Lineage



The 5 Pillars of Data Observability

DATA OBSERVABILITY PILLARS

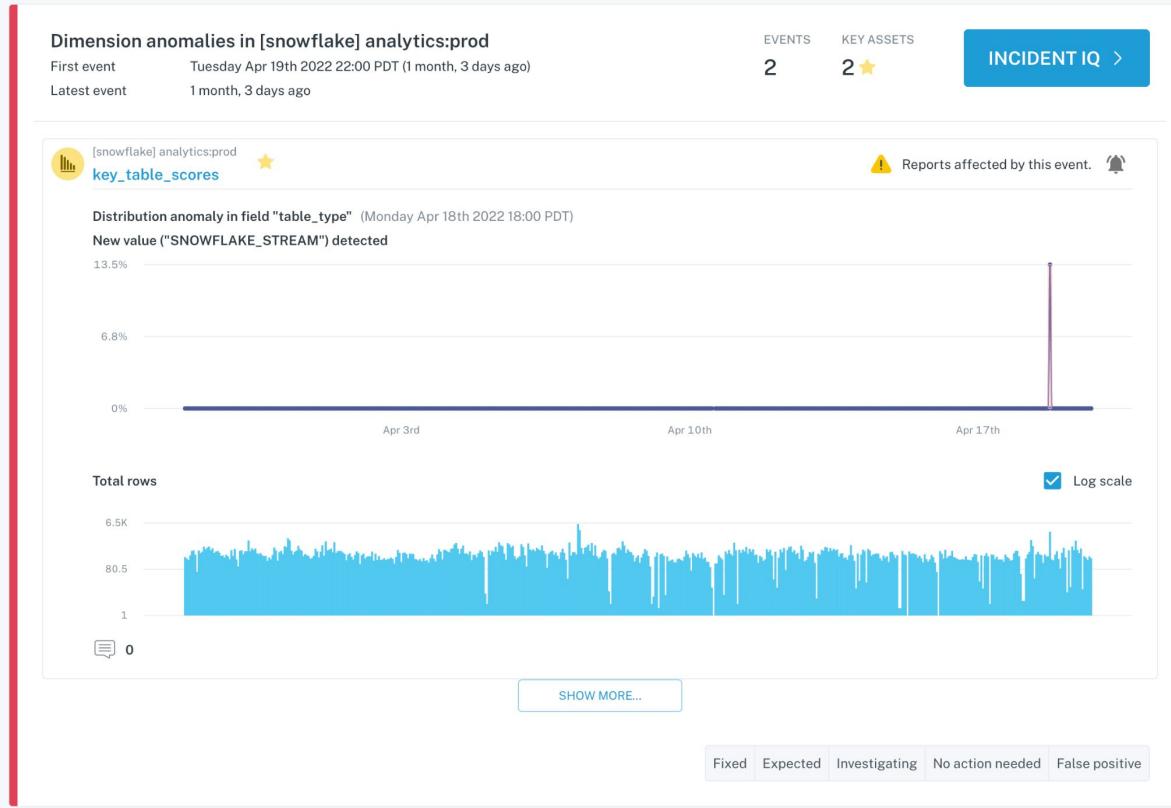
Freshness

Distribution

Volum

Schema

Lineage



The 5 Pillars of Data Observability

DATA OBSERVABILITY PILLARS

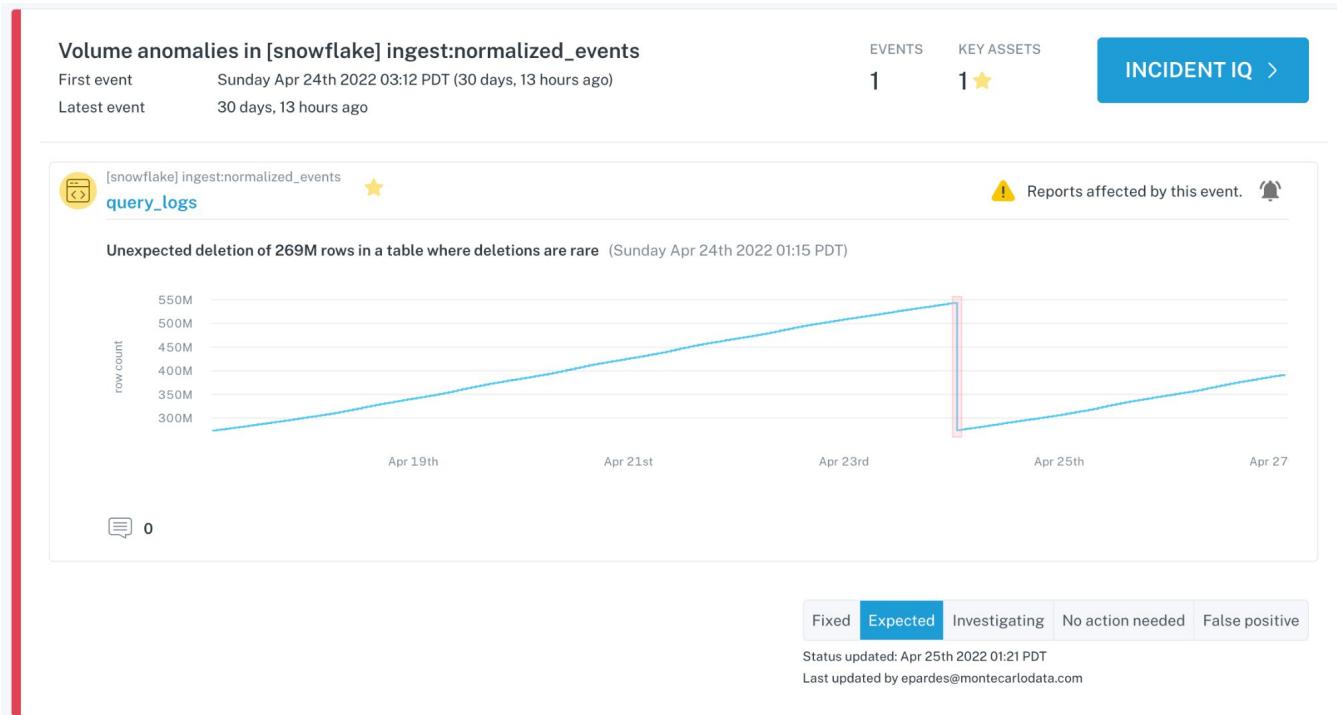
Freshness

Distribution

Volume

Schema

Lineage



The 5 Pillars of Data Observability

Freshness

Distribution

Volume

Schema

Lineage

Schema changes in [snowflake] analytics:prod_detectors

First event Sunday Apr 24th 2022 01:15 PDT (3 days, 4 hours ago)

Latest event 3 days, 2 hours ago

EVENTS

4

KEY ASSETS

3 ★

INCIDENT IQ >



[snowflake] analytics:prod_detectors
freshness_training ★

⚠ Reports affected by this event. 🔔

Schema changes detected (Sunday Apr 24th 2022 01:15 PDT)

new_n days_existing update_pattern were added

reason changed from "varchar(99)" to "varchar(126)"

max_delay_sec changed from "number(18,0)" to "number(24,0)"

0

SHOW MORE...

Fixed

Expected

Investigating

No action needed False positive

Status updated: Apr 25th 2022 01:21 PDT

Last updated by epardes@montecarloldata.com

The 5 Pillars of Data Observability

DATA OBSERVABILITY PILLARS

Freshness

Distribution

Volume

Schema

Lineage

Field Lineage

Column-level lineage across data assets

ANALYTICS:PROD active_monitors	▼
ANALYTICS:PROD_DETECTORS metric_status_indicator	▼
ANALYTICS:PROD_INSIGHTS insight_custom_rules	▼
ANALYTICS:PROD_STAGE custom_rules	▼
ANALYTICS:PROD_DETECTORS numeric_status_indicator	▼
ANALYTICS:PROD table_metadata	▼
ANALYTICS:PROD recent_metrics	▼

ANALYTICS:PROD monitor_issues_and Suggestions	SQL
monitor_status	VARCHAR(30)

ANALYTICS:DBT_APETLINOVICH monitor_issues_and_suggesti...	SQL
ANALYTICS:PROD_SNAPSHOTS monitor_issues_and_suggesti...	SQL
ANALYTICS:PROD_INSIGHTS to_monolith_monitor_issues...	SQL
ANALYTICS:PROD_INSIGHTS insight_monitor_issues_and_s...	SQL
SELECT clause lineage	1
monitor_status	VARCHAR(30)
Non-SELECT lineage	16
account_id	VARCHAR(16777216)
tables	ARRAY
full_table_id	VARCHAR(16777216)
project_name	VARCHAR(16777216)
dataset_name	VARCHAR(16777216)
created_on	TIMESTAMP_NTZ(9)
table_name	VARCHAR(16777216)

Field

agg_type

Type

VARCHAR(16777216)

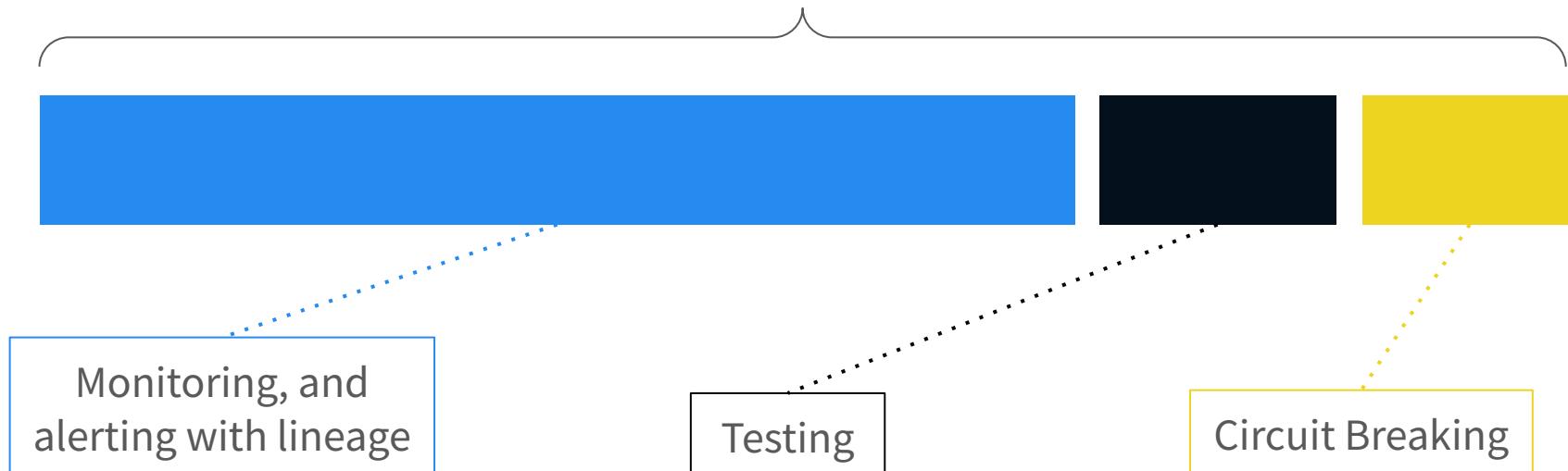
Description

Tags



Common use cases & technologies

Data downtime incidents



Aside: The primary challenge with Airflow monitoring

- Airflow has plenty of interfaces and mechanisms to help with monitoring and logging, but it is not data-aware out of the box.



So what is a **circuit breaker**, really?

- Interrupts current flow to protect equipment if an issue occurs.
- No electricity (an outage) is better than fire.
- Unlike a fuse, can be reset to resume normal operation.

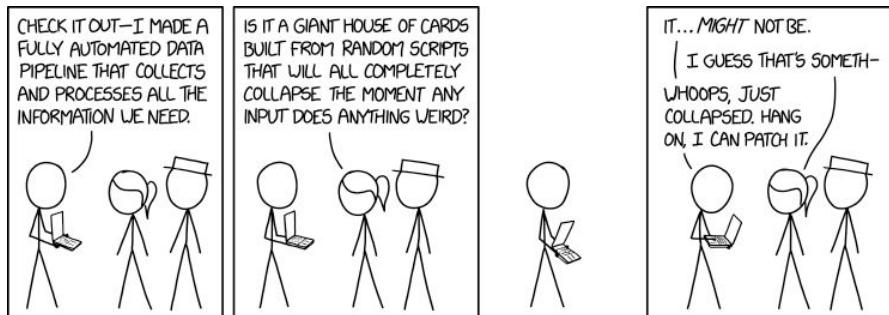


Circuit breakers **vs** testing

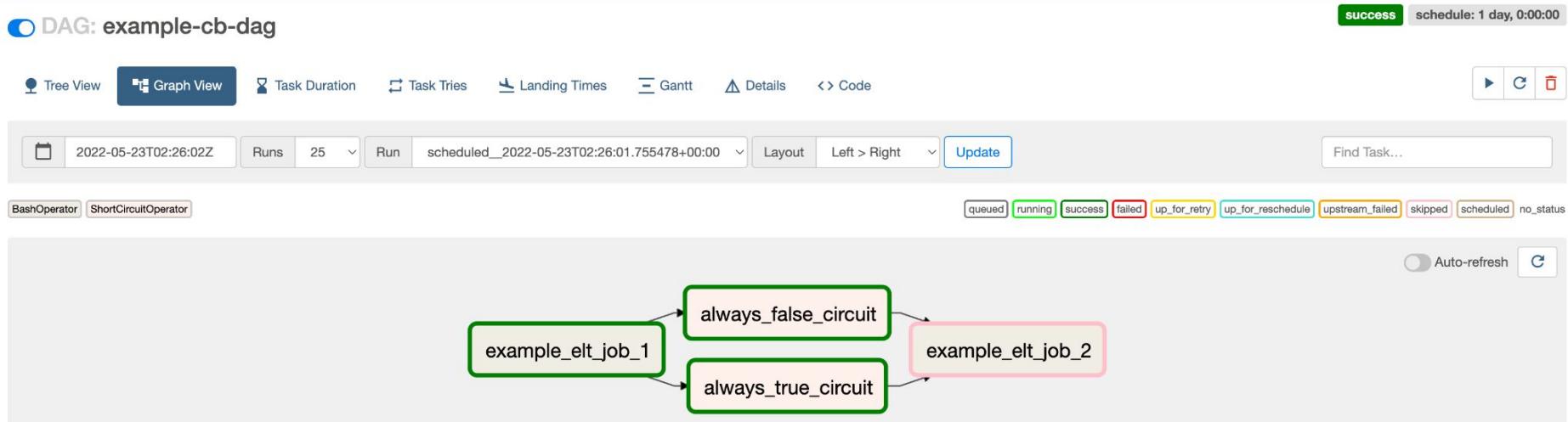
- Circuit breakers should be the **most important test cases**.
- Circuit breakers create a **guarantee** that any insights produced from this dataset meet your expectations and acceptable thresholds.
- Circuit breakers **should be used sparingly**. Testing should be prolific.
- Circuit breakers are **always automated**. No manual verification of results.
- If a circuit is open, teams should be **immediately paged** to resolve.
- Circuit breakers are **always run in production**. Can (and should) also be used in dev and stage.

Challenges with circuit breakers

- Circuit breakers provide an excellent way to **prevent data issues** from occurring in the first place.
- However, circuit breakers can also **wreak havoc on pipelines** when delayed jobs set off a chain reaction of **data failures** downstream.
- Keep the following in mind when considering how you implement circuit breakers:
 - Select only rules where you have a good understanding of how often incidents are triggered.
 - Start small, by implementing circuit breakers in a handful of Airflow DAGs, with a lot of visibility to catch issues.
 - Select rules where the underlying SQL executes quickly to avoid long query timeouts.
 - Understand where and when to fail open.



Now, let's build a circuit breaker!



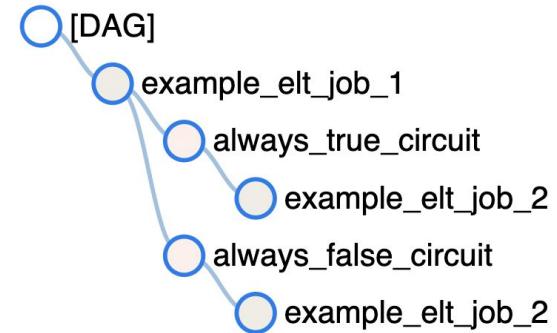
Using the ShortCircuitOperator

Airflow DAG Security Browse Admin Docs

DAG: example-cb-dag

Tree View Graph View Task Duration Task Tries Landing Times Gantt Details < > Code

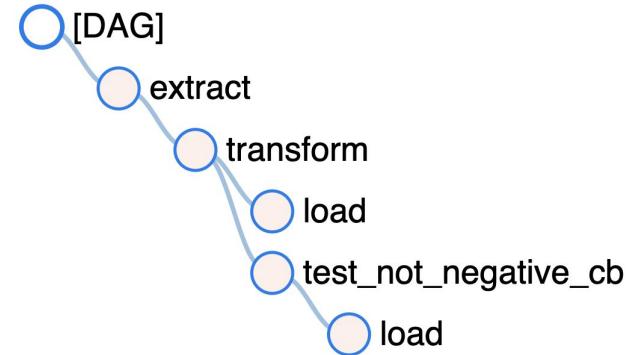
```
1 from datetime import datetime, timedelta
2
3 from airflow import DAG
4 from airflow.operators.bash import BashOperator
5 from airflow.operators.python import ShortCircuitOperator
6
7 with DAG('example-cb-dag', start_date=datetime(2022, 5, 22), catchup=False, schedule_interval=timedelta(1)) as dag:
8     example_elt_job_1 = BashOperator(
9         task_id='example_elt_job_1',
10        bash_command='echo I am transforming a very important table!',
11    )
12
13     always_true_circuit = ShortCircuitOperator(
14         task_id='always_true_circuit',
15         python_callable=lambda: True # Replace with business logic
16     )
17
18     always_false_circuit = ShortCircuitOperator(
19         task_id='always_false_circuit',
20         python_callable=lambda: False # Replace with business logic
21     )
22
23     example_elt_job_2 = BashOperator(
24         task_id='example_elt_job_2',
25         bash_command='echo I am building a very important dashboard from the table created in example_elt_job_1!',
26         trigger_rule='none_failed'
27     )
28
29     example_elt_job_1 >> [always_true_circuit, always_false_circuit] >> example_elt_job_2
```



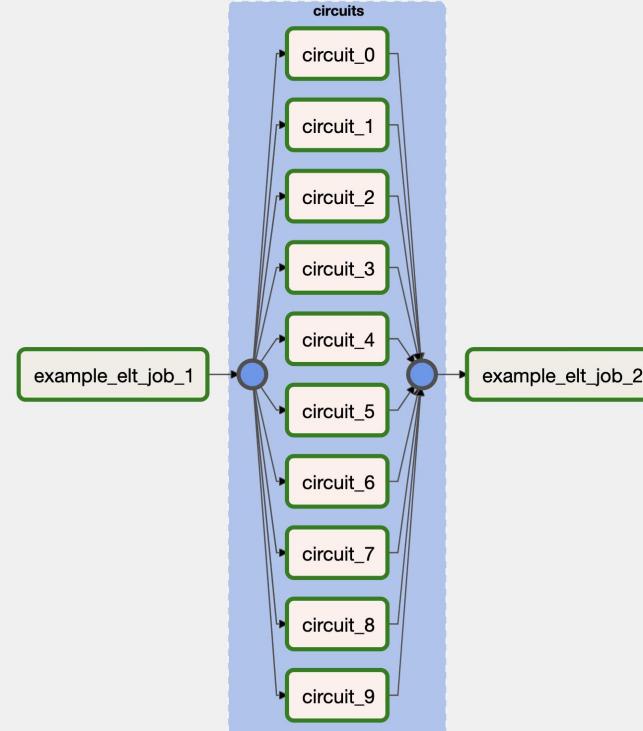
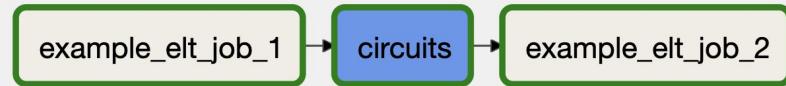
In Airflow 2.X - Using the TaskFlow API paradigm

Airflow DAGs Security Browse Admin Docs

```
1 import json
2 from typing import Dict
3
4 import pendulum
5 from airflow.decorators import dag, task
6 from airflow.operators.python import ShortCircuitOperator
7
8
9 def test_not_negative(**context) -> bool:
10     items_in_usd = context['ti'].xcom_pull(task_ids='transform')
11     for item in items_in_usd.values():
12         if item < 0:
13             return False
14     return True
15
16
17 @dag(schedule_interval=None, start_date=pendulum.datetime(2022, 5, 22, tz='UTC'), catchup=False)
18 def example_cb_etl():
19     @task()
20     def extract() -> Dict:
21         return json.loads('{"item_1": 42, "item_2": 24}') # items in dollars
22
23     @task()
24     def transform(items_in_usd: Dict) -> Dict:
25         usd_to_euro_conversion_rate = -0.937478 # bad conversion rate.
26
27         return {name: round(cost_in_usd * usd_to_euro_conversion_rate, 2) for name, cost_in_usd in items_in_usd.items()}
28
29     test_not_negative_cb = ShortCircuitOperator(
30         task_id='test_not_negative_cb',
31         provide_context=True,
32         python_callable=test_not_negative
33     )
34
35     @task()
36     def load(items_in_euros: Dict) -> None:
37         for name, cost_in_euros in items_in_euros.items():
38             print(f"Item {name} is {cost_in_euros} euros.")
39
40     transformed_items_in_euros = transform(extract())
41     transformed_items_in_euros >> test_not_negative_cb >> load(transformed_items_in_euros)
42
43
44 example_cb_etl_dag = example_cb_etl()
```



In Airflow 2.X - Using Task Groups to improve manageability



Suggestions when adding circuit breakers to your pipeline

- **Do not** limit yourself to one type of operator.
 - You can use any operator, not just the **ShortCircuitOperator** to create a custom circuit breaker.
 - You can leverage many tools like DBT, Great Expectations, or Monte Carlo 😊 as a circuit breaker.
- **Try** raising an **AirflowSkipException** instead of an **AirflowException** when a circuit is closed.
 - This increases visibility and prevents automatic retries.
- **Do not** merge multiple circuits into one operator.
 - This makes it harder to trace exactly what issue or threshold tripped the circuit.
- **Do not** merge a circuit breaker into the operator performing the task to evaluate.
 - This makes it more difficult to determine if the job failed or if the circuit was tripped.
- **Include** a mechanism to bypass or skip a circuit breaker.
- **Generalize** repeated patterns as plugins or custom operators.

Practical thresholds: Using **freshness** in circuit breakers

- A common pattern to circuit break on is to check if your table was actually updated.
- Many warehouse and lakes expose this information via a metastore or log.
 - In Snowflake:
`SELECT CONVERT_TIMEZONE('UTC', last_altered) last_altered FROM information_schema.tables ...`
 - With Delta Lake:
`DESCRIBE DETAIL ...` or via the `deltaLog.snapshot`
- It is worth combining this metric with volume and tracking changes and patterns over time.

Practical thresholds: Using *custom logic* in circuit breakers

- Breach when row count is greater than 0.

Row count match

Two tables have the same number of rows.

```
SQL
WITH t1
AS (SELECT Count(*)
     FROM {{table_1}}),
t2
AS (SELECT Count(*)
     FROM {{table_2}}),
t3
AS (SELECT (SELECT *
              FROM t1) - (SELECT *
                             FROM t2) AS diff_count)
SELECT diff_count
FROM t3
WHERE diff_count != 0
```

Identical tables

Two tables are identical record for record.

```
SQL
(SELECT *
  FROM {{table_1}}
{{except_operator}}
SELECT *
  FROM {{table_2}})
UNION ALL
(SELECT *
  FROM {{table_2}}
{{except_operator}})
SELECT *
  FROM {{table_1}})
```

Referential integrity

All IDs in a child table must be present in a parent table.

```
SQL
WITH child
AS (SELECT {{column}} AS id
     FROM {{table_1}}),
parent
AS (SELECT {{column}} AS id
     FROM {{table_2}})
SELECT *
FROM child
LEFT JOIN parent using (id)
WHERE child.id IS NOT NULL
      AND parent.id IS NULL
```

What happens **AFTER** a circuit breaker is triggered?

- Each data incident is different, but the following **7** steps are good starting point:
 1. **Review** table and field lineage for all impacted downstream assets and **notify** stakeholders.
 2. **Review** logs and any documentation for the circuit that tripped.
 3. **Review** logs and any documentation for the tasks(s) being evaluated.
 4. **Review** field lineage for all upstream transformations for table(s) that could be the root cause.
 5. **Check** for recent changes to the pipeline, like new data sources, streams, or modified queries.
 6. **Sample** table for anomalies and use any tools to review historical versions of the model.
 7. **Review** for any historical incidents or schema changes on the affected table(s).

Check out our [Airflow Provider](#) and [SDK](#)

airflow-mcd 0.0.5

`pip install airflow-mcd` 

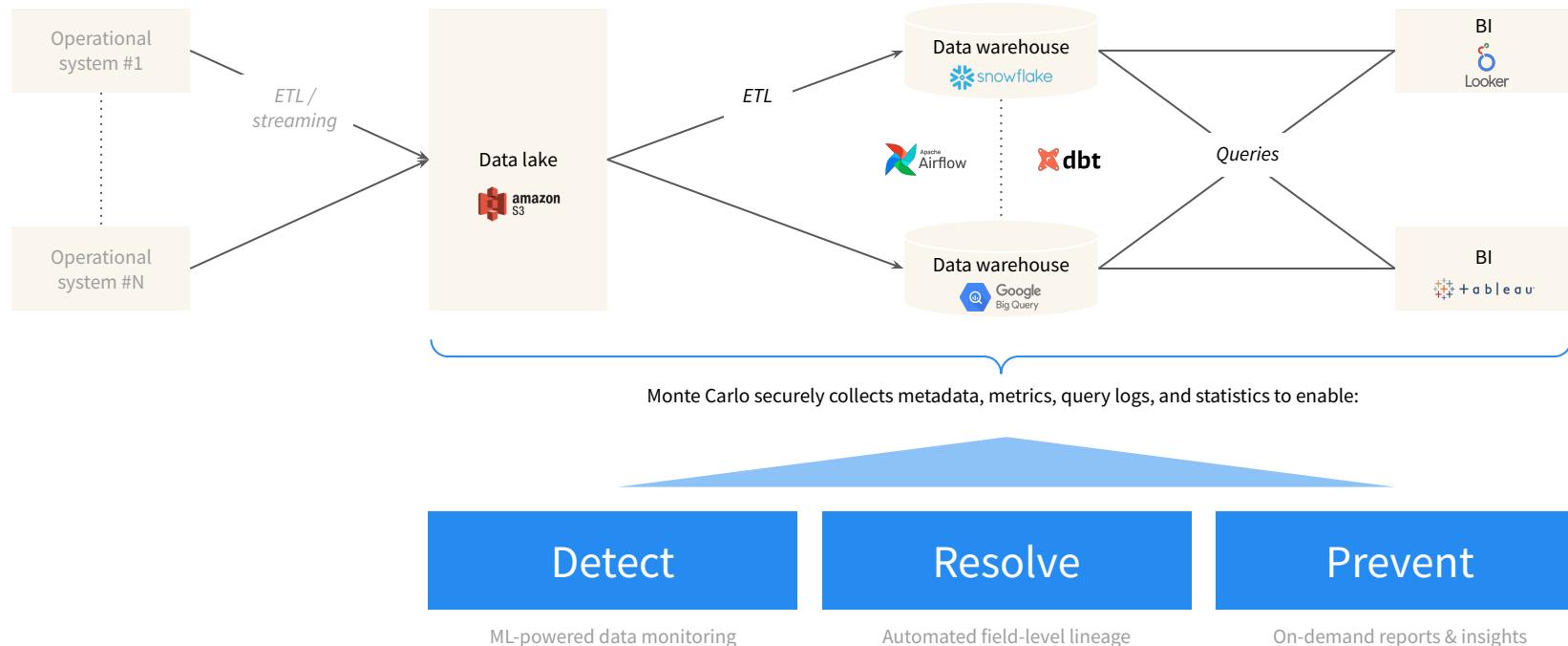
<https://pypi.org/project/airflow-mcd/>

pycarlo 0.1.2

`pip install pycarlo` 

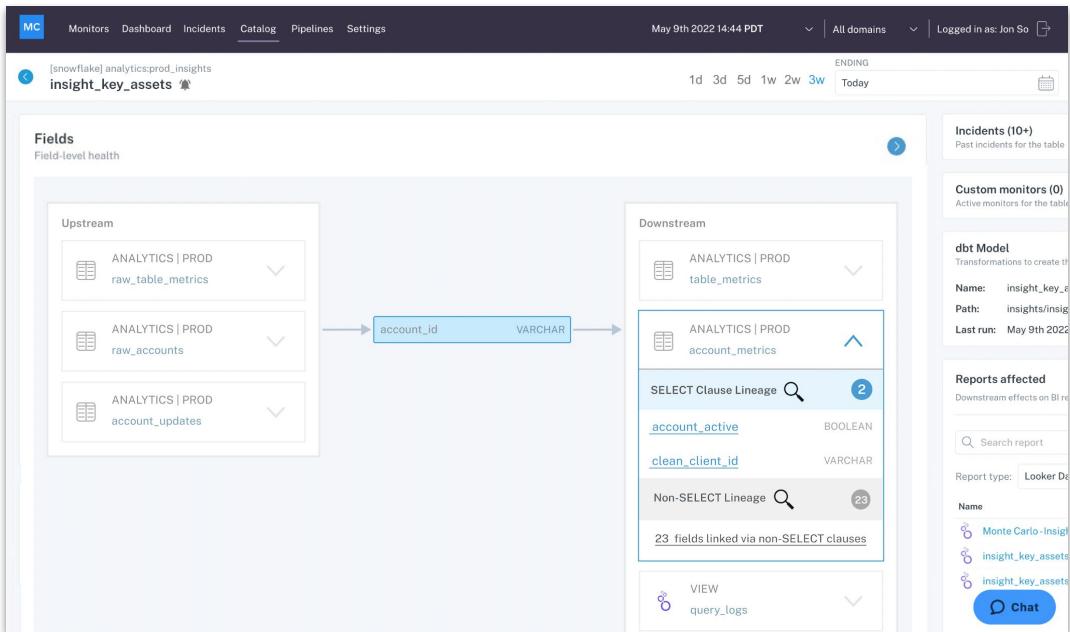
<https://pypi.org/project/pycarlo/>

A complete data observability platform enables your team to detect, resolve, and prevent bad data altogether



Triage, find root cause, and manage incident workflows all in one platform

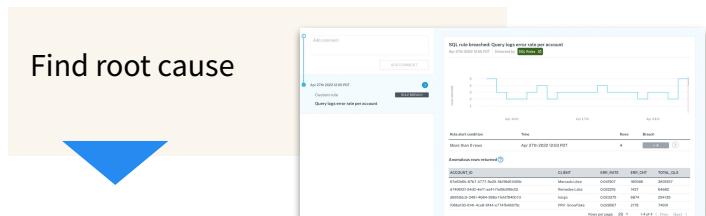
Automated field-level lineage



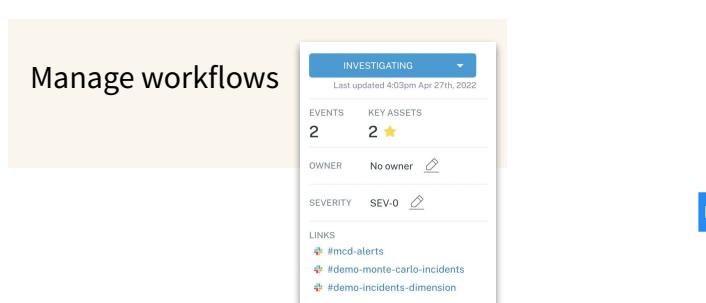
Triage incidents



Find root cause



Manage workflows



Key takeaways

- Data downtime has a measurable **impact**.
- The 5 pillars of data observability can give you a **holistic view of data health**.
- Circuit breaking is **powerful** and can be **easy to implement**, but it is only one tool and must be combined with others like monitoring, alerting, lineage, and testing for maximum effectiveness.
- Circuit breakers are like testing on steroids, so you have to **be careful** to prevent unintended data downtime.
- You have a **lot of options** for how to instrument circuit breakers and what **quality or integrity thresholds** you can use. Experiment and have fun!

Thank you!

www.montecarlodata.com/blog

Prateek Chawla: pchawla@montecarlodata.com
