



`New` Workflow Orchestrator in town:

"Apache Airflow 2.x"



Jarek Potiuk

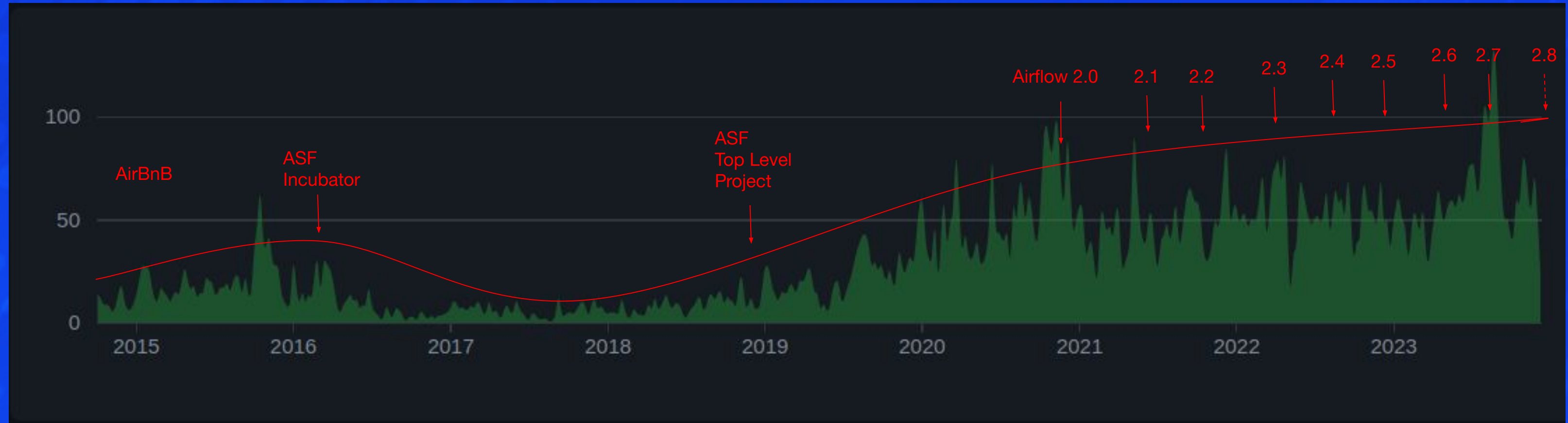
Apache Airflow Committer
<https://github.com/potiuk>

Airflow as ..

- Modern UI
- With great DAG Authoring capabilities
- Being extensible platform
- Being True Open Source with as strong community you can get
- Solid infrastructure
- Shortly - modern orchestrator of your choice :)



Airflow 2 timeline



10th Anniversary next year:
Unofficial
Airflow Summit 2024, September, Bay Area, 1000+ attendees



Modern UI

Grid View

The screenshot displays a task monitoring interface. At the top, there is a header with a timestamp '06/08/2023, 09:10:13 PM', a page number '25', and filter options for 'All Run Types' and 'All Run Sta'. A 'Clear Filters' button and an 'Auto-refresh' toggle are also present.

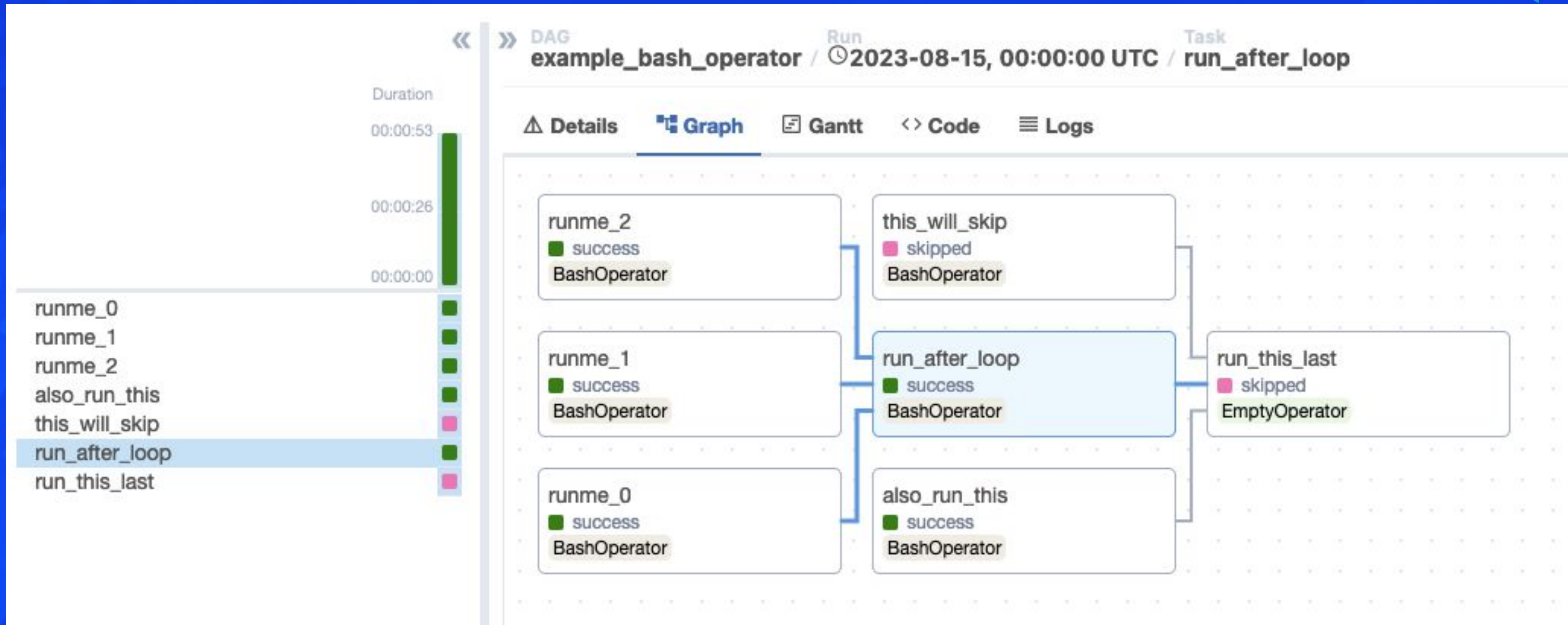
Below the header, a row of status buttons is shown: 'deferred', 'failed', 'queued', 'removed', 'restarting', 'running', 'scheduled', 'shutdown', 'skipped', 'success', 'up_for_reschedule', and 'up_for_retry'. A note indicates 'Press shift + / for Shortcuts'.

The main content area is split into two panels. The left panel shows a 'Grid View' with a bar chart of task durations. The y-axis is labeled 'Duration' and ranges from '00:00:00' to '00:01:35'. The x-axis shows dates from 'May 24, 15:52' to 'Jun 04, 21:42'. A table below the chart lists task sections: 'start', 'section_1', 'section_2', and 'end', each with a corresponding row of green bars representing task durations.

The right panel shows a 'Gantt' view for a task group 'example_task_group' with a run time of '2023-06-07, 21:42:49 UTC'. It includes buttons for 'Clear task', 'Mark state as...', and 'Filter Tasks'. Below these are navigation tabs for 'Details', 'Graph', 'Gantt', and 'Logs'. The Gantt chart shows a timeline from '21:42:50 UTC' to '21:42:55 UTC' with green bars indicating task execution.



Graph View



Log view



09/19/2023, 04:08:53 PM 25 All Run Types All Run States Clear Filters Auto-refresh

Press shift + / for Shortcuts

deferred failed queued removed restarting running scheduled shutdown skipped success up_for_reschedule up_for_retry upstream_failed no_status

example_task_group / 2023-09-19, 14:56:54 UTC / task_2

Clear task Mark state as... Filter Tasks

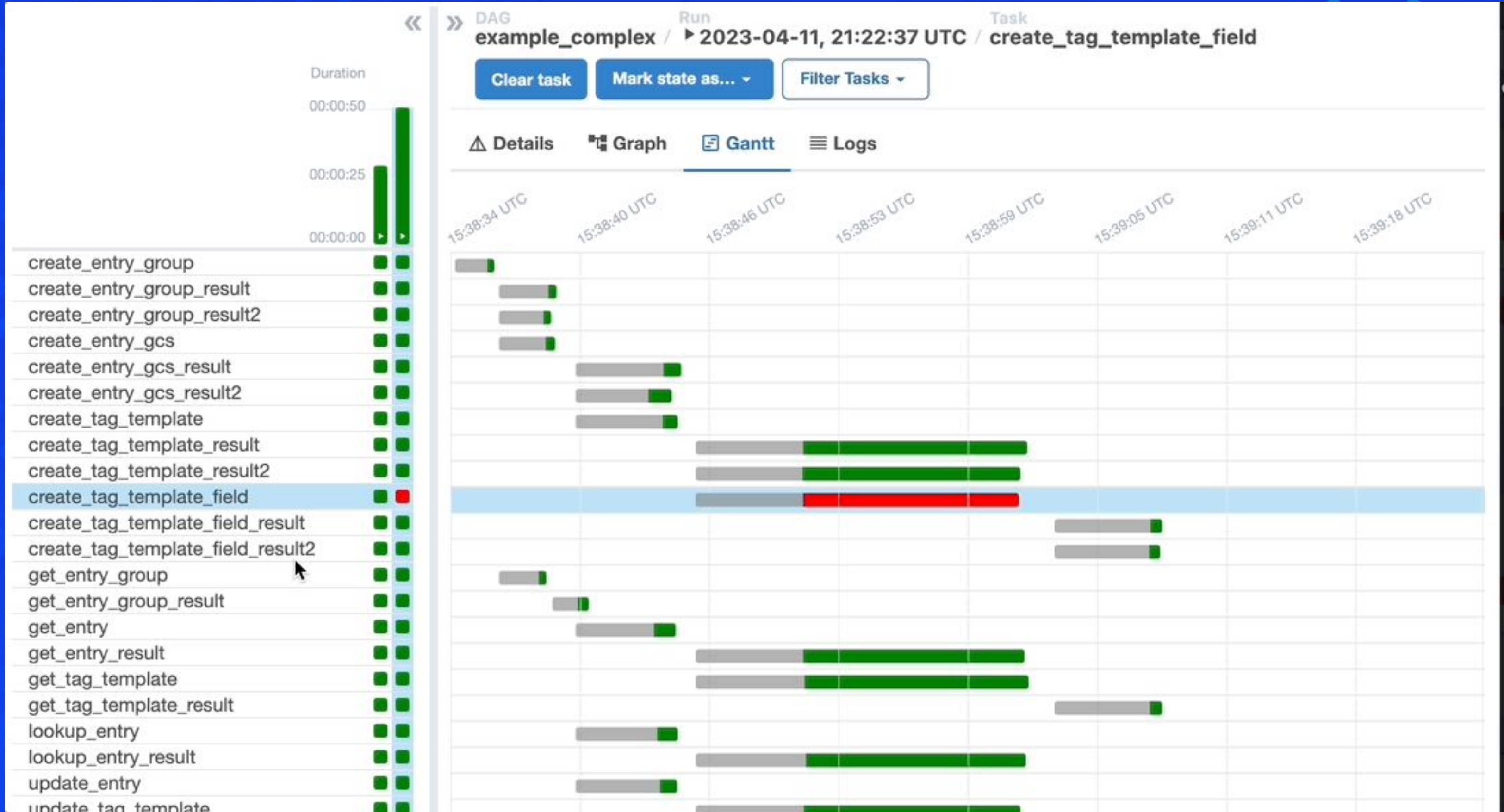
Details Graph Gantt Code Logs

(by attempts) 1

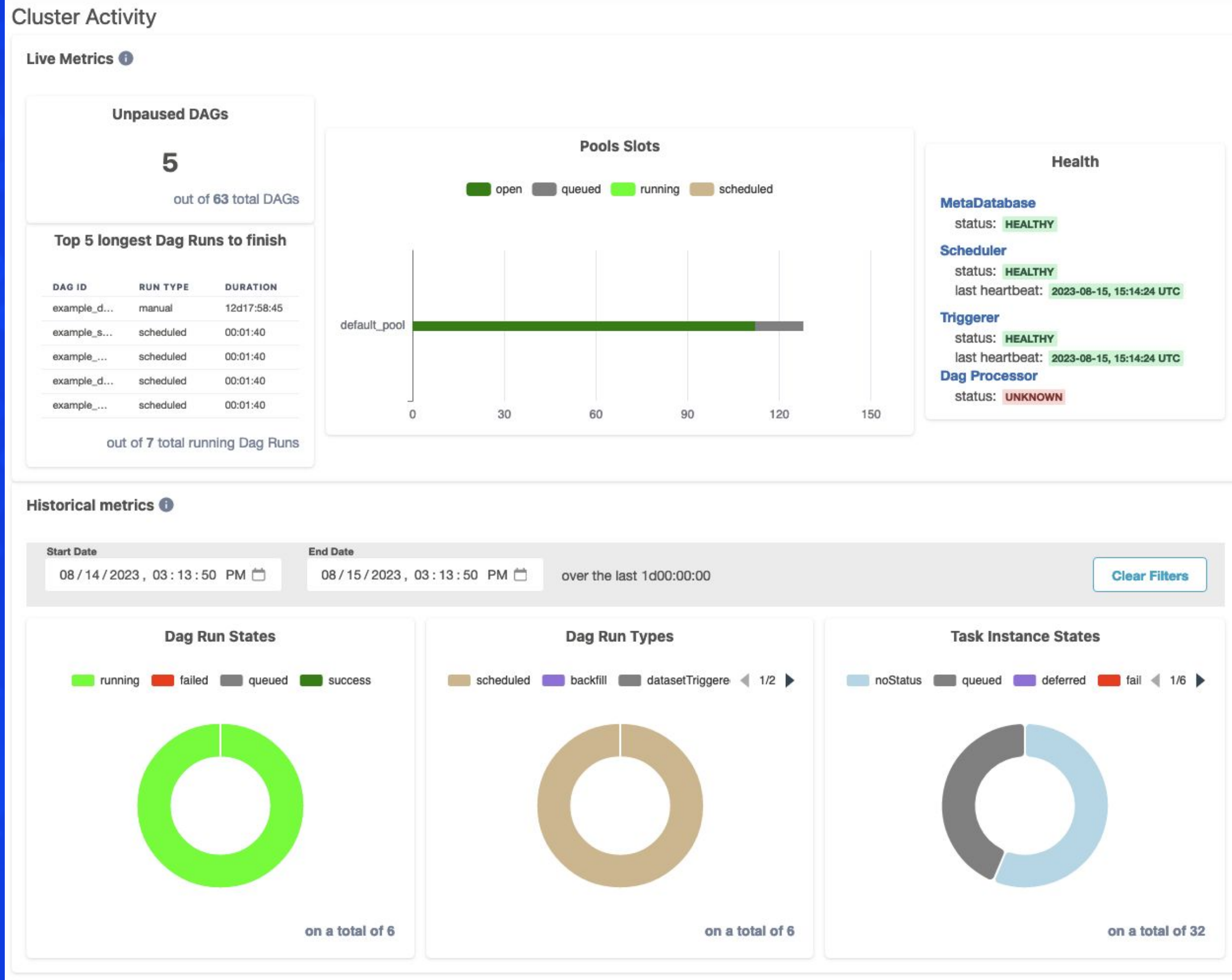
All Levels All File Sources Wrap Download See More

```
172.20.8.62
*** Found logs in s3:
*** * s3://airflow-logs-clmp5782l000501qdew3h97wb/clmp9s2on1497071uyvq4eu8g0p/dag_id=example_task_group/run_id>manual__2023-09-19T14:56:54.216546+00:00/task_id=section_1.task_2/attempt=1.log
[2023-09-19, 14:58:47 UTC] {taskinstance.py:1157} INFO - Dependencies all met for dep_context=non-requeueable deps ti=<TaskInstance: example_task_group.section_1.task_2 manual__2023-09-19T14:56:54.216546+00:00 [queued]>
[2023-09-19, 14:58:47 UTC] {taskinstance.py:1157} INFO - Dependencies all met for dep_context=requeueable deps ti=<TaskInstance: example_task_group.section_1.task_2 manual__2023-09-19T14:56:54.216546+00:00 [queued]>
[2023-09-19, 14:58:47 UTC] {taskinstance.py:1359} INFO - Starting attempt 1 of 1
[2023-09-19, 14:58:47 UTC] {taskinstance.py:1380} INFO - Executing <Task(BashOperator): section_1.task_2> on 2023-09-19 14:56:54.216546+00:00
[2023-09-19, 14:58:47 UTC] {standard_task_runner.py:57} INFO - Started process 35 to run task
[2023-09-19, 14:58:47 UTC] {standard_task_runner.py:84} INFO - Running: ['airflow', 'tasks', 'run', 'example_task_group', 'section_1.task_2', 'manual__2023-09-19T14:56:54.216546+00:00', '--job-id', '9', '--raw', '--subdir', '/usr/local/lib/
[2023-09-19, 14:58:47 UTC] {standard_task_runner.py:85} INFO - Job 9: Subtask section_1.task_2
[2023-09-19, 14:58:47 UTC] {task_command.py:415} INFO - Running <TaskInstance: example_task_group.section_1.task_2 manual__2023-09-19T14:56:54.216546+00:00 [running]> on host 172.20.8.62
[2023-09-19, 14:58:47 UTC] {taskinstance.py:1660} INFO - Exporting env vars: AIRFLOW_CTX_DAG_OWNER='airflow' AIRFLOW_CTX_DAG_ID='example_task_group' AIRFLOW_CTX_TASK_ID='section_1.task_2' AIRFLOW_CTX_EXECUTION_DATE='2023-09-19T14:56:54.2165
[2023-09-19, 14:58:48 UTC] {subprocess.py:63} INFO - Tmp dir root location: /tmp
[2023-09-19, 14:58:48 UTC] {subprocess.py:75} INFO - Running command: ['/bin/bash', '-c', 'echo 1']
[2023-09-19, 14:58:48 UTC] {subprocess.py:86} INFO - Output:
[2023-09-19, 14:58:48 UTC] {subprocess.py:93} INFO - 1
```

Gantt view



Cluster Activity





DAG Authoring



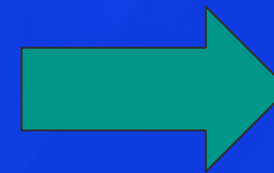
Handling dependencies

- An issue in early Airflow 2.0 days - much less nowadays
- [Multiple options](#) to handle it
 - Python Virtualenv Operator, External Python Operator, Docker Operator, Kubernetes Pod Operator, Multiple Docker Images + Celery Queues
 - Coming soon -> Multi-tenancy setup with per-tenant dependencies
- [Mastering dependencies: The Airflow Way](#) talk from Airflow Summit 2023
- Plays super-well with Task Flow

TaskFlow



```
def choose_mode():  
    accuracy = 6  
    if accuracy > 5:  
        return "accurate"  
    return 'inaccurate'  
  
choose_best_model = BranchPythonOperator(  
    task_id = 'choose_best_model',  
    python_callable = choose_best_model  
)
```



```
@task.branch  
def choose_best_model():  
    accuracy = 6  
    if accuracy > 5:  
        return 'accurate'  
    return 'inaccurate'
```



Task Flow cases

Core:

- @dag
- @task.python
- @task.virtualenv
- @task.external_python
- @task.sensor
- @task.branch
- @task.short_circuit
- @task.bash (coming)

Providers:

- @task.docker
- @task.kubernetes
- @task.sftp_sensor
- ...
- Providers can provide their own
... and
- @task_group

Task Groups

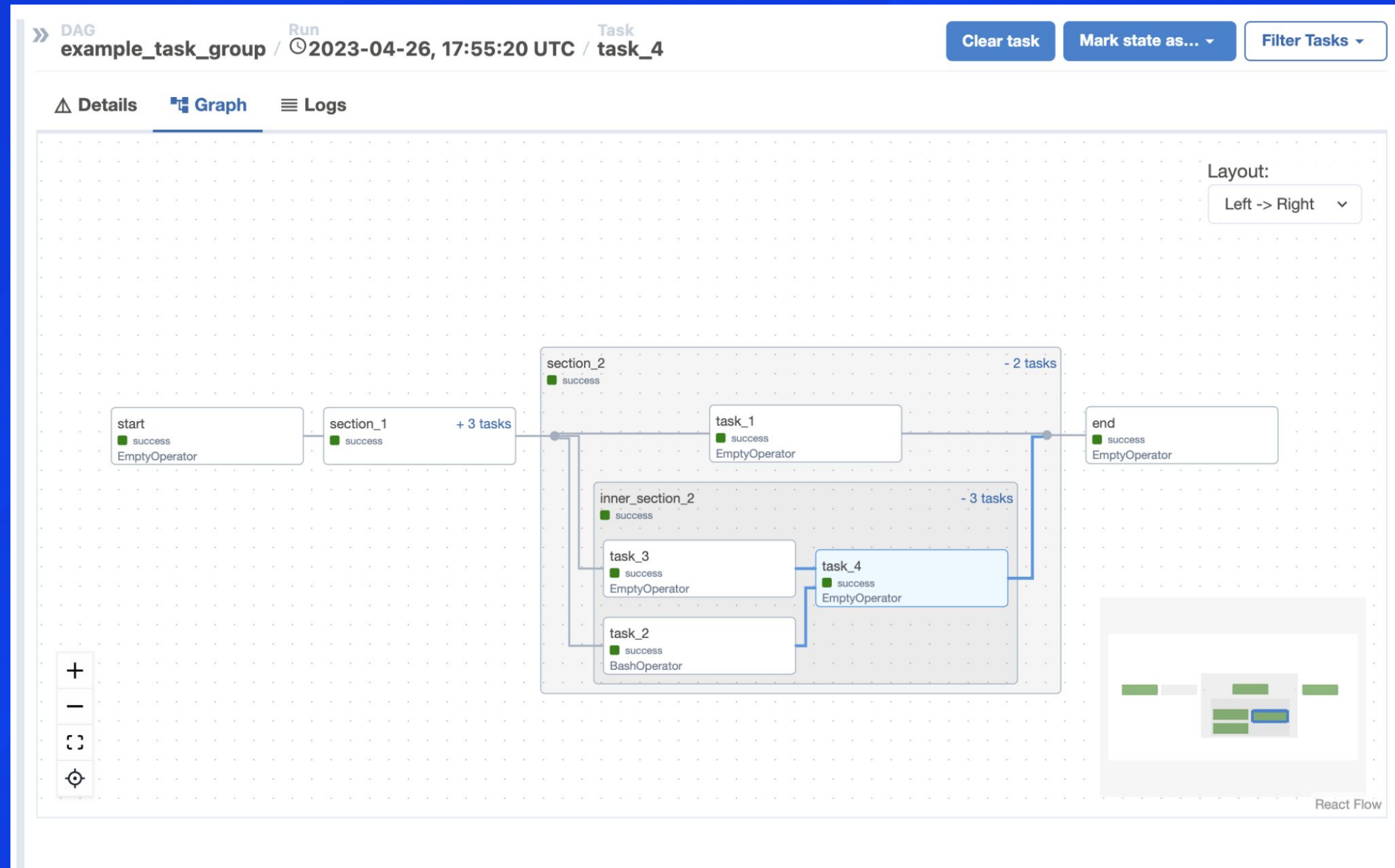
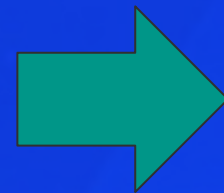


```
@task
def task_1(value: int) -> str:
    """Empty Task1"""
    return f"[ Task1 {value} ]"

@task
def task_2(value: str) -> str:
    """Empty Task2"""
    return f"[ Task2 {value} ]"

@task
def task_3(value: str) -> None:
    """Empty Task3"""
    print(f"[ Task3 {value} ]")

@task_group
def task_group_function(value: int) -> None:
    task_3(task_2(task_1(value)))
```

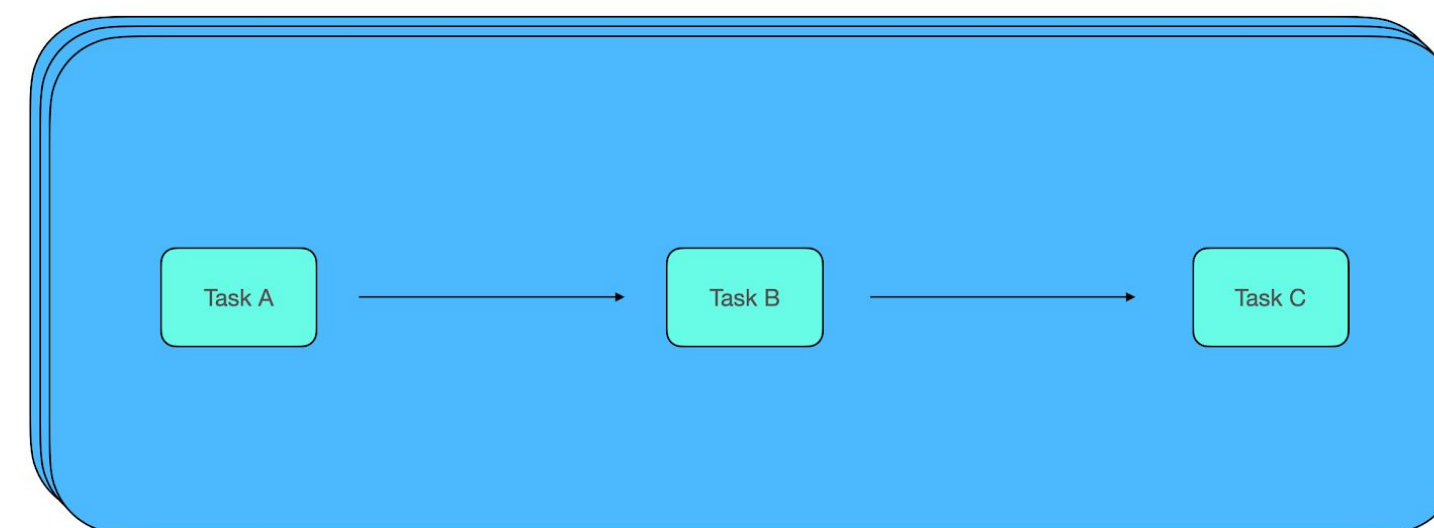


Dynamic Task and Group mapping

Breadth First Execution



Depth First Execution



- Map Reduce - kind of workflows if you want Airflow to also “do stuff”
- You can parallelise even complex workflows

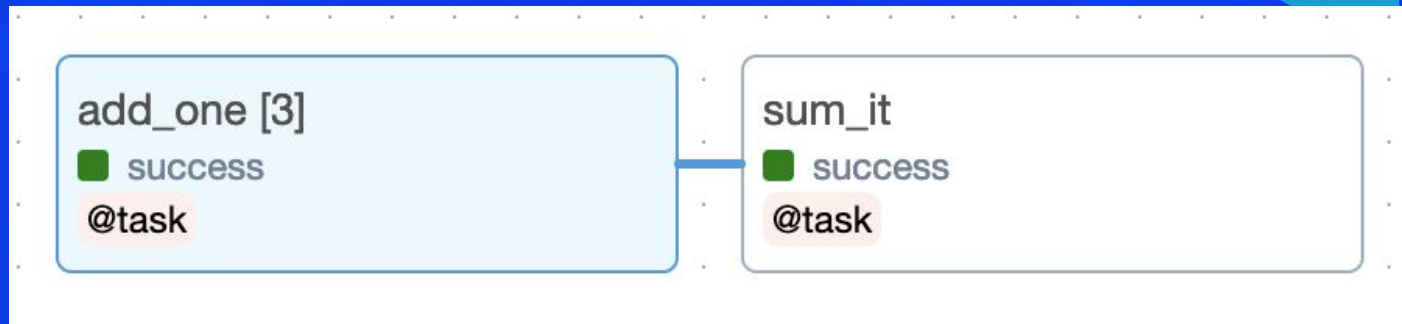
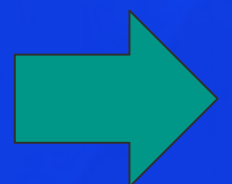
Dynamic Task mapping

```

@task
def add_one(x: int):
    return x + 1

@task
def sum_it(values):
    total = sum(values)
    print(f"Total was {total}")

added_values = add_one.expand(x=[1, 2, 3])
sum_it(added_values)
    
```



DAG: simple_mapping / Run: 2022-03-07, 17:00:00 MST / Task: add_one

Status: success
 Started: 2022-04-08, 17:29:05 MDT
 Ended: 2022-04-08, 17:29:06 MDT

3 Tasks Mapped
 success: 3

Task Id: add_one
 Run Id: scheduled__2022-03-07T00:00:00+00:00
 Operator: _PythonDecoratedOperator
 Duration: 00:00:01

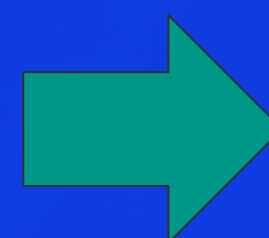
MAP INDEX	STATE	DURATION	START DATE	END DATE
0	success	00:00:00	2022-04-08, 17:29:05 MDT	2022-04-08, 17:29:06 MDT
1	success	00:00:00	2022-04-08, 17:29:05 MDT	2022-04-08, 17:29:06 MDT
2	success	00:00:00	2022-04-08, 17:29:05 MDT	2022-04-08, 17:29:06 MDT



Deferrable (AsyncIO) operators

```
class WaitOneHourSensor(BaseSensorOperator):
    def execute(self, context: Context) -> None:
        self.defer(
            trigger=TimeDeltaTrigger(timedelta(hours=1)),
            method_name="execute_complete"
        )

    def execute_complete(self, ti: TaskInstance) -> None:
        # We have no more work to do here.
        # Mark as complete.
        return
```



80% - 90%
performance
improvements

- no worker slots while waiting (other jobs can run)
- multiple 100s of Deferrable Operators out-of-the-box
- 10s of Triggers available
- you can roll your own Trigger

Setup/Teardown



```
create_cluster.as_setup() >> run_query >> delete_cluster.as_teardown()  
create_cluster >> delete_cluster
```



Notifiers

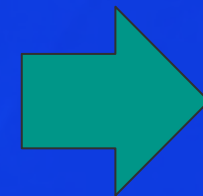


```
from airflow.notifications.basenotifier import BaseNotifier
from my_provider import send_message

class MyNotifier(BaseNotifier):
    template_fields = ("message",)

    def __init__(self, message):
        self.message = message

    def notify(self, context):
        # Send notification here, below is an example
        title = f"Task {context['task_instance'].task_id} failed"
        send_message(title, self.message)
```



```
from datetime import datetime

from airflow.models.dag import DAG
from airflow.operators.bash import BashOperator

from myprovider.notifier import MyNotifier

with DAG(
    dag_id="example_notifier",
    start_date=datetime(2022, 1, 1),
    schedule_interval=None,
    on_success_callback=MyNotifier(message="Success!"),
    on_failure_callback=MyNotifier(message="Failure!"),
):
    task = BashOperator(
        task_id="example_task",
        bash_command="exit 1",
        on_success_callback=MyNotifier(message="Task Succeeded!"),
    )
```

- easily reusable notifiers when your task fails (or not)

Object storage - FsSpec (Coming in Airflow 2.8)



- Open standard
- Integrates with all object storages
- Modern Pythonic way of interacting
 - Pathlib
- Supported by:
 - Pandas, Polars, Parquet, DuckDB, Iceberg, PyArrow
- One way to rule them all

```
fs_base = ObjectStoragePath("s3://airflow-tutorial-data/", conn_id="aws_default")
# ensure the bucket exists
fs_base.mkdir(exist_ok=True)

formatted_date = execution_date.format("YYYYMMDD")
path = fs_base / f"air_quality_{formatted_date}.parquet"

df = pd.DataFrame(response.json()).astype(foelds)
with path.open("wb") as file:
    df.to_parquet(file)

return path
```

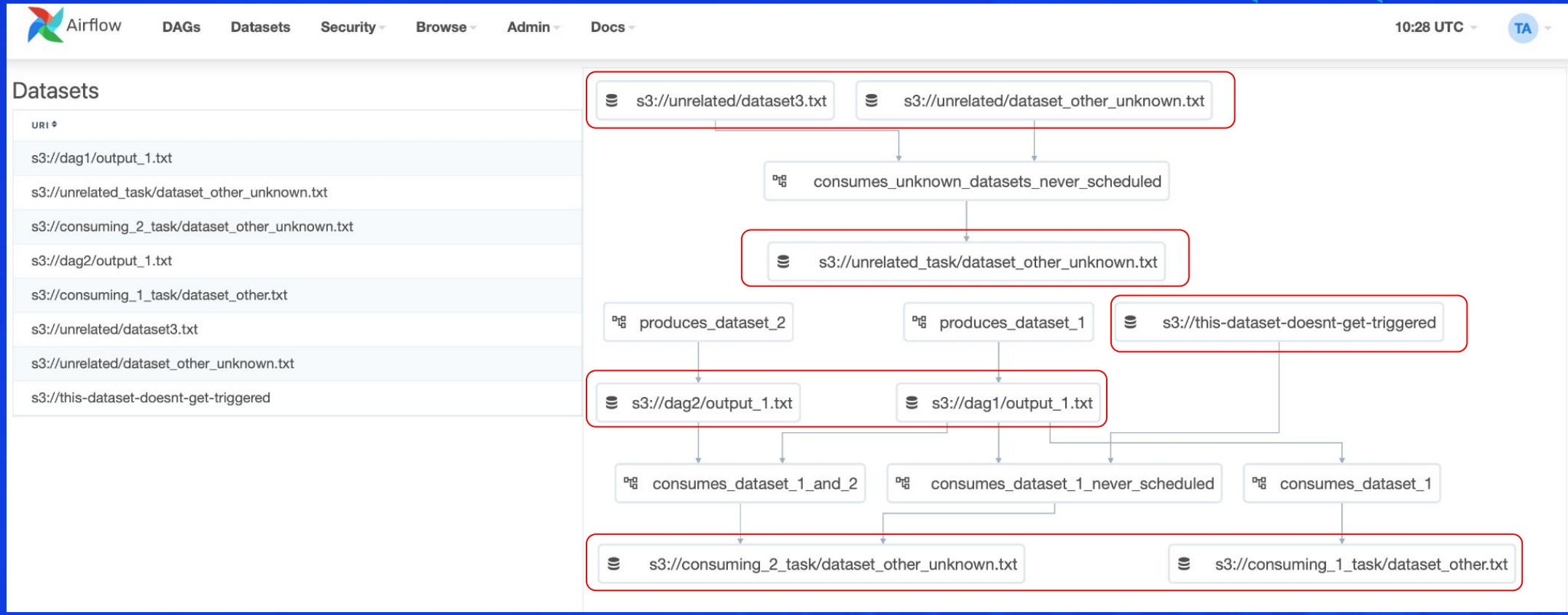
Data-aware scheduling

```

example_dataset = Dataset("s3://dataset/example.csv")

with DAG(dag_id="producer", ...):
    BashOperator(task_id="producer", outlets=[example_dataset], ...)

with DAG(dag_id="consumer", schedule=[example_dataset], ...):
    ...
    
```



- Micropipelines concept
- Still early days
- But mind-boggling things are coming (Object storage integration, Partial Datasets, Data aware triggering, Open Lineage)



LLM Operators

Donated by Astronomer (yay!)

- Open AI
- Cohere
- Weviate
- pgvector
- Pinecone
- OpenSearch

Powering @AskAstro:

<https://ask.astronomer.io/>

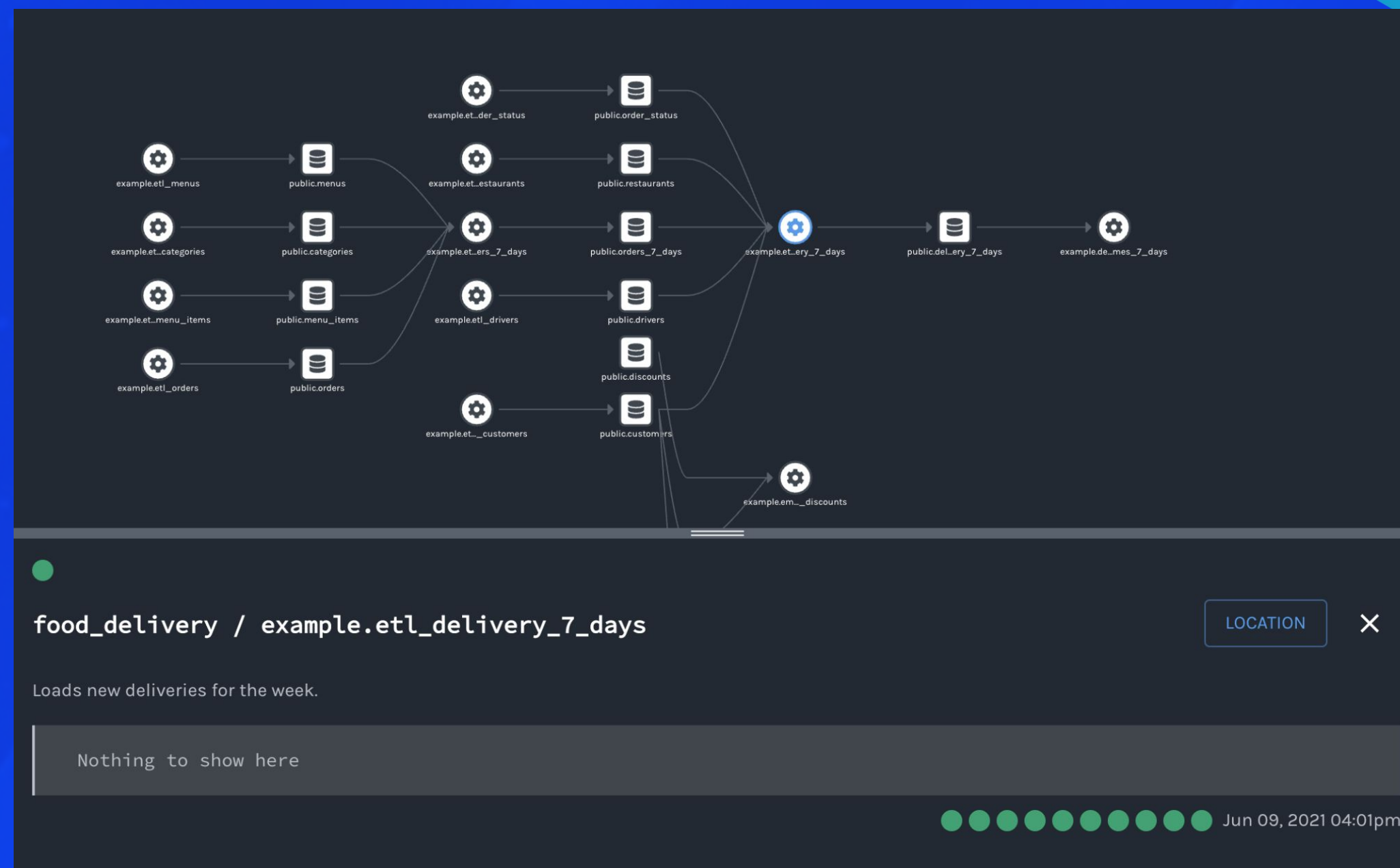




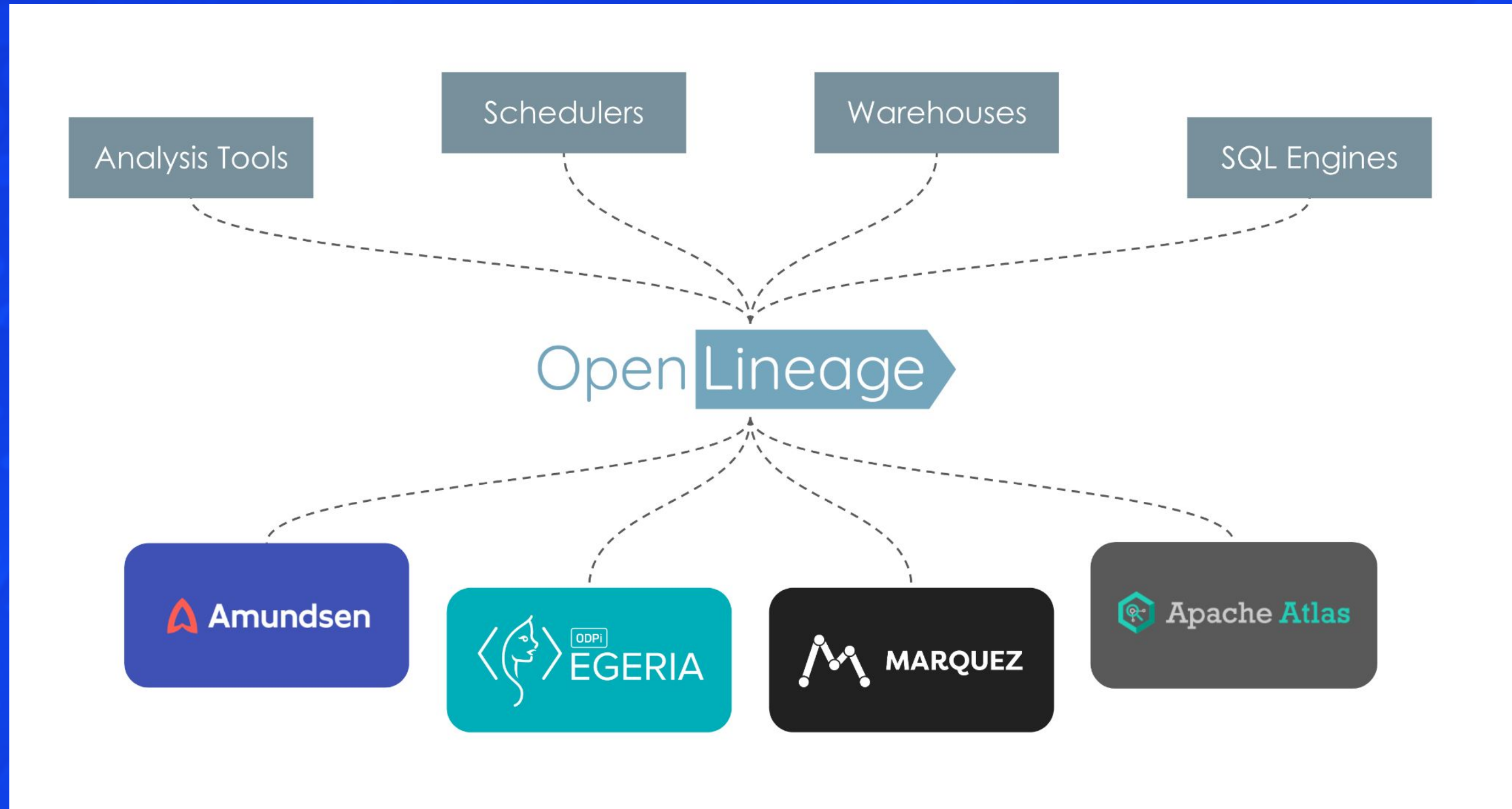
Airflow as a platform

Open Lineage

- Integrated in Airflow
- Column level lineage
- Better TaskFlow support in works
- Great adoption as open standard



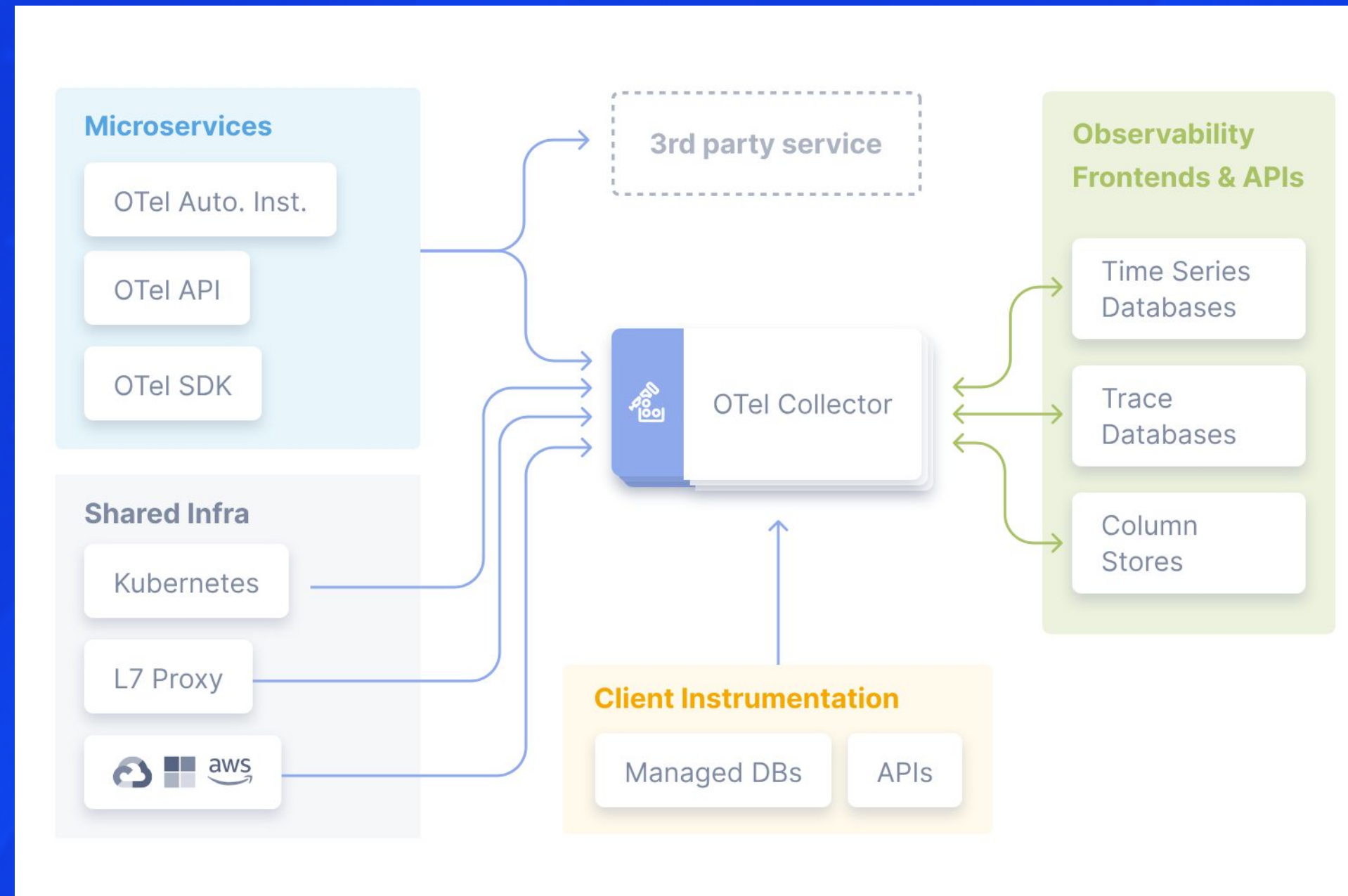
Open Lineage



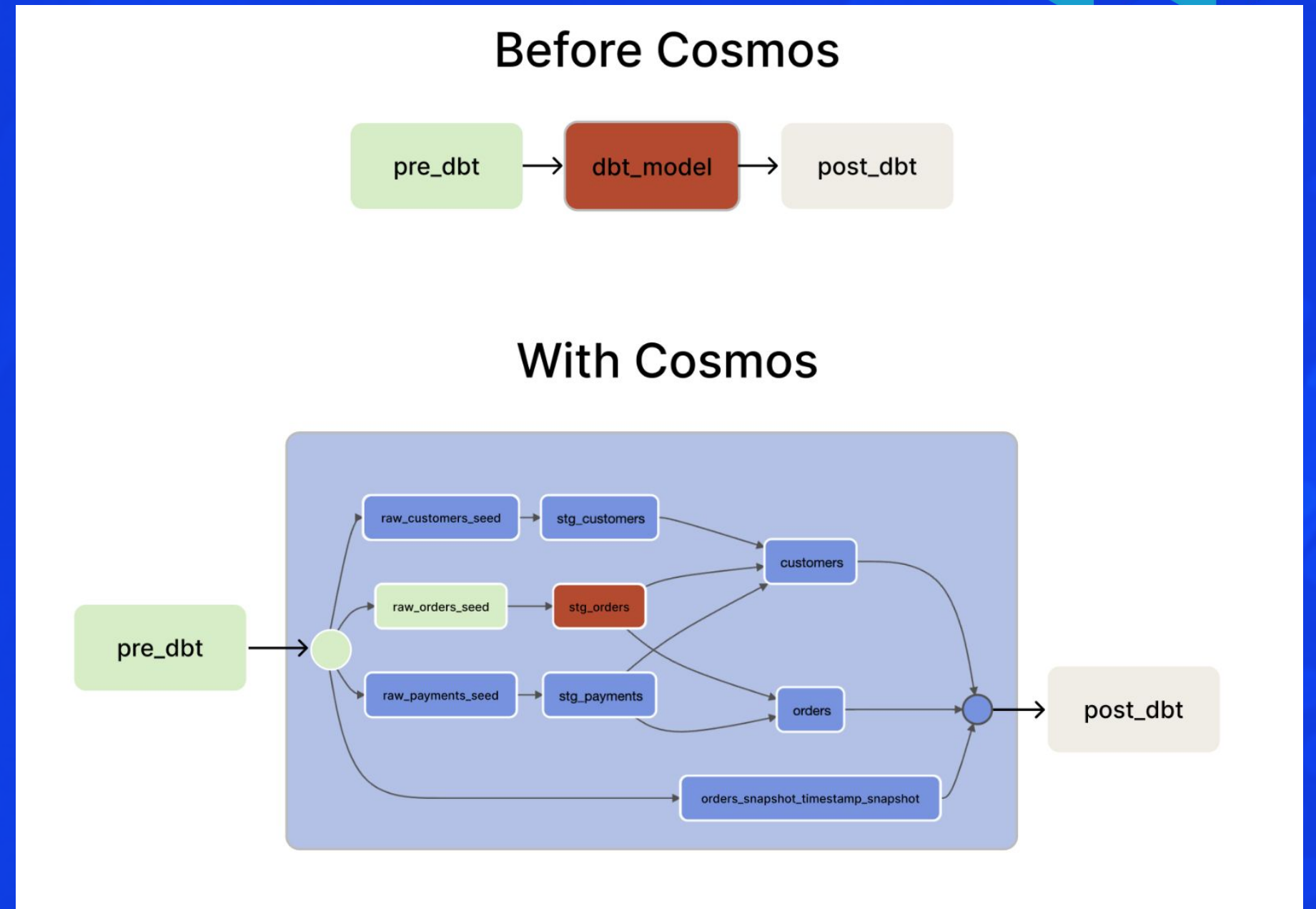
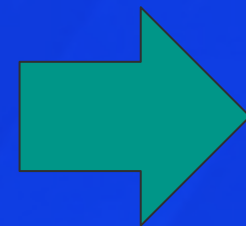
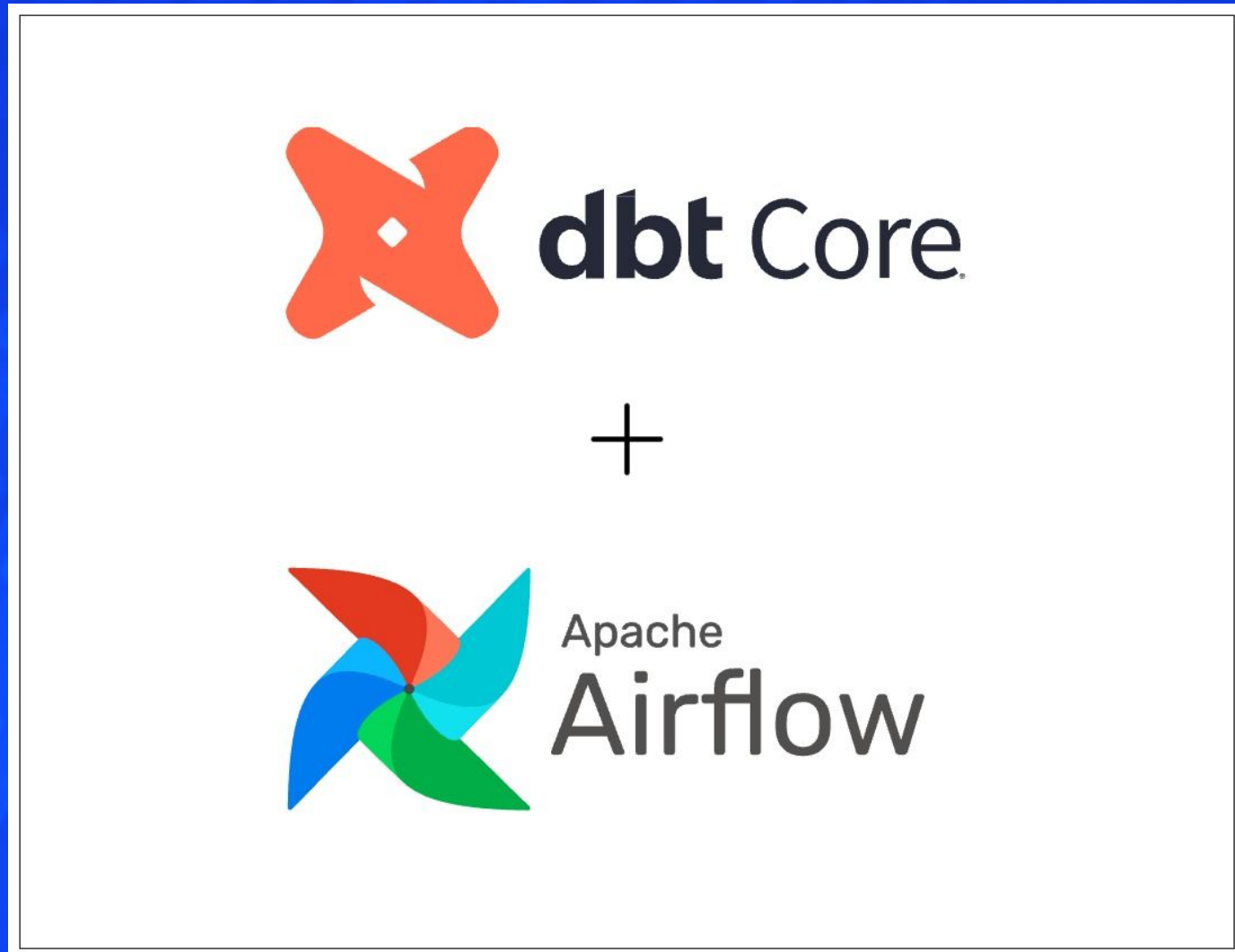
Open Telemetry



- Integrated in Airflow
- Adopted by everyone
- Still early days
- Traces, Log support in the works



Astronomer's Cosmos



Fully fledged REST API

🔍 Search...

- Overview >
- Trying the API >
- Authentication
- Errors >
- Config >
- Connection >
- DAG >
- DAGRun >
- EventLog >

Airflow API (Stable) (2.7.3)

Download OpenAPI specification: [Download](#)

Apache Software Foundation: dev@airflow.apache.org | URL: <https://airflow.apache.org>
License: Apache 2.0

<https://airflow.apache.org/docs/apache-airflow/stable/>

Overview

To facilitate management, Apache Airflow supports a range of REST API endpoints across its objects. This section provides an overview of the API design, methods, and supported use cases.



Engineering friendliness

- Workflow as a code front and center
- Tests
 - airflow task test
 - airflow dag test
 - unit test guidelines
 - system tests support
- Running Airflow locally
 - airflow standalone
 - docker compose
 - airflowctl - by Kaxil, Airflow PMC member

```
> airflowctl --help

Usage: airflowctl [OPTIONS] COMMAND [ARGS]...

Streamline getting started with Apache Airflow™ and managing multiple Airflow projects.

Options
  --install-completion  Install completion for the current shell.
  --show-completion     Show completion for the current shell, to copy it or customize the installation.
  --help                Show this message and exit.

Commands
  airflow  Forward commands to Airflow CLI.
  build    Build an Airflow project. This command sets up the project environment, installs Apache Airflow and its dependencies.
  info     Display information about the current Airflow project.
  init     Initialize a new Airflow project.
  list     List all Airflow projects created using this CLI.
  logs     Continuously display live logs of the background Airflow processes.
  start    Start Airflow.
  stop     Stop a running background Airflow process and its entire process tree.
```



Airflow IS Open Source

(and always will)

Community - part of Apache Software Foundation



- The largest project in ASF (for contributors count) >2700
- Licencing ASF, permissive licence (that will NEVER change)
- Well established, strong governance
- 61 committers, 32 PMC members
- Stakeholders/Managed services/Vendor neutrality
 - Astronomer, Amazon, Google, Microsoft, ...
- Security / Release process / Maintenance certainty

Tools integrating with Airflow

- DAG visual editors
- Declarative DAG authoring
- IDE integration
- CLIs to manage Airflow
- Debugging aids
- UI extensions
- ...

Tools integrating with Airflow

ADA - A microservice created to retrieve analytics metrics from an Airflow database instance.

as-scrapers - An integration with Selenium to build & maintain web scrapers inside Airflow.

afctl - A CLI tool that includes everything required to create, manage and deploy airflow projects faster and smoother.

airflint - Enforce Best Practices for all your Airflow DAGs.

airflow-aws-executors - Run Airflow Tasks directly on AWS Batch, AWS Fargate, or AWS ECS; provisioning less infra is more.

airflow-code-editor - A tool for Apache Airflow that allows you to edit DAGs in browser.

airflow-diagrams - Auto-generated Diagrams from Airflow DAGs

airflow-maintenance-dags - Clairvoyant has a repo of Airflow DAGs that operator on Airflow itself, clearing out various bits of the backing metadata store.

AirflowK8sDebugger - A library for generate k8s pod yaml templates from an Airflow dag using the KubernetesPodOperator.

Airflow Ditto - An extensible framework to do transformations to an Airflow DAG and convert it into another DAG which is flow-isomorphic with the original DAG, to be able to run it on different environments (e.g. on different clouds, or even different container frameworks - Apache Spark on YARN vs Kubernetes). Comes with out-of-the-box support for EMR-to-HDInsight-DAG transforms.

Amundsen - Amundsen is a data discovery and metadata platform for improving the productivity of data analysts, data scientists and engineers when interacting with data. It can surface which Airflow task generates a given table.

Apache-Liminal-Incubating - Liminal provides a domain-specific-language (DSL) to build ML/AI workflows on top of Apache Airflow. Its goal is to operationalise the machine learning process, allowing data scientists to quickly transition from a successful experiment to an automated pipeline of model training, validation, deployment and inference in production.

Astro CLI - The Astro CLI is the easiest way to get a local Airflow server for prototyping and development.

Astro SDK - Astro SDK allows rapid and clean development of Extract, Load, Transform workflows using Python and SQL, powered by Apache Airflow and maintained by Astronomer.

Chartis - Python package to convert Common Workflow Language (CWL) into Airflow DAG.

CWL-Airflow - Python package to extend Apache-Airflow 1.10.11 functionality with CWL v1.2 support.

dag-factory - A library for dynamically generating Apache Airflow DAGs from YAML configuration files.

Dag Dependencies viewer - A tool which creates a view to visualize dependencies between the Airflow DAGs

data-dag - A library for building factories to dynamically generate DAGs from data (such as YAML files)

Databand - Observability platform built on top of Airflow.

DataHub - A metadata platform for the modern data stack. It can automatically collect lineage and other metadata from Airflow.

dbt (data build tool) - Data transformation tool, [dbt jobs can be scheduled using Airflow](#).

Domino - Domino is an open source Graphical User Interface platform for creating data and Machine Learning workflows (DAGs) with no-code, visually intuitive drag-and-drop actions. It is also a standard for publishing and sharing your Python code so it can be automatically used by anyone, directly in the GUI.

Elyra - Elyra provides a visual editor that enables data scientists to create AI pipelines in a low-code/no-code fashion.

GeniumCloud - One-Stop-Shop Platform for rapid build, scheduling and control Airflow workflows via completely new UI. Out of the box comprehensive Airflow infrastructure monitoring, integration with alerting systems and service adoption from small to enterprise organizations. The easiest way to manage complex workflows.

gusty - Create a DAG using any number of YAML, Python, Jupyter Notebook, or R Markdown files that represent individual tasks in the DAG. gusty also configures dependencies, DAGs, and TaskGroups, features support for your local operators, and more. A fully containerized demo is available [here](#).

Marquez - Marquez is an open source metadata service that maintains data provenance, shows how datasets are consumed and produced and centralizes dataset lifecycle management. Marquez can be used with Apache Airflow as an OpenLineage backend.

Meltano - Open source, self-hosted, CLI-first, debuggable, and extensible ELT tool that embraces [Singer](#) for extraction and loading, leverages [dbt](#) for transformation, and [integrates with Airflow](#) for orchestration.

Nexla - Build, transform, and manage data flows to and from databases, APIs, streams, SaaS services, events, and even emails. Use Nexla's Airflow Operator to trigger flows to start in other Operators when your Nexla flow finishes running.

Oozie to Airflow - A tool to easily convert between [Apache Oozie](#) workflows and Apache Airflow workflows.

OpenLineage - An open standard for the collection of data lineage, which can be used to trace the path of datasets as they traverse multiple systems including Apache Airflow.

Panda Patrol - Test and profile your data right within your Airflow DAGs. With dashboards and alerts already pre-built.

Pylint-Airflow - A Pylint plugin for static code analysis on Airflow code.

Redactics - A managed appliance (built on Airflow) installed next to your databases that powers a growing collection of data management workflows.

simple-dag-editor - Zero configuration Airflow tool that let you manage your DAG files.

Viewflow - An Airflow-based framework that allows data scientists to create data models without writing Airflow code.

whirl - Fast iterative local development and testing of Apache Airflow workflows.

ZenML - Run your machine learning specific pipelines on Airflow, easily integrating with your existing data science tools and workflows.

Airflow Vscode Extension This is a VSCode extension for Apache Airflow 2+. You can trigger your DAGs, pause/unpause DAGs, view execution logs, explore source code and do much more.

Airflow Provider Template - Template and commands for creating and testing airflow provider packages.

Airflow Template - Template and commands for creating minimal airflow environments for rapid testing and prototyping.





Solid Infrastructure

Public Interface of Airflow



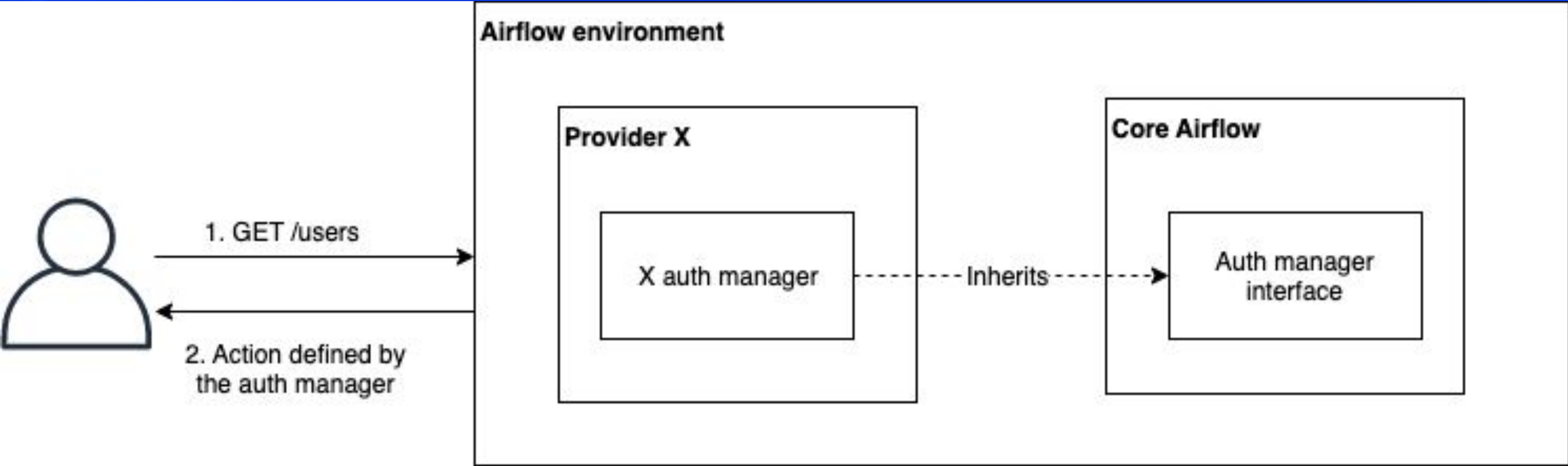
A screenshot of the Apache Airflow documentation website. The page title is 'Public Interface of Airflow'. The main content area contains a paragraph describing the Public Interface as a set of interfaces for developers to interact with Airflow features. Below this is a section titled 'Using Airflow Public Interfaces' which lists three bullet points: extending Airflow classes (Operators and Hooks), writing new Plugins, and bundling custom Operators, Hooks, and Plugins into provider packages. A sidebar on the left shows a navigation menu with 'Public Interface of Airflow' selected. A sidebar on the right lists various sub-topics like DAGs, Operators, Task Instances, etc. The top navigation bar includes links for Community, Meetups, Documentation, Use-cases, Announcements, Blog, and Ecosystem. The Apache Airflow logo is in the top left corner of the page.



Providers

- Can upgrade/downgrade separately
- Can provide:
 - Hooks/Operators/Sensors, Extra-links, Connection types
 - Secret Backends, Triggers, Log Handlers,
 - Executors, Notifications, Configuration, Decorators
 - Filesystems (2.8)
- Full lifecycle of providers defined
 - Approval by community (or not)
 - Support lifecycle for multiple Airflow versions
 - Suspension/Resuming/Removal
- 3rd-party providers and registries

Extensible user management





Security - coming soon for everyone

- Regulations are coming (CRA act just agreed in EU Trilogue)
- Airflow is part of the HackerOne OSS Bounty
- Highly functional Security Team ~50 reports handled
- 4 Airflow contributors: Sovereign Tech Fund funding for security
 - Security Model and Security Policy
 - SBOM generated
 - Securing release process (reproducible builds)
 - Component Isolation (Multi-tenancy in progress)

Summary



- Airflow is a modern, solid orchestrator with strong foundations
- New, slick ways to interact with the Modern Data Stack
- True Open Source
- Community is huge, strong and supportive
- More, exciting things are coming. Fast.



Q&A

<https://github.com/potiuk>

<https://www.linkedin.com/in/jarekpotiuk>