# OpenLineage: from Operators to Hooks

# Maciej Obuchowski

## Astronomer

Maciej is a software engineer, Airflow and OpenLineage commiter. He loves rock climbing, contributing to open source data projects and playing with cats.

## Sessions by Maciej Obuchowski

- OpenLineage: From Operators to Hooks (2024)
- OpenLineage in Airflow: A Comprehensive Guide (2023)
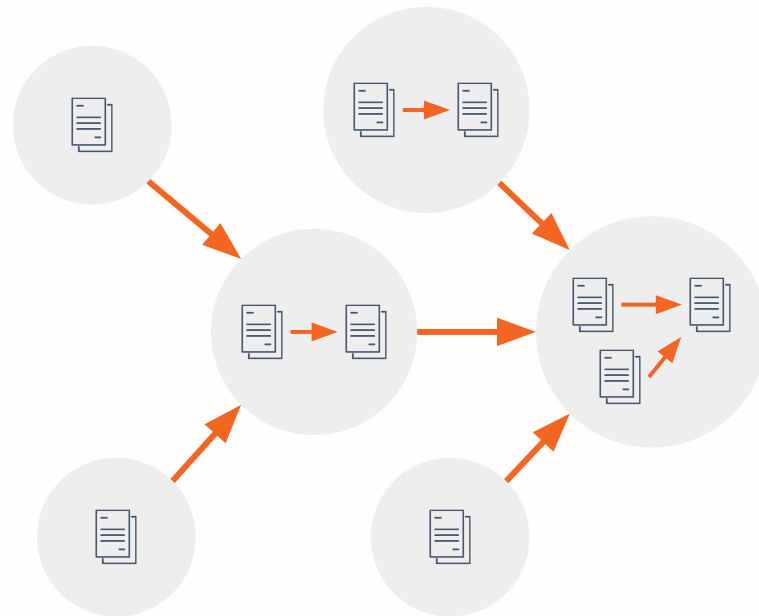- OpenLineage & Airflow - data lineage has never been easier (2022)

# Agenda

- What is OpenLineage
- OpenLineage Airflow Integration
- Getting Lineage From Hooks
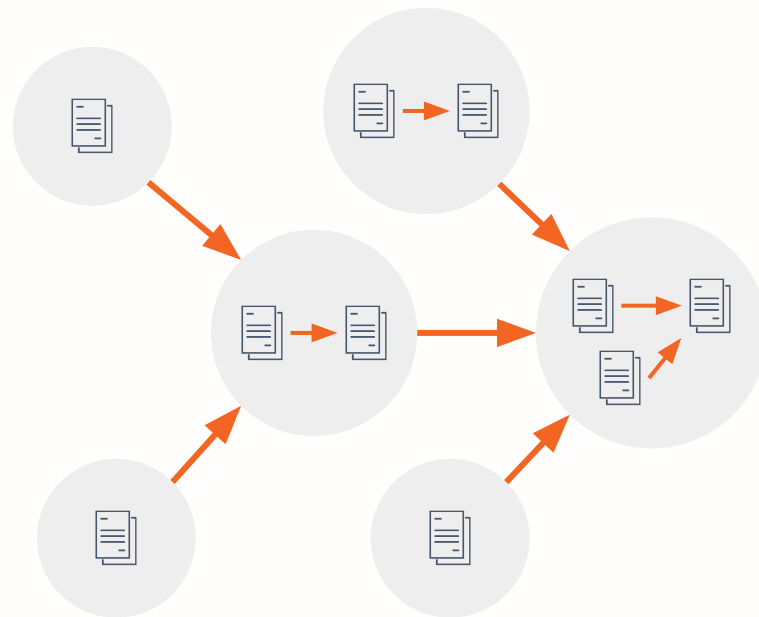- Peek into the future?

# What is Data Lineage?

Data lineage is the set of complex relationships between datasets and jobs in data pipelines.
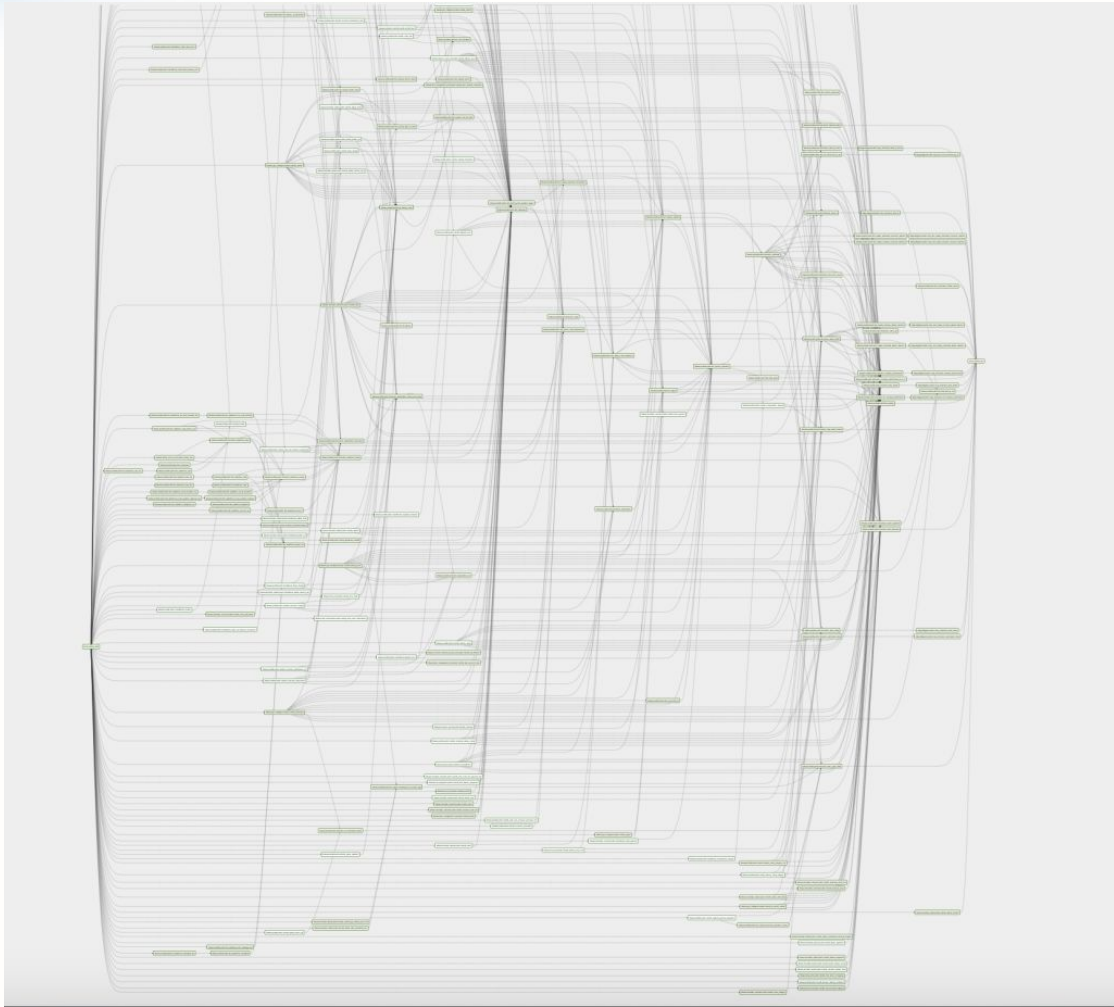
- Producers & consumers of each dataset

- Inputs and outputs of each job

ASTRONOMER

# What problems Data Lineage solves

- Holistic view on data flowing through organization

- Security and Compliance

  - A German bank suffered data breach from vendor - but was wholly unaware that PII data was being send there.

  - Prod data send to dev env

- Impact analysis - pipeline failed, which other datasets it affects

# Working with data in 2024

# Open Lineage

**Mission:**

To define an **open standard** for the collection of lineage metadata from pipelines **as they are running.**

LF AI & DATA

# The Data World Without OpenLineage

# The Data World With OpenLineage



https://openlineage.io/ecosystem

# Why runtime?



You can try to infer the date and location of an image after the fact...



...or you can capture it when the image is originally created!

OpenLineage

ASTRONOMER

# OpenLineage Integrations

**Metadata producers**

**Metadata consumers**

# OpenLineage Contributors

ASTRONOMER

getindata
Part of Xebia

pandas

DECATHLON

Microsoft

APACHE Spark

dbt

Amundsen

Parquet

snowflake

ODPI EGERIA

ICEBERG

Apache Airflow

MARQUEZ

matillion

asana

Natural Intelligence

Bloomberg®

Booking.com

OpenLineage

ASTRONOMER

# Astro Observe

# Astro Observe

# What does it do

# What does it do

# 2.7+ AIP-53 Implementation

- OpenLineage is part of Airflow since 2.7+ introduced in AIP-53
- Part of implementation happens in Operators
- START, COMPLETE, FAIL states are exposed via different get_openlineage_facets_* methods returning OperatorLineage class

Ex. on BigQueryToGCSOperator

```python
def get_openlineage_facets_on_complete(self, task_instance):
    parsed_uri = urllib.parse.urlparse(self.destination_cloud_storage_uris[0])
    return OperatorLineage(
        inputs=[Dataset(
            namespace="bigquery",
            name=f"{self.conf['datasetId']}.{self.conf['projectId']}.{self.conf['tableId']}")
        ],
        outputs=[
            # Caveat: single URI
            Dataset(namespace=f"gs://{parsed_uri.hostname}", name=parsed_uri.path)
        ]
    )
```

# How does it work

- OpenLineage Provider has listener that gets called by Airflow by Listener API
  - On worker, before task runs, and after task completes or fails
- Checks if Operator has implemented lineage methods, calls them
- Enrich data returned from Operator with common Airflow facets
- More on that in last year's presentation https://youtu.be/SZBVgREqets

# Supported Operators

https://airflow.apache.org/docs/apache-airflow-providers-openlineage/stable/supported_classes.html
Caveat - does not make sense for OpenLineage to cover all Operators

# Is this enough?

# Airflow Survey 2023 results

**Which features would you like to see in Airflow? (multiple choice)**

**52.2%** (391)
DAG versioning

**34.4%** (258)
More data lineage

**30%** (225)
Multi-Tenancy

**28.6%** (214)
Submitting new DAGs externally via API

**26.4%** (198)
Better security (isolation)

**25.1%** (188)
More support for datasets
and data-driven scheduling

**24.2%** (181)
Data cataloguing

**22.6%** (169)
Support for native cloud executors
(AWS/GCP/Azure etc.)

# How people use Airflow?

what operators are most popular in Apache Airflow?

Here are some of the most popular operators in Apache Airflow:

## 1. PythonOperator

- **Description**: Executes Python functions or callables, offering great flexibility and is widely used for various tasks.

- **Use Cases**: Data processing, API integrations, and executing custom Python logic.

# TaskFlow API

```python
from airflow.decorators import dag, task
from airflow.utils.dates import days_ago
from airflow.providers.amazon.aws.hooks.s3 import S3Hook

@dag(schedule_interval=None, start_date=days_ago(1), catchup=False)
def s3_data_processing_dag():

    @task
    def read_data_from_s3(bucket_name: str, key: str) -> str:
        s3_hook = S3Hook(aws_conn_id='my_aws_conn')
        data = s3_hook.read_key(key=key, bucket_name=bucket_name)
        return data

    @task
    def process_data(data: str) -> str:
        processed_data = data.upper()
        return processed_data

    @task
    def write_data_to_s3(bucket_name: str, key: str, data: str):
        s3_hook = S3Hook(aws_conn_id='my_aws_conn')
        s3_hook.load_string(string_data=data, key=key, bucket_name=bucket_name, replace=True)

    raw_data = read_data_from_s3(bucket_name='my_bucket', key='input_data.txt')
    processed_data = process_data(raw_data)
    write_data_to_s3(bucket_name='my_bucket', key='output_data.txt', data=processed_data)
```

# Object Storage

```python
from __future__ import annotations
from datetime import datetime

from airflow.decorators import dag, task
from airflow.io.path import ObjectStoragePath


base = ObjectStoragePath("s3://aws_default@openlineage-test/")


@dag(
    schedule=None,
    start_date=datetime(2022, 1, 1),
    catchup=False,
)
def objectstorage_lineage():
    @task
    def read_s3_data(**kwargs):
        from pandas import DataFrame as pd
        storage_path = base / "input" / "dataset"
        with storage_path.open("r") as file:
            data = file.read()

        df = pd.DataFrame(data)
        return df

    @task
    def transform_data(df):
        ...

    @task
    def write_to_s3(df):
        storage_path = base / "output" / "dataset"
        with storage_path.open("w") as file:
            df.to_parquet(file)

    data = read_s3_data()
    transformed_data = transform_data(data)
    write_to_s3(transformed_data)
```

# Let's look at another DAG

# Let's look at another DAG

# The problem: arbitrary code

- People love writing their own code!
- Non-code operators: we generally know what they do (*)
- Code operators: can be anything!
- Is there anything we can do?

```python
from airflow import DAG
from airflow.operators.python import PythonOperator
from datetime import datetime
from my_awesome_code import do_the_work

def extract_load_transform(**kwargs):
    # What happens there?
    do_the_work()

default_args = {
    'owner': 'airflow',
    'start_date': datetime(2023, 9, 1),
    'retries': 1,
}

with DAG(
    dag_id='example_python_operator',
    default_args=default_args,
    schedule_interval='@daily',
    catchup=False,
) as dag:
    the_task = PythonOperator(
        task_id='extract_load_transform',
        python_callable=extract_load_transform,
        provide_context=True,
    )
```

ASTRONOMER

# Instrument the hooks

# What are hooks?

- Hooks - a mechanism for communicating with external systems
- Operations on hooks closely resemble those exposed by external systems
- Ex. copy file between object storage buckets, execute SQL on some database, or even create a Kubernetes job

```python
def copy(
    self,
    source_bucket: str,
    source_object: str,
    destination_bucket: str | None = None,
    destination_object: str | None = None,
) -> None:
    """
    Copy an object from a bucket to another, with renaming if requested.

    destination_bucket or destination_object can be omitted, in which case
    source bucket/object is used, but not both.

    :param source_bucket: The bucket of the object to copy from.
    :param source_object: The object to copy.
    :param destination_bucket: The destination of the object to copied to.
        Can be omitted; then the same bucket is used.
    :param destination_object: The (renamed) path of the object if given.
        Can be omitted; then the same name is used.
    """
    destination_bucket = destination_bucket or source_bucket
    destination_object = destination_object or source_object

    if source_bucket == destination_bucket and source_object == destination_object:
        raise ValueError(
            f"Either source/destination bucket or source/destination object must be different, "
            f"not both the same: bucket={source_bucket}, object={source_object}"
        )
    if not source_bucket or not source_object:
        raise ValueError("source_bucket and source_object cannot be empty.")

    client = self.get_conn()
    source_bucket = client.bucket(source_bucket)
    source_object = source_bucket.blob(source_object)  # type: ignore[attr-defined]
    destination_bucket = client.bucket(destination_bucket)
    destination_object = source_bucket.copy_blob(  # type: ignore[attr-defined]
        blob=source_object, destination_bucket=destination_bucket, new_name=destination_object
    )

    self.log.info(
        "Object %s in bucket %s copied to object %s in bucket %s",
        source_object.name,  # type: ignore[attr-defined]
        source_bucket.name,  # type: ignore[attr-defined]
        destination_object.name,  # type: ignore[union-attr]
        destination_bucket.name,  # type: ignore[union-attr]
    )
```

ASTRONOMER

# How to instrument hooks?

- We can't copy the "pull" approach we have for Operators
- One method vs many methods
- We can "push" it from the instrumented method
- Where to?

```python
def copy(
    self,
    source_bucket: str,
    source_object: str,
    destination_bucket: str | None = None,
    destination_object: str | None = None,
) -> None:
    """
    Copy an object from a bucket to another, with renaming if requested.

    destination_bucket or destination_object can be omitted, in which case
    source bucket/object is used, but not both.

    :param source_bucket: The bucket of the object to copy from.
    :param source_object: The object to copy.
    :param destination_bucket: The destination of the object to copied to.
        Can be omitted; then the same bucket is used.
    :param destination_object: The (renamed) path of the object if given.
        Can be omitted; then the same name is used.
    """
    destination_bucket = destination_bucket or source_bucket
    destination_object = destination_object or source_object

    if source_bucket == destination_bucket and source_object == destination_object:
        raise ValueError(
            f"Either source/destination bucket or source/destination object must be different, "
            f"not both the same: bucket={source_bucket}, object={source_object}"
        )
    if not source_bucket or not source_object:
        raise ValueError("source_bucket and source_object cannot be empty.")

    client = self.get_conn()
    source_bucket = client.bucket(source_bucket)
    source_object = source_bucket.blob(source_object)  # type: ignore[attr-defined]
    destination_bucket = client.bucket(destination_bucket)
    destination_object = source_bucket.copy_blob(  # type: ignore[attr-defined]
        blob=source_object, destination_bucket=destination_bucket, new_name=destination_object
    )

    self.log.info(
        "Object %s in bucket %s copied to object %s in bucket %s",
        source_object.name,  # type: ignore[attr-defined]
        source_bucket.name,  # type: ignore[attr-defined]
        destination_object.name,  # type: ignore[union-attr]
        destination_bucket.name,  # type: ignore[union-attr]
    )
```

# What does the instrumentation look like?

- Added HookLineageCollector on Worker with .add_input_dataset and .add_output_dataset methods
- Those methods accept dataset_kwargs that are used to construct the compliant dataset later
- Methods in hooks that modify datasets call those methods to register dataset changes
- Data is then deduplicated - we don't want to see hundreds of writes to same dataset

```python
@unify_bucket_name_and_key
@provide_bucket_name
def load_file(
    self,
    filename: Path | str,
    key: str,
    bucket_name: str | None = None,
    replace: bool = False,
    encrypt: bool = False,
    gzip: bool = False,
    acl_policy: str | None = None,
) -> None:
    filename = str(filename)
    if not replace and self.check_for_key(key, bucket_name):
        raise ValueError(f"The key {key} already exists.")

    extra_args = self.extra_args
    if encrypt:
        extra_args["ServerSideEncryption"] = "AES256"
    if gzip:
        with open(filename, "rb") as f_in:
            filename_gz = f"{f_in.name}.gz"
            with gz.open(filename_gz, "wb") as f_out:
                shutil.copyfileobj(f_in, f_out)
                filename = filename_gz
    if acl_policy:
        extra_args["ACL"] = acl_policy

    client = self.get_conn()
    client.upload_file(filename, bucket_name, key, ExtraArgs=extra_args, Config=self.transfer_conf
    get_hook_lineage_collector().add_input_dataset(
        context=self, scheme="file", dataset_kwargs={"path": filename}
    )
    get_hook_lineage_collector().add_output_dataset(
        context=self, scheme="s3", dataset_kwargs={"bucket": bucket_name, "key": key}
    )
```

# What does the instrumentation look like?

- Object Storage: track reads and writes across file-like objects
- No additional instrumentation needed for different types of paths

```python
44  class TrackingFileWrapper(LoggingMixin):
45      """Wrapper that tracks file operations to intercept lineage."""
46
47      def __init__(self, path: ObjectStoragePath, obj):
48          super().__init__()
49          self._path: ObjectStoragePath = path
50          self._obj = obj
51
52      def __getattr__(self, name):
53          attr = getattr(self._obj, name)
54          if callable(attr):
55              # If the attribute is a method, wrap it in another method to intercept the call
56              def wrapper(*args, **kwargs):
57                  self.log.debug("Calling method: %s", name)
58                  if name == "read":
59                      get_hook_lineage_collector().add_input_dataset(context=self._path, uri=str(self._path))
60                  elif name == "write":
61                      get_hook_lineage_collector().add_output_dataset(context=self._path, uri=str(self._path))
62                  result = attr(*args, **kwargs)
63                  return result
64
65              return wrapper
66          return attr
67
68      def __getitem__(self, key):
69          # Intercept item access
70          return self._obj[key]
71
72      def __enter__(self):
73          self._obj.__enter__()
74          return self
75
76      def __exit__(self, exc_type, exc_val, exc_tb):
77          self._obj.__exit__(exc_type, exc_val, exc_tb)
```

ASTRONOMER

# Works only if you use it

- Not OpenLineage specific - you can write your own plugin, register HookLineageReader and use the gathered hook lineage for your own purposes
- If there are no HookLineageReaders, the .add_input_datasets and .add_output_datasets just send data to /dev/null

```python
class HookLineageReader(LoggingMixin):
    """Class used to retrieve the hook lineage information collected by HookLineageCollector."""

    def __init__(self, **kwargs):
        self.lineage_collector = get_hook_lineage_collector()

    def retrieve_hook_lineage(self) -> HookLineage:
        """Retrieve hook lineage from HookLineageCollector."""
        hook_lineage = self.lineage_collector.collected_datasets
        return hook_lineage


def get_hook_lineage_collector() -> HookLineageCollector:
    """Get singleton lineage collector."""
    global _hook_lineage_collector
    if not _hook_lineage_collector:
        from airflow import plugins_manager

        plugins_manager.initialize_hook_lineage_readers_plugins()
        if plugins_manager.hook_lineage_reader_classes:
            _hook_lineage_collector = HookLineageCollector()
        else:
            _hook_lineage_collector = NoOpCollector()
    return _hook_lineage_collector
```

# AIP-60 compliant datasets

- Operators don't always work on datasets they know - GCSToS3Operator
- Object Storage does not explicitly know what the file is, it operates on abstract files
- AIP-60 Dataset is an URI, OpenLineage dataset is name+namespace
- Providers know the dataset abstractions they own
- Providers can provide good factory method for AIP-60 URI, and how to translate it to OpenLineage dataset

```
dataset-uris:
  - schemes: [s3]
    handler: airflow.providers.amazon.aws.datasets.s3.sanitize_uri
    to_openlineage_converter: airflow.providers.amazon.aws.datasets.s3.convert_dataset_to_openlineage
    factory: airflow.providers.amazon.aws.datasets.s3.create_dataset
```

# Same Object Storage DAG in Airflow 2.10

# Same Object Storage DAG in Airflow 2.10

# Lineage For Python Tasks

- It does not matter how you author your Python tasks, it captures hook lineage whether you use PythonOperator, TaskFlow or Custom Operators
- However, it has to be the same process as worker - python-based KubernetesOperator job won't work

# Peek into the future

As in, what's happening with OL in Airflow 3?

# ~~Datasets~~ Assets

- Despite all those changes, not everything can be detected.
- OpenLineage integration will be able to take advantage of that and expose annotated lineage

Pages / Airflow Home / Airflow Improvement Proposals                                        ...

## AIP-73 Expanded Data Awareness

Created by Constance Martineau, last modified by Tzu-ping Chung on Aug 02, 2024

### Status

| State | Accepted |
| --- | --- |
| **Discussion Thread** | https://lists.apache.org/thread/6rp4jhflwg3czhtvjszoctdry85vfv8r |
| **Vote Thread** | https://lists.apache.org/thread/9570fr5b2jv6hb2fd5z43jmsws42ls1z |
| **Vote Result Thread** | https://lists.apache.org/thread/f0rvbdbpq7bylt3kv0v6gn90qr3f98ng |
| **Progress Tacking (PR/GitHub Project/Issue Label)** | |
| **Date Created** | 2024-06-26 |
| **Version Released** | |
| **Authors** | @Constance Martineau   @Tzu-ping Chung |

### Motivation

Airflow has become the standard for orchestrating complex data workflows. However, it operates with limited visibility into the actual data it processes or produces. While it understands task execution order and attributes like operators and parameters in use, it lacks insight into the nature of data inputs and outputs. This link between data and tasks is fundamental to data engineering, and vital for providing insights into the state and health of data as they move through the workflow. Orchestrators are the heart of data platforms, and if they can understand this link, they can make orchestration decisions based on data quality and freshness, while also providing data engineers with insights about system and data reliability in one place.

# ~~Datasets~~ Assets

## Seeing Clearly with Airflow: Bridging Task-Centric and Data-Aware Orchestration

Speaker(s):

*Sep-12 12:30-13:15 in Elizabethan A+B*   📅 **Add to Calendar**

As Apache Airflow evolves, a key shift is emerging: the move from task-centric to data-aware orchestration. Traditionally, Airflow has focused on managing tasks efficiently, with limited visibility into the data those tasks manipulate. However, the rise of data-centric workflows demands a new approach—one that puts data at the forefront.

This talk will explore how embedding deeper data insights into Airflow can align with modern users' needs, reducing complexity and enhancing workflow efficiency. We'll discuss how this evolution can transform Airflow into a more intuitive and powerful tool, better suited to today's data-driven environments.

Constance Martineau

Tzu Ping Chung

# Synergy

- Building on Airflow - using Connections, Hooks, Assets, Object Storage gives you much more than using Airflow as a "dumb" scheduler

# AIP-72 Task Execution Interface

- Decoupling access to database from Worker requires us to rework approach for integration
- Below diagram won't be accurate: worker after completion will send metadata using AIP-72 interface, and OL will process it asynchronously

# AIP-72 Task Execution Interface

- Gains:
  - Async event emission: does not block a worker slot
  - Isolated execution: won't affect worker if some bug happen
  - Possible to enrich the data more due to asynchronous execution
- However, it's possible that contributing OpenLineage Operator implementations will get more complex

ASTRONOMER

# What AIP-72 enables

## Running Airflow Tasks Anywhere, in any Language

Speaker(s):



Ash Berlin-Taylor



Vikram Koka

*Sep-10 12:00-12:45 in Elizabethan A+B*   📅 Add to Calendar

Imagine a world where writing Airflow tasks in languages like Go, R, Julia, or maybe even Rust is not just a dream but a native capability. Say goodbye to BashOperators; welcome to the future of Airflow task execution.

Here's what you can expect to learn from this session:

- Multilingual Tasks: Explore how we empower DAG authors to write tasks in any language while retaining seamless access to Airflow Variables and Connections.
- Simplified Development and Testing: Discover how a standardized interface for task execution promises to streamline development efforts and elevate code maintainability.
- Enhanced Scalability and Remote Workers: Learn how enabling tasks to run on remote workers opens up possibilities for seamless deployment on diverse platforms, including Windows and remote Spark or Ray clusters. Experience the convenience of effortless deployments as we unlock new avenues for Airflow usage.

Join us as we embark on an exploratory journey to shape the future of Airflow task execution. Your insights and contributions are invaluable as we refine this vision together. Let's chart a course towards a more versatile, efficient, and accessible Airflow ecosystem.

OpenLineage

ASTRONOMER

# Related event: OpenLineage Meetup

OpenLineage

Get more details and the signup link at **https://openlineage.io/blog**

## Where

Astronomer
8 California St.

## When

Thursday, 9/12
6-9 pm

**Join us for in-depth talks & discussion over dinner!**

## Agenda

- **Unlocking Data Products with OpenLineage at Astronomer**: Julian LaNeve and Jason Ma, Astronomer.
- **OpenLineage: From Operators to Hooks** by Maciej Obuchowski, Astronomer+GetInData/Xebia.
- **Activating Operational Metadata with Airflow, Atlan and Openlineage** by Kacper Muda, GetInData/Xebia.
- **Hamilton, a Scaffold for all Your Python Platform Concerns (and a New OpenLineage Producer)** by Stefan Krawczyk
- **Lightning Talk on New Marquez Features and the Marquez Project Roadmap** by Willy Lulciuc, Marquez Lead, and Peter Hicks, Marquez Committer.

ASTRONOMER

# Thank you!
# Any questions?

# Bonus slides

# Consistent API introduction

- Adding new APIs is easy, right?
  - Add new method to Airflow Core
  - Make providers use it
  - Release it all
  - Then it breaks on older versions of Airflow

# Consistent API introduction

- Adding new APIs is easy, right?
  - Add new method to Airflow Core
  - Make providers use it
  - Release it all
  - Then it breaks on older versions of Airflow
- Make provider check Airflow version, and perform an action only if version matches the requirement
  - Requires potentially a lot of duplicated work across providers

# Consistent API introduction

- Adding new APIs is easy, right?
  - Add new method to Airflow Core
  - Make providers use it
  - Release it all
  - Then it breaks on older versions of Airflow
- Make provider check Airflow version, and perform an action only if version matches the requirement
  - Requires potentially a lot of duplicated work across providers
- Airflow's solution: common.compat provider
  - Introduce a method within compat provider
  - Method should perform version check, and have fallback for older version