# A GenAI Pipeline for Content Generation with Apache Airflow

May 30th 2024

ASTRONOMER

# Agenda

- **Quick introduction:** Airflow and Astronomer

- **Challenges** of GenAI pipelines and how Airflow addresses them

- **Key Airflow features** used in the demo
    - Focus on the latest 2.9 features:
        - Advanced Dataset Scheduling
        - Dynamic Task Mapping

- **Demo:** Fine-tune GPT for an RAG pipeline for content generation
        https://github.com/astronomer/gen-ai-fine-tune-rag-use-case

ASTRONOMER

# Airflow is the open standard for Workflow Management.

**22M+**
monthly downloads

**46k+**
Slack members

**2800+**
contributors

**1600+**
building blocks

GROWTH

COMMUNITY

INNOVATION

ECOSYSTEM

# More and more people are using Airflow for ML/AI



Ingestion and ETL/ELT related to analytics — **90%**

Ingestion and ETL/ELT related to business operations — **68%**

Training, serving, or generally manage MLOps — **28%**

Spinning up and spinning down infrastructure — **13%**

Other — **3%**

0%  25%  50%  75%  100%

*Source: 2023 Apache Airflow Survey, n=797*

# 28%

Of Airflow survey respondents in 2023 said they use Airflow for at least one ML/AI related use case.

ASTRONOMER

# Challenges when creating GenAI pipelines

The prototype works great - but production is a different beast

- API outages and rate limits
- Need to keep training data up to date
  - Your data is what sets you apart from competitors!
- Changing tools and APIs - new models coming out every day
- Complex pipeline structures
- Need the ability to determine which data went into training (compliance!)
- Scalability
- Reliability
- …

# Challenges when creating GenAI pipelines

The prototype works great - but production is a different beast

- API outages and rate limits -> Automatic retries
- Need to keep training data up to date -> Airflow already the standard
- Changing tools and APIs -> Airflow is tool agnostic, TaskFlow API
- Complex pipeline structures -> Datasets, dynamic task mapping, branching
- Need the ability to determine which data went into training (compliance!) -> Observability + OpenLineage integration
- Scalability -> Pluggable compute
- Reliability -> Battle tested + it is all code: CI/CD and DevOps best practices

Your data **+ your orchestration** is what sets you apart from competitors!

# Key Airflow features for GenAI

These features build a good foundation for best practice GenAI pipelines

- TaskFlow API
- Automatic retries
- Branching
- Deferrable operators
- Data-driven scheduling using Datasets
- Dynamic task mapping
- Alerts and notifications
- Setup and teardown tasks
- Backfills and reruns

ASTRONOMER

# TaskFlow API ⇒ Airflow decorators

The pythonic way to write Airflow DAGs

```python
from airflow.operators.python import PythonOperator

def say_hi_func(name: str = "") -> str:
    return f"Hi {name}!"

say_hi_obj = PythonOperator(
    task_id="say_hi_2",
    python_callable=say_hi_func,
    op_args=["Astra"],
)
```

⇒

```python
from airflow.decorators import task

@task
def say_hi_1(name: str = "") -> str:
    return f"Hi {name}!"

say_hi_obj = say_hi_1("Astra")
```
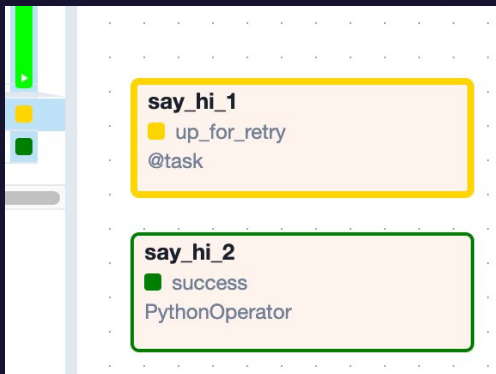
You can mix traditional operators and Airflow decorators!

There are many decorators: @dag, @task.kubernetes,
@task.branch, @task.bash etc… see: https://astronomer.io/docs/learn/airflow-decorators

ASTRONOMER

# Automatic retries in Airflow

Protects pipelines against rate-limits and API failures



You can configure:
- Number of retries
- Delay between retries
- Exponential backoff
- Maximum delay

Ways to configure:
- Airflow config
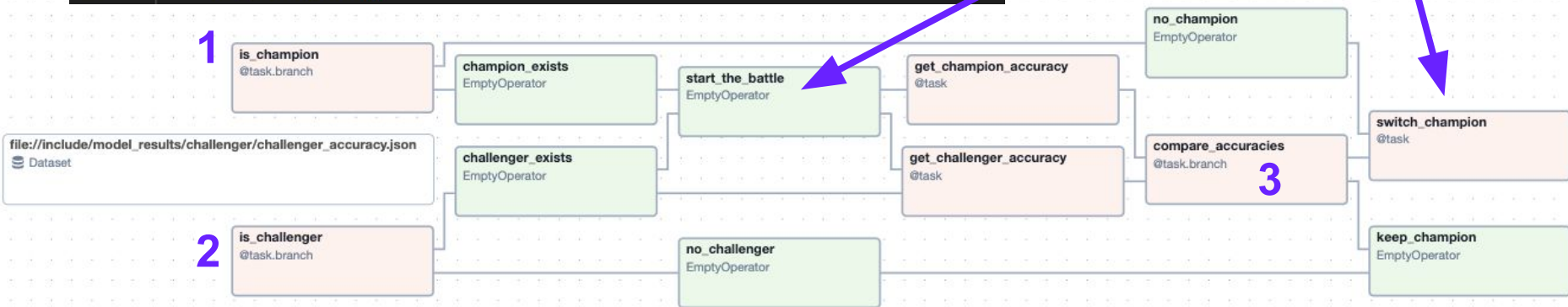- `default_args` in DAGs
- Individual tasks

**Best practice**: Always set retries in production, unless a task has a reason not to

Details: https://astronomer.io/docs/learn/rerunning-dags#automatically-retry-tasks

ASTRONOMER

# Branching in Airflow

```python
@task.branch
def is_champion() -> str:
    is_champion = # logic to determine if champion exists
    if is_champion:
        return "champion_exists"
    else:
        return "no_champion"

is_champion()
```

Careful with downstream trigger rules!

ASTRONOMER

# Deferrable operators

- Deferrable operators can start async processes in the **Triggerer** component.
- Use case:
    - Waiting for a long running process to finish (e.g. model training)
    - Waiting for an event to occur in an external system (like a sensor)
- Advantage:
    - The worker slot is released = resource use optimization

- **Best practice**: Use deferrable operators whenever possible for longer tasks.



Details: https://astronomer.io/docs/learn/deferrable-operators

# Dataset scheduling



**DAG with producer task**

**Dataset**

**Consumer DAG + producer task**

**Next Dataset**

# Datasets in the Airflow 2.9 UI

**Consumer DAG**



**Dataset**

**Producer task**

**Next Dataset**

# Advanced Dataset scheduling

Airflow 2.9 additions:

- Schedule on logical dataset expressions
  - Use AND (&) / OR (|) to create dataset logic


- Schedule on both time and datasets
  - `DatasetOrTimeSchedule` takes a `timetable` and a `dataset` argument


- REST API endpoint to update Datasets
  - Use for cross-deployment dependencies

Details: https://astronomer.io/docs/learn/airflow-datasets

# Dynamic Task Mapping

- Create a variable number of copies of the same task based on input at runtime!
- Define parameters that stay the same (`.partial()`) and parameters that change in between task instances (`.expand()` / `.expand_kwargs()`)

- **Best practice:**
  - Use dynamic tasks when possible over dynamic DAGs
  - Customize the map index (Airflow 2.9)

# Dynamic Task Mapping

Basic:

- `.partial(a=2)` → all parameters that stay the **same** for each mapped instance

- `.expand(b=[0,1])` → the parameter that changes as a list. Naming the kwarg is mandatory!

- `map_index_template` → customize the map index displayed in the **UI** (2.9)

Advanced:

- .expand_kwargs([{"a":1}]) → map over sets of keyword arguments
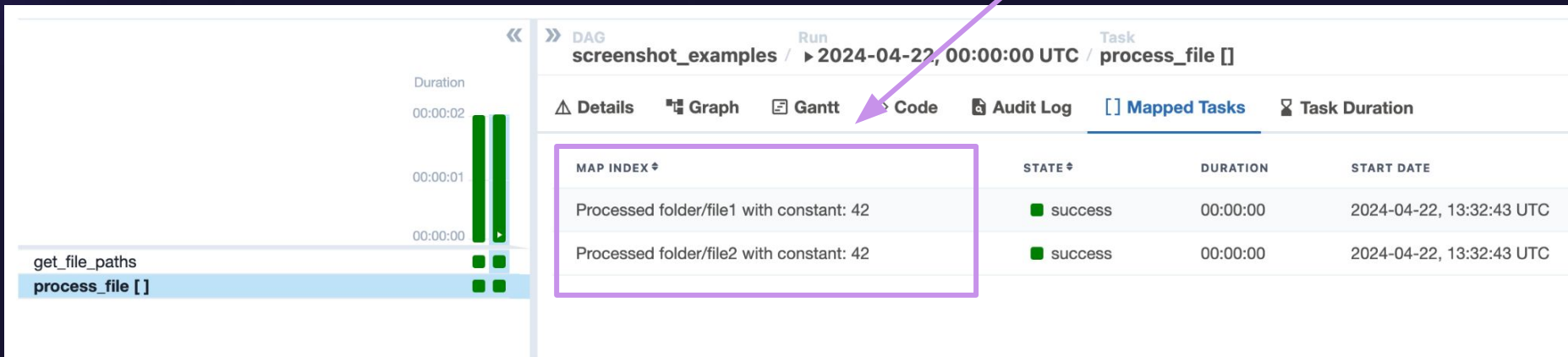- .map(lambda x: x) → transform the output of an upstream task before mapping over it

Details: https://astronomer.io/docs/learn/dynamic-tasks

ASTRONOMER

# Dynamic Task Mapping - Simple example

```python
44    @task
45    def get_file_paths() -> str:
46        # logic to get file paths. (potentially)
47        # results in differnet number of files each run
48        return ["folder/file1", "folder/file2"]
49
50    @task(map_index_template="{{ my_custom_map_index }}")
51    def process_file(constant: int, file: str) -> None:
52        # logic to process file
53
54        # create the custom map index
55        from airflow.operators.python import get_current_context
56
57        context = get_current_context()
58        context["my_custom_map_index"] = f"Processed {file} with constant: {constant}"
59
60    file_paths = get_file_paths()
61    processed_files = process_file.partial(constant=42).expand(file=file_paths)
62
```

# Dynamic Task mapping custom index
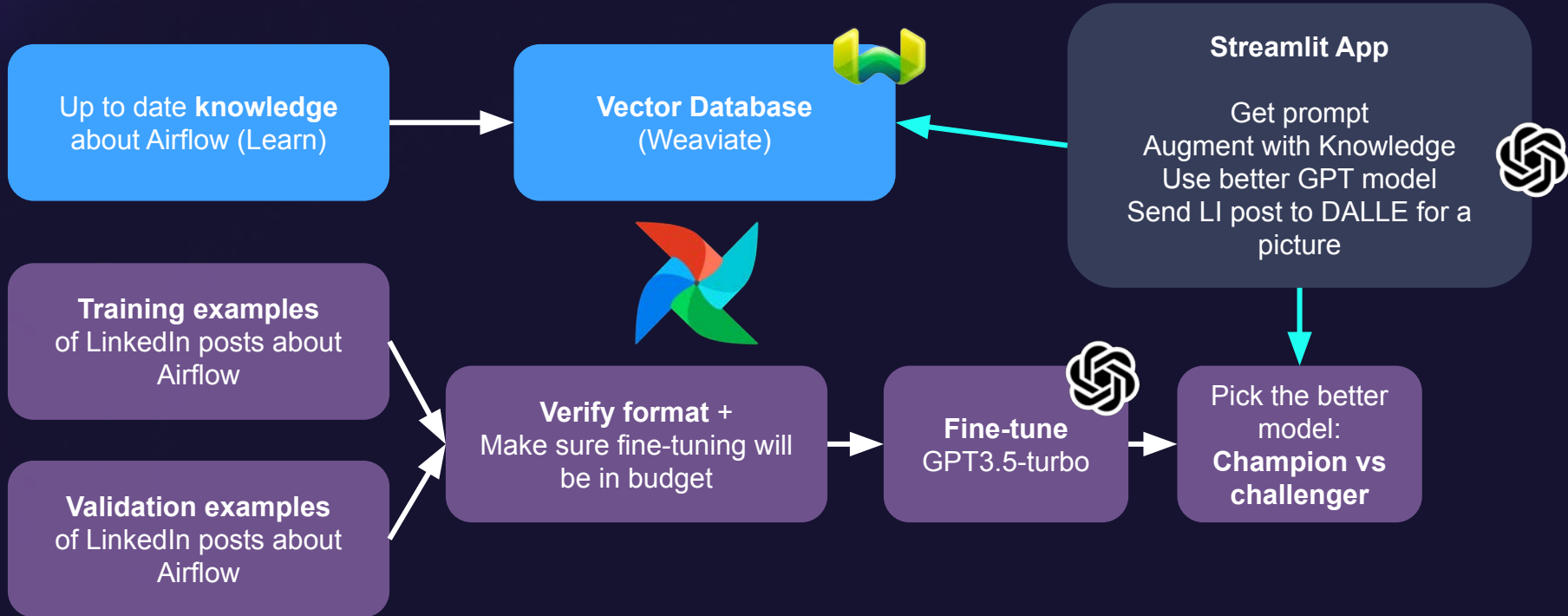
**Custom map index (2.9)**



Details: https://astronomer.io/docs/learn/dynamic-tasks

# Demo repository: Content Generation

https://github.com/astronomer/gen-ai-fine-tune-rag-use-case

# Demo

https://github.com/astronomer/gen-ai-fine-tune-rag-use-case

**Take Home Message:**
**Your data + your orchestration with Airflow is what sets you apart from competitors when creating GenAI applications!**

# Appendix

# Feature focus: Advanced Dataset scheduling (2.9)

Conditional Dataset Scheduling

```
5    @dag(
6        start_date=datetime(2024, 3, 1),
7        schedule=(
8            (Dataset("dataset1") | Dataset("dataset2"))
9            & (Dataset("dataset3") | Dataset("dataset4"))
10       ),  # Use () instead of [] to be able to use conditional dataset scheduling!
11       catchup=False
12   )
13   def downstream2_one_in_each_group():
```

(Dataset 1 **OR** Dataset 2) **AND** (Dataset 3 **OR** Dataset 4)

ASTRONOMER

# Feature focus: Advanced Dataset scheduling (2.9)

## Time + Dataset Scheduling

```python
10    from airflow.timetables.datasets import DatasetOrTimeSchedule
11    from airflow.timetables.trigger import CronTriggerTimetable
```

```python
55    @dag(
56        dag_display_name="📚 Ingest Knowledge Base",
57        start_date=datetime(2024, 4, 1),
58        schedule=DatasetOrTimeSchedule(
59            timetable=CronTriggerTimetable("0 0 * * *", timezone="UTC"),
60            datasets=reduce(
61                lambda x, y: Dataset(x) | Dataset(y), _KNOWLEDGE_BASE_DATASET_URIS
62            ),
63        ),
64        catchup=False,
```