Chương 8 Apache Airflow: Orchestrating Data Workflows

What is Apache Airflow?

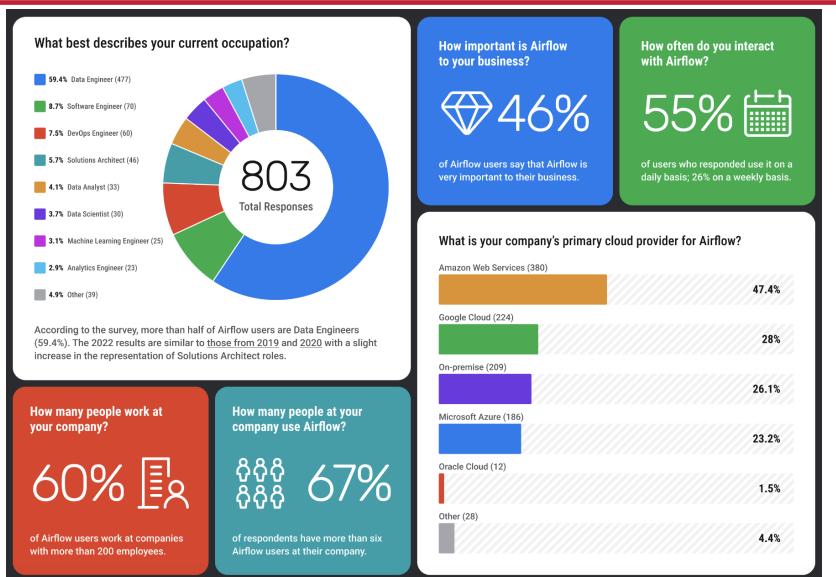
- An open-source platform for orchestrating complex computational workflows and data processing pipelines
- Pain points in data engineering
 - Dependency management between tasks
 - Scheduling and retries
 - Monitoring and error handling
 - Scalability of data pipelines

Airflow's Approach to Workflow Management

- Introduces the concept of Directed Acyclic Graphs (DAGs) for workflow representation
- Benefits
 - Clear visualization of task dependencies
 - Easier management of complex workflows
 - Improved reliability and error handling



https://airflow.apache.org/survey/





Average number of active DAGs in largest Airflow instance

Average maximum number of tasks in a single DAG

Average number of production Airflow instances

Average number of schedulers in largest Airflow instance

64%

have Airflow instances dedicated to Dev and Staging environments

65%

do not follow security announcements and advisories from apache.org

86%

have not made any customizations to **Airflow**

What use cases do you use Airflow for?

Ingestion and ETL/ELT related to analytics (771)

89.5%

Ingestion and ETL/ELT related to business operations (539)

67.3%

Training, serving, or generally managing MLOps (227)

28.3%

Spinning up and spinning down infrastructure (107) 13.4%

Other (25)

2.6%

How has your team's use cases with Airflow grown over time?

A lot. We initially started with just a few use cases, but we have many more now. (354)

46.5%

A little, but we're mostly still using it for the original reason we started using Airflow. (305)

40.1%

None, we're only using it for what we started doing with it. (74)

9.7%

It's decreasing as we're moving towards other orchestrators. (19)

2.5%

Other (9)

1%



Airflow's Workflow Management Processes

- DAG Definition
- Scheduling
- Execution
- Monitoring
- Error Handling

DAG Definition: Workflows as Python Code

- DAGs are defined using Python, allowing for dynamic and programmatic creation
- Each DAG is a separate Python file in the DAGs folder
- DAGs consist of a series of tasks and their dependencies

```
from airflow import DAG
from airflow.operators.python import PythonOperator
from datetime import datetime, timedelta
def task_1():
    print("Executing task 1")
def task_2():
    print("Executing task 2")
with DAG('example_dag',
         start_date=datetime(2024, 1, 1),
         schedule_interval=timedelta(days=1)) as dag:
    t1 = PythonOperator(task id='task 1', python callable=task 1)
    t2 = PythonOperator(task_id='task_2', python_callable=task_2)
    t1 >> t2 # Set task dependency
```

DAG Scheduling

- Cron-based Scheduling
 - Uses cron expressions for precise time-based scheduling
 - Example: 0 0 * * * (run daily at midnight)
- Interval-based Scheduling
 - Defines intervals between DAG runs
 - Example: timedelta(days=1) for daily runs
- External Triggers
 - Manual triggers via UI or API
 - Triggered by external systems or events

```
# Cron-based
dag = DAG('cron_dag', schedule_interval='0 0 * * *')
# Interval-based
dag = DAG('interval_dag', schedule_interval=timedelta(hours=6))
# Externally triggered
dag = DAG('on_demand_dag', schedule_interval=None)
```



Task Execution Based on Dependencies

Execution Process

- DAG parsing and task instantiation
- Dependency resolution
- Task queuing
- Worker assignment and execution

Types of Dependencies

- Temporal: Based on execution time or data intervals
- Logical: Based on task relationships within the DAG

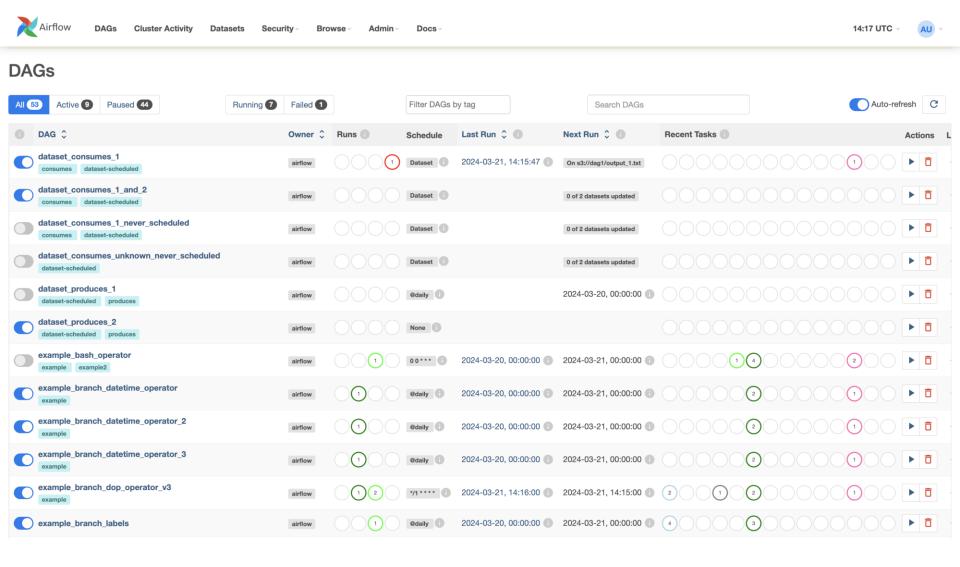
Execution Modes

- Sequential: One task at a time
- Parallel: Multiple tasks simultaneously (based on dependencies)

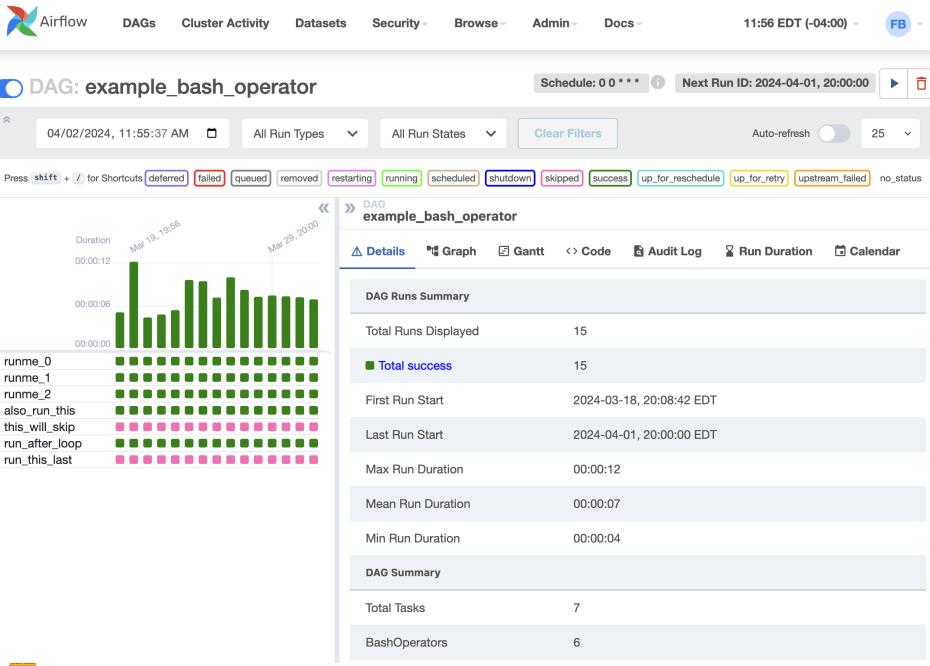
Example

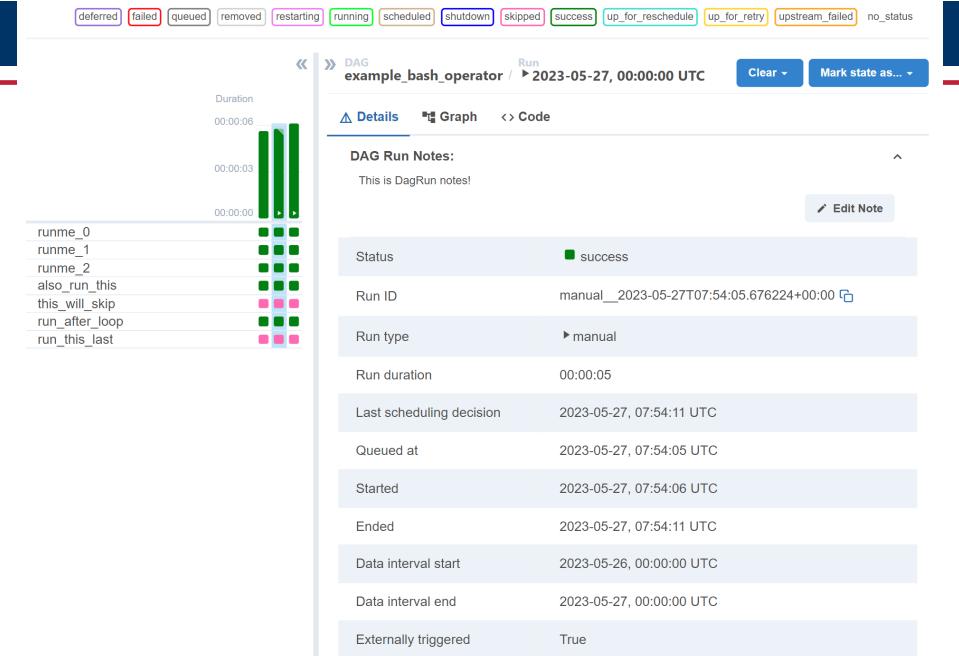
```
# Option 1: Using bitshift operators
task_1 >> task_2 >> task_3
# Option 2: Using set_upstream/set_downstream
task_3.set_upstream(task_2)
task_2.set_upstream(task_1)
# Option 3: Using list notation [task_1, task_2] >> task_3
```

Real-time Monitoring through Web UI and Logging



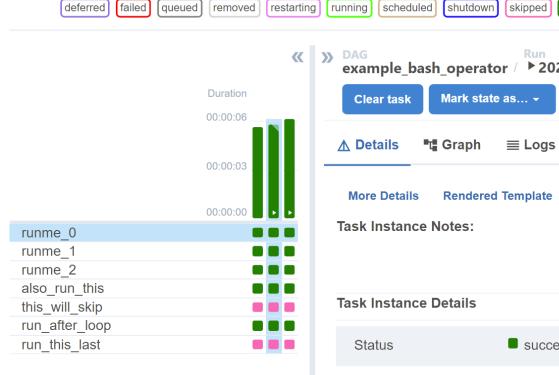


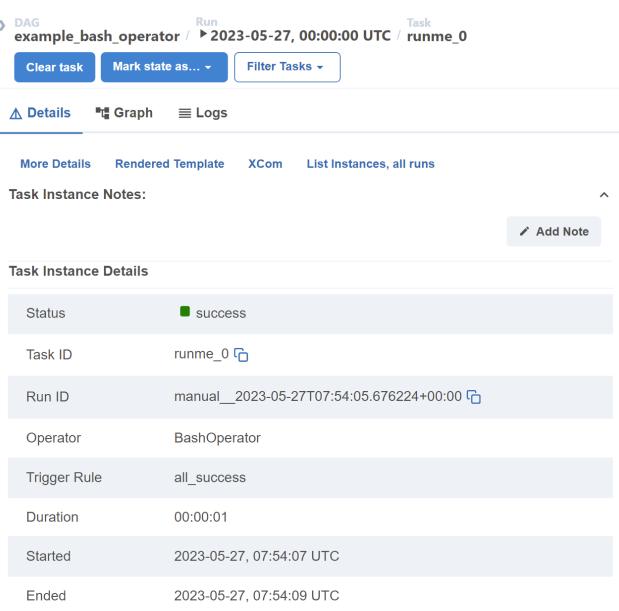




None

Run config





success

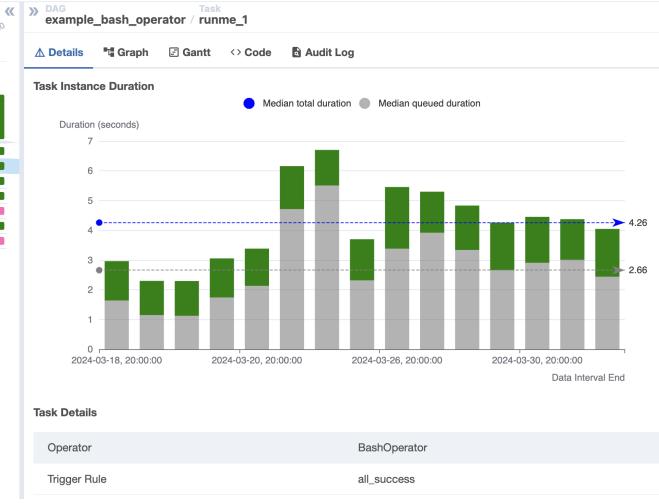
up for reschedule

up_for_retry

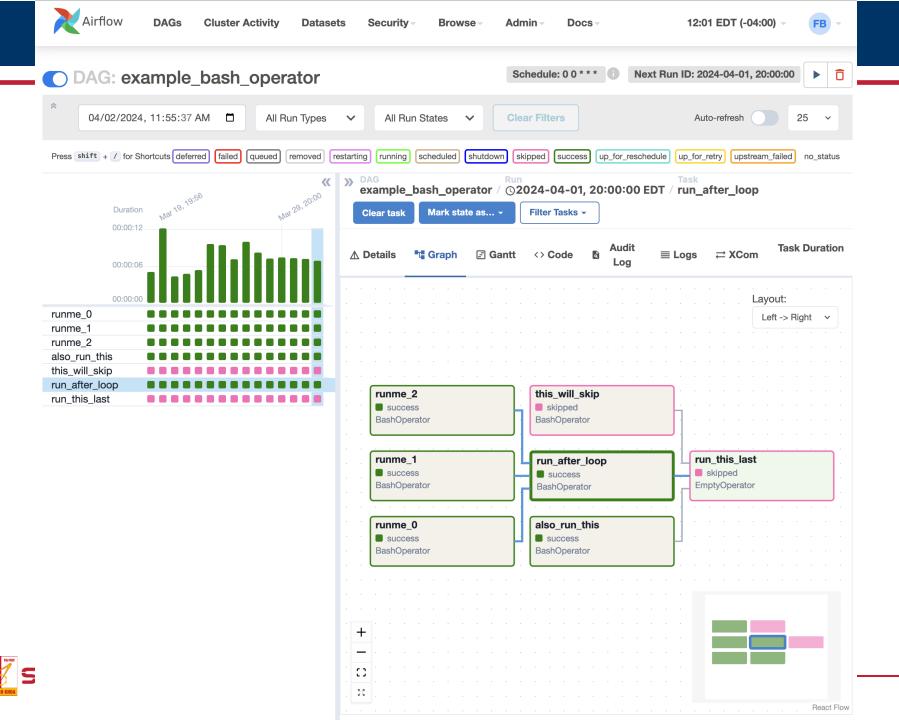
upstream failed

no status

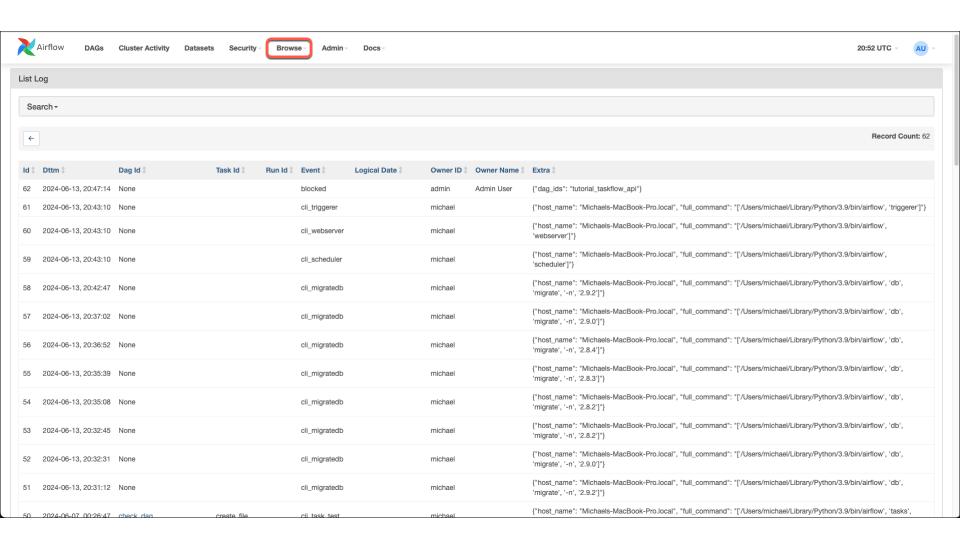








Airflow logging





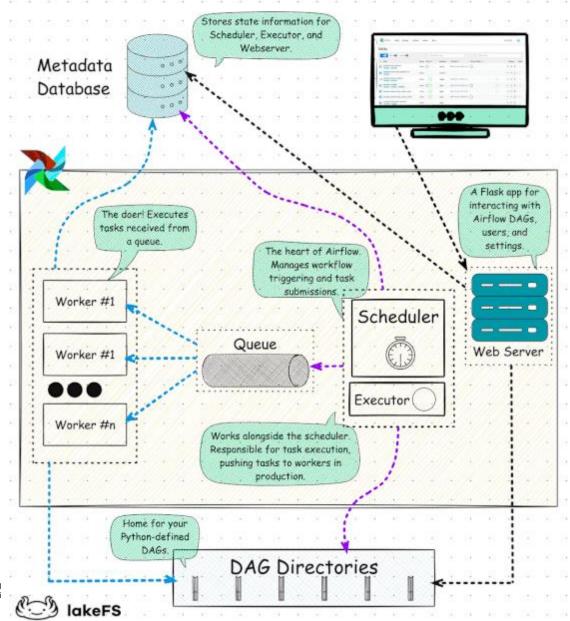
Airflow DAG notifications

```
@task(
          email=["noreply@astronomer.io", "noreply2@astronomer.io"],
          email on failure=True,
          email on retry=True
def t1():
          return "hello"
from airflow.utils.email import send email
def success email function(context):
          dag run = context.get("dag run")
          msg = "DAG ran successfully"
          subject = f"DAG {dag run} has completed"
          send email(to=your_emails, subject=subject, html_content=msg)
@dag(
          start date=datetime(2023, 4, 26),
          schedule="@daily",
          catchup=False,
          on success callback=success email function
```

Robust Error Handling and Retry Mechanisms

- Automatic Retries
 - Configurable retry attempts and intervals
 - Exponential backoff support
- Failure Notifications
 - Email alerts on task failures
 - Customizable notification rules
- Conditional Execution
 - Trigger specific tasks on failure of others

Airflow- Architecture





Airflow Executors

SequentialExecutor

- Default executor
- Runs one task at a time
- Uses SQLite as the Metadata DB
- Pros: Simple, easy to set up
- Cons: Not suitable for production or parallel execution
- Use case: Development and testing

LocalExecutor

- Runs tasks on the same machine as the Scheduler
- Can run tasks in parallel using multiple processes
- Pros: Easy to set up, good for medium workloads
- Cons: Limited by single machine resources
- Use case: Small to medium-scale deployments



Airflow Executors (2)

CeleryExecutor

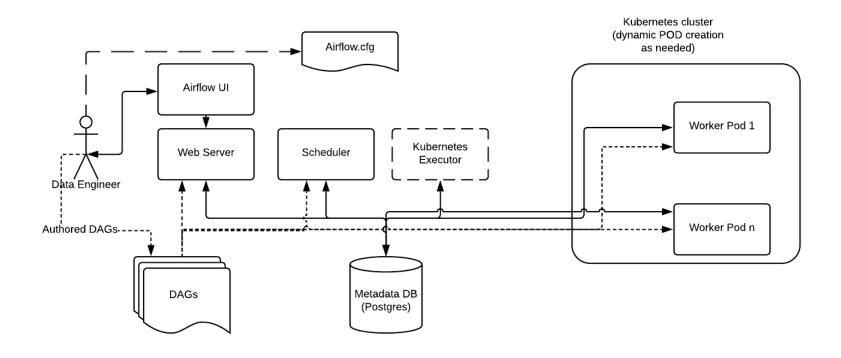
- Uses Celery to run distributed tasks
- Allows for horizontal scaling across multiple worker machines
- Pros: Highly scalable, good for large workloads
- Cons: More complex setup, requires additional components (Redis/RabbitMQ)
- Use case: Large-scale production deployments

KubernetesExecutor

- Runs each task in a separate Kubernetes pod
- Allows for dynamic allocation of resources
- Pros: Highly scalable, efficient resource utilization
- Cons: Requires Kubernetes cluster, complex setup
- Use case: Cloud-native environments, need for isolation between tasks



KubernetesExecutor





Thank you for your attention! Q&A