

Streamlining Enterprise Data Pipelines with an Automated DAG Factory for Airflow Orchestration in Cloud Environments using YAML Templates and JSON - Serialized Variables

Ramamurthy Valavandan¹, Balakrishnan Gothandapani², Savitha Ramamurthy³

^{1,2,3}Nature Labs - Research Centre, Namakkal, Tamil Nadu, India

¹Corresponding author Email: [find\[at\]ngosys.com](mailto:find[at]ngosys.com)

Abstract: *Airflow is an open - source platform for creating, scheduling, and monitoring data pipelines. Its Directed Acyclic Graph (DAG) factory provides a mechanism for creating and managing DAGs in a programmatic way. However, the current implementation of the DAG factory in Airflow requires writing Python code, which can be time - consuming and error - prone. In this research paper, we propose a YAML - based DAG factory automation framework for Airflow, which provides a simple and intuitive way to define DAGs in YAML format. We describe the design and implementation of the framework and provide examples of how it can be used to automate the creation and management of DAGs in a cloud environment. We also evaluate the performance and scalability of the framework using real - world datasets and compare it to the existing Python - based DAG factory in Airflow. Our results demonstrate that the YAML - based DAG factory automation framework provides a more efficient and flexible way to create and manage DAGs in Airflow, especially in large - scale data processing scenarios.*

Keywords: Airflow, Directed Acyclic Graph, DAG factory, YAML, automation, Python, CLI tool, schema file, GCP, Composer, JSON, dictionary, task status, DAG tasks, template generation, variable

1. Introduction

This research paper investigates the development and implementation of a Python script that automates the creation of Airflow DAG files using a YAML - based DAG factory and JSON - serialized variables. This tool streamlines the process of generating DAG files, enhancing the efficiency and effectiveness of data pipeline orchestration in cloud environments. [1]

The innovative DAG factory automation framework presented in this paper has a wide range of potential use cases for enterprises that rely on Airflow [2] for data pipeline orchestration. One use case is the quick generation of DAG files for complex data pipelines involving multiple data sources and destinations, reducing the time and effort required to create these pipelines. [2]

Additionally, this framework can help standardize the process of creating and managing DAG files across multiple teams and projects, reducing the risk of errors and inconsistencies. This standardization also facilitates the sharing of DAG files between teams and projects, enabling better collaboration and faster development times. [3]

Another use case for this DAG factory automation framework is the streamlined deployment of new data pipelines. With its ability to quickly generate DAG files using YAML templates and JSON - serialized variables, [4] organizations can adapt to changing business requirements and data sources more quickly, ultimately driving better business outcomes. [5]

The use of YAML templates and JSON - serialized variables

also makes it easy to make changes to DAG files and propagate those changes across multiple pipelines. This feature can be particularly beneficial in large, complex data pipeline environments where changes are common and time is of the essence. [6]

Finally, the DAG factory automation framework can be used to develop custom data connectors for Airflow, allowing organizations to easily integrate with new data sources and destinations as they become available. [7] This capability can help enterprises stay ahead of the curve in terms of data processing and analysis, ultimately leading to better business outcomes. [8]

In summary, the key benefits of this automated DAG factory include increased efficiency, reduced errors, and faster deployment times for complex data pipelines in cloud environments. These benefits have numerous potential applications and use cases for enterprises seeking to streamline their data pipeline orchestration processes and achieve their data processing and analysis objectives more effectively and efficiently.

2. Approach in data pipeline

Streamlining Enterprise Data Pipelines with an Automated DAG Factory for Airflow Orchestration in Cloud Environments using YAML Templates and JSON - Serialized Variables". This paper explores an innovative approach to enhancing data pipeline orchestration with Airflow by automating the creation of DAG files. Using a Python script that leverages YAML templates and JSON - serialized variables, this automated DAG factory offers a streamlined and efficient process for generating DAG files

with minimal human intervention. [9]

This approach has a wide range of potential applications in enterprise data pipeline management. By reducing the risk of errors and inconsistencies, standardizing DAG file creation and management, and accelerating the deployment of new data pipelines, organizations can achieve their data processing and analysis objectives more effectively and efficiently. Additionally, the use of YAML templates and JSON - serialized variables makes it easy to make changes to DAG files and propagate those changes across multiple pipelines.

In summary, this innovative approach to Airflow orchestration offers significant benefits for organizations seeking to streamline their data pipeline management processes in cloud environments

3. Data Pipeline Automation Approach

This research paper investigates a novel approach to streamlining Airflow orchestration in cloud environments for enterprise data pipelines. [10]

By automating the generation of YAML template DAG factories and JSON - serialized variables for Airflow, this approach enhances the efficiency and effectiveness of data pipeline management. We delve into the key features and benefits of this automated DAG factory and explore its potential use cases for enterprises. Additionally, we provide examples of how this approach can be used to accelerate the development and deployment of complex data pipelines, ultimately helping organizations to achieve their data processing and analysis objectives more efficiently and effectively.

3.1 Automation of DAG factory

Automating Airflow DAG Factory Generation with YAML Templates and JSON - Serialization for Streamlined Data Pipeline Orchestration in Cloud Environments.

Enhancing Enterprise Data Pipelines with an Automated Airflow DAG Factory using YAML Templates and JSON - Serialization in Cloud Environments [11]

Simplifying Airflow Orchestration in Cloud Environments with an Automated DAG Factory using YAML Templates and JSON - Serialization for Enterprise Data Pipelines

Boosting Data Pipeline Efficiency with an Automated Airflow DAG Factory using YAML Templates and JSON - Serialization in Cloud Environments for Enterprise Data Orchestration

A Comprehensive Solution for Streamlined Airflow Orchestration in Cloud Environments with Automated DAG Factory Generation using YAML Templates and JSON - Serialization for Enterprise Data Pipelines.

3.2 Airflow DAG Generation in Cloud Environments.

An approach on streamlining the Airflow orchestration in

cloud environments in enterprise data pipelines by automating the generation of YAML template DAG factory and JSON - serialized variables for Airflow:

Assess the current state of Airflow DAG file creation and management within the organization.

Evaluate the benefits and potential drawbacks of automating DAG file creation using YAML templates and JSON - serialized variables. [12]

Determine the specific use cases and requirements for the automated DAG factory, such as supporting multiple data sources and destinations, standardizing DAG file creation across teams and projects, and streamlining deployment.

Develop a Python script that leverages YAML templates and JSON - serialized variables to automate the creation of Airflow DAG files.

Test and refine the automated DAG factory script to ensure that it is efficient, accurate, and reliable.

Train relevant stakeholders and teams on how to use the automated DAG factory and integrate it into their data pipeline orchestration workflows.

Monitor and analyze the performance and impact of the automated DAG factory, gathering feedback and making improvements as needed.

Continuously iterate and improve the automated DAG factory to meet evolving business needs and data pipeline requirements.

4. An overview of Airflow and its role in data pipeline orchestration

4.1 Introduction to Data Pipeline Orchestration

Data pipeline orchestration is the process of integrating multiple data sources and transforming them into a unified format for analysis and reporting. It involves designing, deploying, and managing a series of interconnected data processing tasks that run in a specific sequence to achieve a desired outcome. Data pipeline orchestration is a critical component of data management, especially for organizations dealing with large volumes of data. [13]

One popular tool for data pipeline orchestration is Apache Airflow. Airflow is an open - source platform that enables users to define, schedule, and monitor workflows. Airflow uses directed acyclic graphs (DAGs) to define and manage workflows. DAGs are a collection of tasks that are executed in a specific order. Airflow provides a web - based user interface that enables users to monitor the status of their workflows and diagnose issues. [14]

Airflow has become an essential tool for organizations looking to streamline their data processing workflows. Its flexibility and extensibility have made it a popular choice for data engineers, data analysts, and data scientists. With Airflow, organizations can create complex data pipelines

that integrate multiple data sources and destinations, transforming data into actionable insights.

In the following sections, we will explore the key features and benefits of Airflow in data pipeline orchestration, as well as its limitations and challenges. We will also discuss how the use of YAML templates and JSON - serialized variables can enhance the efficiency and effectiveness of Airflow DAG generation and data pipeline orchestration in cloud environments.

4.2 An Introduction to Apache Airflow

In this subsection, the focus can be on introducing Apache Airflow and its features. It can cover the history of Airflow, its architecture, and how it works.

Apache Airflow is an open - source platform used for programmatically authoring, scheduling, and monitoring workflows or data pipeline. It allows users to create directed acyclic graphs (DAGs) of tasks, which can be orchestrated and executed in parallel. With Airflow, users can define, schedule, and monitor complex workflows with ease, making it a popular choice for data pipeline orchestration. [15]

It was created by Maxime Beauchemin in 2014 as an alternative to the existing workflow management tools that were available at the time.

Airflow allows users to define their workflows as code using Python, making it easy to version control and maintain them. It also provides a rich set of operators, such as BashOperator and PythonOperator, which can be used to build complex data pipelines. [16]

One of the key features of Airflow is its use of Directed Acyclic Graphs (DAGs) to represent workflows. DAGs allow users to define the dependencies between tasks in their workflows and specify the order in which they should be executed. This makes it easy to visualize the entire workflow and understand the dependencies between tasks. [17]

Another important feature of Airflow is its ability to integrate with a variety of technologies and services. Airflow comes with a wide range of built - in connectors for popular services such as Amazon S3, Google Cloud Storage, and Hadoop Distributed File System (HDFS), as well as a Python API that makes it easy to create custom connectors. [18]

Overall, Airflow provides a powerful and flexible platform for data pipeline orchestration that can be used in a wide range of use cases, from simple ETL pipelines to complex machine learning workflows.

4.3 DAG Factory for Airflow

DAG Factory is a Python script that automates the creation of Airflow DAG files. By utilizing YAML templates and JSON - serialized variables, this tool streamlines the process of generating DAG files, enhancing the efficiency and effectiveness of data pipeline orchestration in cloud

environments. [19] DAG Factory can quickly generate DAG files for complex data pipelines involving multiple data sources and destinations, standardizing the process of creating and managing DAG files across multiple teams and projects, reducing the risk of errors and inconsistencies. [20]

4.4 Requirements and Challenges in YAML Generation

To generate YAML templates for Airflow DAGs, there are certain requirements and challenges that need to be considered. The YAML templates must be designed to allow for flexibility and customization while maintaining consistency across the DAGs. Variables need to be defined and organized in a way that allows for easy management and reuse across multiple DAGs. Furthermore, proper validation and error handling must be incorporated into the generation process to ensure that the YAML templates are correct and meet the required specifications. [21 - 25]

Overall, the use of DAG Factory with YAML templates and JSON - serialized variables provides a powerful and flexible approach to Airflow DAG automation in the cloud, streamlining the data pipeline orchestration process and allowing organizations to achieve their data processing and analysis objectives more effectively and efficiently.

5. Airflow Components and use cases

This subsection can delve into the various components that make up an Airflow instance, such as the scheduler, web server, and workers. It can also cover how these components work together to orchestrate data pipelines.

Airflow is a comprehensive platform that includes various components, each playing an important role in data pipeline orchestration. Some of the key components of Airflow include:

DAGs (Directed Acyclic Graphs): DAGs are the fundamental building blocks of Airflow. They represent a series of tasks that need to be executed in a specific order. DAGs can be defined in Python files, and Airflow uses these files to create and manage the pipeline. [26]

Operators: Operators are the individual units of work within a DAG. Each operator performs a specific task, such as extracting data from a source or transforming data in a particular way. [27]

Sensors: Sensors are similar to operators, but instead of performing an action, they wait for a specific event or condition to occur before proceeding with the next task.

Hooks: Hooks are a way for Airflow to interact with external systems, such as databases or APIs. Each hook provides a connection to a specific system and a set of methods for interacting with that system. [28]

Executors: Executors determine how tasks are executed within Airflow. There are several types of executors available, including Local Executor, Sequential Executor, and Celery Executor. [29]

5.1 Some of the common use cases of Airflow

ETL (Extract, Transform, Load) pipelines: Airflow is often used to orchestrate ETL pipelines, which involve extracting data from various sources, transforming it into a desired format, and loading it into a destination system. [30 - 31]

Data warehousing: Airflow can be used to automate the process of building and updating data warehouses, which involve aggregating and storing data from various sources for analysis and reporting. [32]

Machine learning workflows: Airflow can be used to orchestrate complex machine learning workflows, which involve training models on large datasets and deploying them in production environments. [33]

Data streaming: Airflow can also be used to manage real - time data streaming pipelines, which involve processing and analyzing data as it flows in from various sources. [33 - 35]

Overall, Airflow provides a flexible and powerful platform for data pipeline orchestration, making it an ideal choice for enterprises that need to manage complex data workflows in the cloud. However, generating YAML templates for Airflow can be challenging, which is why automated DAG factories can be a valuable tool for streamlining the process.

5.2 DAGs and Operators

DAGs and Operators are the core concepts in Apache Airflow that enable the creation and management of data pipeline workflows. A Directed Acyclic Graph (DAG) is a collection of tasks that are linked together in a specific order, forming a workflow. DAGs are used to define the dependencies and relationships between tasks in a workflow, making it possible to orchestrate complex data processing and analysis tasks.

Operators in Airflow are the building blocks for tasks within a DAG. Operators define what actions need to be taken at each step of the workflow. Each operator represents a specific task or action, such as running a SQL query, moving data between systems, or sending an email notification. Airflow provides a rich set of built - in operators for common tasks, such as BashOperator, PythonOperator, DataProcOperator, and EmailOperator.

To illustrate the use of DAGs and Operators in Airflow, consider an example where we need to run a data processing workflow on a daily basis. The DAG for this workflow, with ID "Workflow_Orchestration, " could be defined as follows:

```
from airflow import DAG
from airflow.operators.bash_operator import BashOperator
from airflow.operators.python_operator import PythonOperator
from airflow.operators.email_operator import EmailOperator
from datetime import datetime, timedelta

default_args = {
    'owner': 'data - engineering',
    'depends_on_past': False,
    'start_date': datetime(2022, 1, 1),
```

```
'email_on_failure': False,
'email_on_retry': False,
'retries': 1,
'retry_delay': timedelta(minutes=5),
}

dag = DAG (
    dag_id='Workflow_Orchestration',
    default_args=default_args,
    schedule_interval=timedelta(days=1),
)

t1 = BashOperator (
    task_id='download_data',
    bash_command='python Python_Application.py
download_data',
    dag=dag,
)

t2 = BashOperator (
    task_id='transform_data',
    bash_command='python Python_Application.py
transform_data',
    dag=dag,
)

t3 = PythonOperator (
    task_id='validate_data',
    python_callable=validate_data,
    dag=dag,
)

t4 = DataProcOperator (
    task_id='run_analysis',
    dataproc_cluster='my - dataproc - cluster',
    main_jar='analysis.jar',
    arguments= ['input_data', 'output_data'],
    dag=dag,
)

t5 = EmailOperator (
    task_id='send_email',
    to='john.doe[at]example.com',
    subject='Workflow_Orchestration Succeeded',
    html_content='The data processing workflow has completed
successfully.',
    dag=dag,
)

t1 >> t2 >> t3 >> t4 >> t5
```

In this example, the workflow consists of five tasks, each represented by an operator. The BashOperator is used to download and transform data, while the PythonOperator is used to validate the data. The DataProcOperator is used to run an analysis job on a Dataproc cluster, and the EmailOperator is used to send a notification email when the workflow completes.

The DAG is scheduled to run once per day, and each task is defined with a unique task ID and a command or callable that specifies what action needs to be taken. The dependencies between tasks are defined using the >> operator, which indicates the order in which tasks should be executed.

Overall, DAGs and Operators provide a powerful way to

orchestrate complex data processing workflows in Airflow, making it easier to manage and automate data pipeline tasks.

5.3 Airflow UI and CLI

In this subsection, the focus can be on how to interact with Airflow through its user interface (UI) and command - line interface (CLI). It can cover the various features of the Airflow UI, such as monitoring and managing DAGs, as well as the different commands available through the CLI. [36]

The Airflow user interface (UI) provides a web - based interface to monitor, manage and visualize DAGs, tasks, and logs. It allows users to perform various tasks such as creating and editing DAGs, setting task dependencies, and monitoring task statuses. The Airflow UI can be accessed through a web browser and provides a graphical representation of DAGs with the ability to view task instances, logs, and statistics. [37]

Airflow also provides a command - line interface (CLI) that enables users to interact with Airflow from the terminal. The CLI offers a range of commands to perform tasks such as starting and stopping the Airflow scheduler, listing available DAGs, and triggering specific tasks. The CLI also provides the ability to create, delete, and update DAGs and tasks. [38]

Using the Airflow UI and CLI, users can perform tasks such as:

Monitoring and managing DAGs: The Airflow UI provides a dashboard that displays the status of all DAGs, tasks, and runs. Users can view the DAGs' graphical representation, track task dependencies, and monitor the execution status of each task. Additionally, users can start, stop, pause, and unpause DAGs.

Setting up alerts and notifications: Airflow provides an EmailOperator that sends email notifications when certain conditions are met, such as a task failing or a DAG taking longer than expected to complete. These notifications can help users identify issues and take corrective action promptly.

Debugging and troubleshooting: Airflow stores task logs and metadata in a database, which can be accessed through the UI. This feature enables users to view the task's output, errors, and status, aiding in debugging and troubleshooting. [39]

Automating data pipelines: With the Airflow CLI, users can automate data pipelines by scheduling DAGs and tasks to run at specified intervals. Additionally, users can create DAGs programmatically, which enables them to integrate Airflow with other tools and platforms. [40]

Overall, the Airflow UI and CLI provide users with powerful tools to monitor, manage, and automate data pipelines.

Airflow Architecture for Cloud Environments

This subsection can cover how Airflow architecture can be adapted for cloud environments, such as AWS, GCP, or

Azure. It can delve into how Airflow can be deployed on these cloud platforms and how it can be integrated with other cloud services. [41]

Cloud Deployment Options:

Airflow can be deployed on various cloud platforms, including Amazon Web Services (AWS), Google Cloud Platform (GCP), and Microsoft Azure. Each cloud platform offers different deployment options for Airflow, such as using managed services like AWS Elastic Kubernetes Service (EKS), GCP Kubernetes Engine, or Azure Kubernetes Service (AKS), or deploying Airflow on virtual machines using infrastructure as code tools like Terraform or CloudFormation. The choice of deployment option depends on factors like scalability, ease of maintenance, and cost. [42]

Integration with Cloud Services:

Airflow can be integrated with various cloud services to build data pipelines, including data storage services like AWS S3, GCP Cloud Storage, or Azure Blob Storage, data processing services like AWS Lambda, GCP Cloud Functions, or Azure Functions, and data warehousing services like AWS Redshift, GCP BigQuery, or Azure Synapse Analytics. Airflow provides operators for each of these cloud services, making it easier to build pipelines that utilize these services. [43]

Scaling Airflow in the Cloud:

One of the advantages of using Airflow in the cloud is the ability to scale horizontally by adding or removing worker nodes based on demand. Cloud platforms like AWS and GCP provide auto - scaling capabilities that can be used to automatically add or remove worker nodes based on the number of pending tasks in the Airflow queue. Additionally, Airflow can be run on Kubernetes clusters, which provide built - in scaling capabilities. [44]

Best Practices for Cloud - based Airflow:

When deploying Airflow in the cloud, there are some best practices that can help ensure optimal performance and minimize downtime. These include:

Designing Airflow DAGs to be idempotent and resilient to failure

Monitoring and alerting for Airflow components like the scheduler, web server, and worker nodes

Implementing a disaster recovery plan for Airflow components and data

Utilizing cloud - native services like AWS RDS or GCP Cloud SQL for the Airflow metadata database

Using security best practices like encrypting sensitive data and restricting access to Airflow components. [45]

5.4 Best Practices for Airflow DAG Orchestration

In this subsection, the focus can be on some best practices for DAG orchestration in Airflow, such as using the right operators for specific tasks, creating modular DAGs, and managing dependencies between tasks. [46]

Airflow DAGs (Directed Acyclic Graphs) are a powerful tool for data pipeline orchestration, and their proper use can

greatly enhance data processing efficiency. However, DAG orchestration can be complex, and it is important to follow best practices to ensure effective workflow management. Here are some best practices to consider when working with Airflow DAGs:

Use the Right Operators: Airflow comes with a wide range of operators, each of which is designed for a specific purpose. It is essential to choose the right operator for the task at hand to ensure that your DAG runs smoothly. For example, the `BashOperator` can be used to run a Bash command, while the `PythonOperator` can be used to run a Python script.

Create Modular DAGs: Modular DAGs are easier to manage and maintain than monolithic ones. Consider breaking your DAG into smaller, more manageable components, each responsible for a specific task. This approach allows for easier testing, debugging, and modification of individual components, without affecting the entire DAG. [47]

Use Sensible Default Arguments: The default arguments of a DAG define its behavior, including its start date, end date, and schedule interval. It is important to define these arguments in a way that makes sense for your workflow. A sensible schedule interval ensures that your DAG runs at the right frequency, while a suitable start and end date ensures that your DAG runs within a specific time frame.

Manage Dependencies: Airflow tasks can have dependencies on other tasks, and it is important to manage these dependencies carefully to ensure that your DAG runs efficiently. Consider using the `trigger_rule` parameter to define how tasks should behave if their dependencies fail or succeed.

Monitor and Debug Your DAGs: Proper monitoring and debugging can help you identify and resolve issues with your DAGs quickly. Use Airflow's built-in monitoring tools to keep track of task progress, and use logging statements to output information about each task's execution.

By following these best practices, you can ensure that your Airflow DAGs are well-structured, efficient, and easy to manage.

5.5 Airflow and Other Orchestration Tools

Here, the focus can be on how Airflow compares to other data pipeline orchestration tools in the market, such as Luigi, Oozie, and Azkaban. It can cover the strengths and weaknesses of each tool and how they differ in terms of features, ease of use, and scalability.

Airflow is not the only data pipeline orchestration tool available in the market. Other tools such as Luigi, Oozie, and Azkaban are also commonly used by organizations. Each tool has its own strengths and weaknesses, and it's important to evaluate them carefully to choose the best tool for your specific use case.

Luigi, developed by Spotify, is another open-source workflow management tool. Like Airflow, it is written in

Python and provides a simple interface to define workflows. Luigi is often used for batch processing and data pipeline automation. [48]

Oozie, developed by Apache, is another workflow scheduler for Hadoop that enables users to create directed acyclic graphs of tasks. It supports multiple Hadoop jobs, such as Java MapReduce, Pig, Hive, and Sqoop.

Azkaban, developed by LinkedIn, is another open-source workflow management tool that is often used for Hadoop job orchestration. It allows users to define and schedule workflows using a web interface, and it can be integrated with other Hadoop tools such as Pig, Hive, and Hadoop MapReduce.

When comparing these tools with Airflow, some common factors to consider include ease of use, scalability, extensibility, and community support. While all of these tools have their strengths, Airflow is often praised for its flexibility, extensibility, and active community support. Additionally, Airflow's use of DAGs provides a simple, intuitive way to define and manage complex workflows.

5.6 Use Cases of Airflow in Industry

In this subsection, the focus can be on some real-world use cases of Airflow in different industries, such as finance, healthcare, and e-commerce. It can cover how Airflow has helped these industries to streamline their data pipeline workflows and achieve better business outcomes.

Here are some examples of how Airflow is used in different industries:

Finance: In the finance industry, Airflow is used for tasks such as data ingestion, ETL processing, and reporting. Airflow can help financial institutions to consolidate data from multiple sources, perform calculations on large datasets, and generate timely reports for compliance and regulatory purposes.

Healthcare: In the healthcare industry, Airflow is used for tasks such as data integration, patient monitoring, and clinical decision-making. Airflow can help healthcare providers to manage patient data from multiple sources, track patient outcomes, and identify patterns and trends that can inform clinical decisions.

E-commerce: In the e-commerce industry, Airflow is used for tasks such as order processing, inventory management, and marketing campaign management. Airflow can help e-commerce companies to automate their order fulfillment processes, keep track of inventory levels, and run targeted marketing campaigns based on customer behavior and preferences.

Media and Entertainment: In the media and entertainment industry, Airflow is used for tasks such as content management, video processing, and analytics. Airflow can help media companies to manage their content libraries, process large video files, and analyze viewer behavior to improve their content offerings.

Telecommunications: In the telecommunications industry, Airflow is used for tasks such as network management, billing, and customer analytics. Airflow can help telecommunications companies to monitor network performance, generate accurate billing reports, and analyze customer behavior to improve their service offerings. [49]

6. Dynamic Dags and automation need

6.1 Dynamic Dags

Dynamic DAGs are a feature of Airflow that allows you to dynamically generate DAGs at runtime. This is useful when you have a large number of tasks that follow a similar pattern, but differ in some parameters, such as input data or parameters for a machine learning model.

In Airflow, you can define a Python function that generates a DAG object. This function can accept arguments, such as the name of the DAG, the start date, or the schedule interval. Inside the function, you can use Python control flow statements, such as loops or conditionals, to generate the tasks and dependencies of the DAG dynamically.

For example, let's say you have a set of data files, each containing a table with the same schema, and you want to run a set of SQL queries on each table. You can define a Python function that reads the list of files, generates a DAG object, and creates a task for each file:

```
from datetime import datetime, timedelta
from airflow import DAG
from airflow.operators.bash_operator import BashOperator

def create_dag (dag_id, start_date, schedule_interval,
data_path):
    dag = DAG (dag_id, start_date=start_date,
schedule_interval=schedule_interval)

    for filename in os.listdir (data_path):
        task_id = f"process_{filename}"
        task = BashOperator (
            task_id=task_id,
            bash_command=f"process_data. sh {os.path.join (data_path,
filename) }",
            dag=dag,
        )

        if filename != os.listdir (data_path) [0]:
            # The first file does not have dependencies
            task.set_upstream (prev_task)

        prev_task = task

    return dag

dag_id = "process_data"
start_date = datetime (2023, 5, 1)
schedule_interval = timedelta (days=1)
data_path = "/path/to/data"

dag = create_dag (dag_id, start_date, schedule_interval,
data_path)
```

In this example, the `create_dag` function reads the list of files in the `data_path` directory, creates a task for each file using the `BashOperator` operator, and sets the dependencies between the tasks based on the order of the files. The resulting DAG is returned by the function and can be scheduled by Airflow.

Note that the `create_dag` function can be customized to generate DAGs for different sets of data or different processing pipelines by changing the arguments or the control flow statements inside the function. This makes dynamic DAGs a powerful tool for automating complex data processing workflows.

6.2 Airflow Variable JSON serialization

To use JSON variables in Airflow, you can create a JSON file containing the variable values and load it into your workflow using the `Variable` class. For example, you could create a `config.json` file containing a set of database connection parameters:

```
{
  "database": {
    "host": "localhost",
    "port": 5432,
    "username": "my_user",
    "password": "my_password",
    "database_name": "my_database"
  }
}
```

And then load the values into your DAG using the `Variable` class:

```
from airflow.models import Variable

database_config = Variable.get ('config',
deserialize_json=True) ['database']

# use database_config values in your DAG tasks
# . . .
```

In this example, we're using the `Variable.get` method to load the config variable and deserialize it as a JSON object. We're then extracting the database object from the JSON and using its values in our DAG tasks.

By combining dynamic DAGs and JSON variables in Airflow, you can create workflows that are more flexible, scalable, and maintainable, and you can also automate the process of generating and executing these workflows.

7. Dynamic Dags and tools

Dynamic DAGs are Airflow workflows that are created dynamically at runtime, based on some external event or input. This allows for a more flexible and adaptive approach to workflow management, as the DAG can be customized to respond to changing data sources, business requirements, or other factors.

There are several tools and techniques that can be used to create dynamic DAGs in Airflow, including:

Jinja Templating: Jinja is a powerful templating engine that allows for dynamic generation of code. Jinja templates can be used to generate DAGs based on variables or data inputs, allowing for dynamic workflows that can adapt to changing conditions.

External Triggers: External triggers can be used to initiate DAGs based on events or inputs from external sources. For example, a new file arriving in a specific folder could trigger a DAG that processes that file.

Parameterized DAGs: Parameterized DAGs allow for the dynamic configuration of DAGs at runtime, based on user input or other factors. This can be useful for workflows that require customization or adaptation based on changing conditions.

Airflow Variables: Airflow provides a built-in feature called Variables, which can be used to store and retrieve dynamic values and configuration settings. Variables can be used to create dynamic DAGs that respond to changing conditions or inputs.

Overall, the use of dynamic DAGs and related tools can greatly enhance the flexibility and adaptability of Airflow workflows, allowing for more efficient and effective data pipeline orchestration.

7.1 DAG Factory

DAG Factory is a Python script that automates the creation of Airflow DAG files using YAML templates and JSON - serialized variables. It is an innovative approach to streamlining the process of generating DAG files, enhancing the efficiency and effectiveness of data pipeline orchestration in cloud environments. [50]

The DAG Factory automation framework has a wide range of potential use cases for enterprises that rely on Airflow for data pipeline orchestration. One of the key benefits of this automated DAG factory is that it can help to quickly generate DAG files for complex data pipelines involving multiple data sources and destinations. Moreover, it can standardize the process of creating and managing DAG files across multiple teams and projects, reducing the risk of errors and inconsistencies.

Another benefit of the DAG Factory automation framework is that it can help to streamline the deployment of new data pipelines, allowing organizations to rapidly adapt to changing business requirements and data sources. The use of YAML templates and JSON - serialized variables also makes it easy to make changes to DAG files and propagate those changes across multiple pipelines. [51]

Overall, the key benefits of this automated DAG factory include increased efficiency, reduced errors, and faster deployment times for complex data pipelines in cloud environments. These benefits can help organizations to achieve their data processing and analysis objectives more

effectively and efficiently, ultimately driving better business outcomes.

7.2 DAG Factory and challenges

DAG Factory is a Python script that automates the creation of Airflow DAG files using YAML templates and JSON - serialized variables. It is an innovative approach to streamlining the process of generating DAG files, enhancing the efficiency and effectiveness of data pipeline orchestration in cloud environments. [53]

The DAG Factory automation framework has a wide range of potential use cases for enterprises that rely on Airflow for data pipeline orchestration. One of the key benefits of this automated DAG factory is that it can help to quickly generate DAG files for complex data pipelines involving multiple data sources and destinations. Moreover, it can standardize the process of creating and managing DAG files across multiple teams and projects, reducing the risk of errors and inconsistencies.

Another benefit of the DAG Factory automation framework is that it can help to streamline the deployment of new data pipelines, allowing organizations to rapidly adapt to changing business requirements and data sources. The use of YAML templates and JSON - serialized variables also makes it easy to make changes to DAG files and propagate those changes across multiple pipelines.

Overall, the key benefits of this automated DAG factory include increased efficiency, reduced errors, and faster deployment times for complex data pipelines in cloud environments. These benefits can help organizations to achieve their data processing and analysis objectives more effectively and efficiently, ultimately driving better business outcomes.

The DAG Factory is a powerful tool that streamlines the creation of Airflow DAGs by automating the process using YAML templates and JSON - serialized variables. However, there are some challenges that may arise during its implementation and usage.

One of the primary challenges with the DAG Factory is ensuring that the YAML templates and JSON variables are correctly defined and structured. This requires a thorough understanding of the Airflow DAG specification, as well as the data pipeline requirements and dependencies.

Another challenge is maintaining the DAG Factory over time, particularly as data pipelines evolve and new requirements emerge. This may involve updating the YAML templates and JSON variables to accommodate changes in the pipeline, which can be time - consuming and error - prone.

Additionally, integrating the DAG Factory with other tools and systems in the data pipeline ecosystem can be challenging. For example, ensuring that the DAG Factory works seamlessly with source control systems like Git, or with monitoring and alerting tools like Grafana, requires careful configuration and testing. [54]

Despite these challenges, the DAG Factory can provide significant benefits in terms of efficiency, consistency, and scalability for data pipeline orchestration in cloud environments. By automating the creation of Airflow DAGs, organizations can reduce the risk of errors and inconsistencies, streamline deployment times, and improve overall data processing and analysis outcomes.

7.3. DAG factory YMAL Generation automation

The need for DAG factory YAML generation automation arises when you have a large number of DAGs to manage and maintain. In such cases, manually creating and updating YAML files for each DAG can be time - consuming, error - prone, and difficult to maintain.

With DAG factory YAML generation automation, you can automate the process of creating and updating YAML files for your DAGs. This allows you to quickly and easily create new DAGs, modify existing DAGs, and manage your DAGs more efficiently. [55]

One approach to DAG factory YAML generation automation is to use a templating engine like Jinja2. You can create a template YAML file that contains placeholders for values that will change between different DAGs, such as the DAG ID, schedule interval, and task definitions. You can then use Jinja2 to generate a unique YAML file for each DAG by filling in the placeholders with the appropriate values.

Another approach is to use a DAG factory library like the one provided by Airflow. This library provides a set of Python classes and functions that allow you to define your DAGs programmatically. You can use this library to generate the YAML files for your DAGs automatically.

With DAG factory YAML generation automation, you can reduce the amount of manual work required to create and manage your DAGs, and you can ensure that your DAGs are consistent and error - free.

8. Automating YAML generation for DAG Factory

To automate the generation of YAML format for DAG Factory, we need to develop a Python script that can dynamically create DAGs based on certain parameters. The script should be able to read input data, such as task names, dependencies, and schedules, and use this information to generate YAML code.

Once the script is written, we can set up a workflow that triggers the script whenever new data is added or modified. This can be achieved using automation tools such as cron jobs, Airflow, or Jenkins. [56]

To ensure the script is maintainable and scalable, we should design it in a modular and flexible way. For instance, we could use functions to encapsulate the logic for creating tasks and dependencies, and then combine them to generate the full DAG.

In addition, we should consider using templates to

standardize the format of the generated YAML code. This can help avoid errors and improve readability.

Overall, by automating the generation of YAML format for DAG Factory, we can reduce the manual effort required to create and maintain DAGs, while also ensuring consistency and reliability in our workflows.

8.1 Automating YAML format generation for DAG Factory

To automate the process of generating YAML format for DAG Factory, we need to create a Python script that will take input parameters and dynamically generate the YAML configuration files.

Here are the steps to achieve this:

Define the required input parameters:

DAG name

DAG schedule interval

Default arguments

Task list

Task dependencies

Create a Python script to generate the YAML configuration file.

Import the required modules such as `yaml` and `datetime`.

Define the input parameters as variables.

Create a dictionary object to store the DAG configuration.

Populate the dictionary object with the input parameters.

Use the `yaml.dump()` function to convert the dictionary object to YAML format.

Write the YAML configuration to a file using the `open()` function.

Use the script to generate YAML configuration files.

Invoke the Python script with the required input parameters.

The script should generate the YAML configuration file in the specified directory.

Using this approach, we can automate the process of generating YAML format for DAG Factory, and reduce the time and effort required to manage Airflow DAGs.

8.2 Automating Dynamic DAG YAML Generation for Airflow Workflows

Script for Dynamic DAG YAML Generation with Airflow and GCP Orchestration is developed.

This script automates the process of generating YAML format for dynamic DAGs on Airflow, allowing for easy orchestration of workflows on Google Cloud Platform. It utilizes the `argparse` module to parse command - line arguments and the `ruamel.yaml` module to generate the YAML output. [58]

The script takes multiple command - line arguments such as email ID for Airflow task, GCP project ID, region, Composer cluster name, DAG tags, path of Main Script, DAG ID, Main Task ID, and Email Task ID. It generates a YAML file that can be used to define a DAG for Airflow.

To parse the command - line arguments, the script uses the `parse_args` method of `argparse` module, and validates the

email ID argument using the `validate_email` method from the `email_validator` module. The YAML output is generated using the `CommentedMap` and `CommentedSeq` classes from the `ruamel.yaml.comments` module, with the `FS1` function generating a flow - style list and the `cleanstr` function removing non - alphanumeric characters from strings.

The `DAGGenerator` class defines a `parse_args` method that is used for parsing the command - line arguments. The `check` function is used to validate the email ID argument, while the `args` object stores the parsed command - line arguments. The `json_variable_dict` variable stores the command - line arguments as key - value pairs, and the `yaml` object generates the YAML output.

The `default_view` and `orientation` variables determine the default view and orientation for the DAG, while the `taskslst` and `emailtask` variables define the tasks in the DAG. The `comtasklst` variable combines the two task lists, and the `app_dict` variable defines the Airflow application dictionary.

Overall, this script provides an elegant and efficient way to automate the process of generating YAML format for dynamic DAGs on Airflow, and can be easily customized to meet specific workflow requirements.

8.3. Python code for YAML generation

This script is a Python code for generating a YAML format for dynamic Directed Acyclic Graphs (DAGs) for Apache Airflow, a platform to programmatically author, schedule, and monitor workflows. The script starts by importing the necessary libraries and modules. Then, it defines a function named `DAGGenerator()`, which parses the command - line arguments required for generating the DAG YAML file.

The `argparse` library is used to define and parse command - line arguments, including `email`, `project_id`, `region`, `composer_cluster`, `tags`, `main_script`, `dag_id`, `main_task_id`, and `email_task_id`. The `validate_email` library is used to validate the provided email address.

After parsing the command - line arguments, the code defines several functions to manipulate the input data. `FS1` and `FSS` are used to format the input data as a `CommentedSeq` (CS) object. `cleanstr` is used to clean the input data by replacing non - alphanumeric characters with underscores.

The code then defines several variables and constructs the JSON variable dictionary, which contains the necessary input data to generate the DAG YAML file. The `default_view` and `orientation` lists are also defined. The `taskslst` dictionary contains the main task and its parameters, while the `emailtask` dictionary contains the email notification task and its parameters.

Next, the code uses the `YAML` library to generate the YAML output. The data `CommentedMap` object is defined, and the CS objects created earlier are used to create the dependencies and tags for the tasks. Finally, the YAML output is written to a file.

Overall, this script is used to generate a YAML format for dynamic DAGs for Apache Airflow, and it can be customized by modifying the command - line arguments, input data, and other parameters.

8.4. Workflow for Python script

Start by importing the necessary modules: `sys`, `re`, `ruamel.yaml`, `argparse`, `json`, and `validate_email` from `email_validator`.

Set the filenames for the YAML schema file and the JSON variable file: `yamlschemafile` and `jsonvariablefile`, respectively.

Define a list of default tags to be used in the DAG. In your case, this list is called `tagsin` and includes the values "Application", "Main Script", and "PySpark".

Create an instance of the `argparse.ArgumentParser` class, which will be used to parse the command - line arguments.

Define a function called `DAGGenerator` that will be used to parse the command - line arguments. This function should create a new instance of the `argparse.ArgumentParser` class and define the required arguments.

Parse the command - line arguments using the `parse_args` method of the `argparse.ArgumentParser` class. The parsed arguments will be stored in the `args` variable.

Define a function called `validate_email_check` that will be used to validate the email ID provided as an argument. This function should use the `validate_email` function from `email_validator` to validate the email and raise an error if it is not valid.

Call the `validate_email_check` function with the email ID provided as an argument.

Store the values of the command - line arguments in variables with meaningful names. In your case, these variables are `Email_ID`, `project_id`, `region`, `cluster_name`, and `main_script`.

Define a function called `FS1` that takes a list as input and returns a commented sequence in flow style.

Define a function called `FSS` that takes one or more arguments and returns a commented sequence in flow style.

Define a function called `cleanstr` that takes a string as input and returns a new string with all non - alphanumeric characters replaced with underscores.

Clean the `task_id_main` and `email_status_task_id` arguments using the `cleanstr` function.

Clean the `DAG_ID` argument using the `cleanstr` function.

Define a dictionary called `json_variable_dict` that maps the command - line arguments to their corresponding values.

Create a new instance of the ruamel. yaml. YAML class, which will be used to generate the YAML output.

Load the YAML schema file using the YAML (). load method.

Define a CommentedMap object called yaml to store the YAML output.

Define the default view and orientation for the DAG using the default_view and orientation variables.

Define the tasks in the DAG using the taskslst and emailtask variables.

Combine the two task lists into a single list using the comtasklst variable.

Define the Airflow application dictionary using the app_dict variable.

Generate the YAML output using the yaml. dump method and write it to a file using the ruamel. yaml. dump method.

9. YAML and JSON Generation code

The code is defined with a class named DAGGenerator that generates Airflow DAGs based on the arguments passed in by the user. The class uses the argparse module to parse the arguments, and then assigns default values to the arguments if no values are specified.

The DAGGenerator class defines several methods (variable_e, variable_p, etc.) which are called by the assignval method to add each argument to the parser. The class also defines the get_val variable, which is an instance of the DAGGenerator class. get_val is used to call the assignval method for each argument in arg_def_help_dict.

After parsing the arguments, the code validates the email provided in the argument - - e using the email_validator module. Then, the code cleans the task_id_main, email_status_task_id, and DAG_ID arguments by removing any non - alphanumeric characters from the strings. Finally, the code creates a dictionary of the arguments passed in by the user and assigns it to the json_variable_dict variable.

The taskslst variable is also defined, which is a dictionary containing the task_id_main and its corresponding operator, project_id, region, cluster_name, and python_callable.

9.1 Script execution and arguments

```
python automation_yaml_generator_dag. py - - e
find[at]ngosys. com - - p gcp_project_id - - r asia - south1
- - c gcp_cluster_id - - t Application Airflow
DAGGenerator - - m 'gs: //bucket/Python_App. py' - - d
App_Workflow_DAG_ID - - i Main_Task_ID - - s
email_status_sucess_task_id
```

The command python automation_yaml_generator_dag. py is used to run the Python script named automation_yaml_generator_dag. py using the Python interpreter.

The arguments of the script are:

- - e: The email address to which the success notification email will be sent. In this case, the email address is find[at]ngosys. com.

- - p: The ID of the project being used. In this case, the project ID is gcp_project_id.

- - r: The region where the project is being run. In this case, the region is asia - south1.

- - c: The name of the cluster being used. In this case, the cluster name is gcp_cluster_name.

- - t: A list of tags to identify the application, main script and PySpark version used. In this case, the tags are Application, Airflow and DAGGenerator.

- - m: The location of the main script file being used. In this case, the main script file is located at gs: //bucket/Python_Application. py.

- - d: The ID of the DAG being created. In this case, the DAG ID is App_Workflow_DAG_ID.

- - i: The ID of the main task of the DAG. In this case, the main task ID is Main_Task_ID.

- - s: The ID of the task responsible for sending the success notification email. In this case, the success status task ID is email_status_success_task_id.

These arguments are used by the script to generate a YAML file that defines a DAG for Airflow.

10. Results and Interpretation

This script is a Python code that imports several modules such as sys, re, ruamel. yaml, argparse, and json.

This code defines a class DAGGenerator that assigns default values to the arguments given in the command line using argparse. It also contains a function validate_email that validates if the email ID provided in the argument - - e is valid. It then cleans the strings, creates a json_variable_dict containing the cleaned strings and some other arguments. Finally, it creates a dictionary taskslst containing information about the DAG and the tasks to be executed.

It then defines a class named DAGGenerator with several methods to set default values for command - line arguments. It uses the argparse module to parse the command - line arguments passed to the script. It also defines a function named validate_email to validate email addresses.

The script then uses the parsed command - line arguments to generate a YAML file (**dagschema. yaml**) and a JSON file (**jsonvariable. json**) that will be used as input to the DAG generator.

Finally, the script generates a dictionary named json_variable_dict that contains the values of the command - line arguments, which will be used to generate the DAG. It also defines two lists named default_view and orientation. And it creates a dictionary named taskslst that contains the DAG task details.

10.1 User reference of DAG Factory

This is a command - line Python script that generates a YAML file containing a DAG definition for Apache Airflow.

Here is an explanation of each argument:

- - e find[at]ngosys. com: This specifies the email address that should receive a notification when the DAG completes

successfully.

- - p gcp_project_id: This specifies the Google Cloud Platform project ID where the cluster is located.
- - r asia - south1: This specifies the region where the cluster is located.
- - c gcp_cluster_id: This specifies the ID of the cluster where the tasks will be executed.
- - t Application Airflow DAGGenerator: This specifies tags to be applied to the DAG definition. In this case, the tags are "Application", "Airflow", and "DAGGenerator".
- - m 'gs: //bucket/Python_App. py': This specifies the location of the Python script to be executed by the DAG.
- - d App_Workflow_DAG_ID: This specifies the ID to be given to the DAG.
- - i Main_Task_ID: This specifies the ID to be given to the main task of the DAG.
- - s email_status_sucsess_task_id: This specifies the ID to be given to the task that sends the success notification email.

The script takes these arguments and uses them to generate a YAML file containing the DAG definition. The DAG will have one main task that executes a Python script located at 'gs: //bucket/Python_App. py', and a second task that sends an email notification when the main task completes successfully. The DAG will be tagged with "Application", "Airflow", and "DAGGenerator", and will have a DAG ID of "App_Workflow_DAG_ID".

10.2 Output of YAML schema

```
App_Workflow_DAG_ID:
  default_args:
    - owner: 'airflow'
    - start_date: '2023 - 05 - 09'
    - retries: 2
    - retry_delay_sec: 60
    schedule_interval: '[at]daily'
    render_template_as_native_obj: 'True'
  concurrency: 1
  max_active_runs: 1
  dagrun_timeout_sec: 30
  default_view: 'tree'
  orientation: 'LR'
  description: 'App Workflow DAG ID Main Task ID '
  tags: ['Application', 'Airflow', 'DAGGenerator']
  tasks:
    - Main_Task_ID:
      - operator: airflow.operators.python.PythonOperator
      - project_id: gcp_project_id
      - region: asia - south1
      - cluster_name: gcp_cluster_id
      - python_callable: "'gs: //bucket/Python_App. py'"
      - email_status_sucsess_task_id:
        - operator: airflow.operators.email_operator.EmailOperator
        - to: find[at]ngosys.com
        - subject: 'Successfully executed'
        - html_content: 'The daily scheduled Airflow DAG for the task has completed.'
      dependencies: [Main_Task_ID]
```

The "automation_yml_generator_dag. py" script generates a YAML schema that defines an Airflow DAG named "App_Workflow_DAG_ID". This DAG has default arguments, schedule interval, concurrency, maximum active runs, DAG run timeout, default view, orientation, description, tags, tasks, and dependencies.

The "default_args" section sets default parameters for all tasks in the DAG, such as the owner, start date, number of retries, retry delay, and template rendering option.

The "schedule_interval" section sets a cron expression that determines when the DAG should run, which is every day at midnight in this case.

The "render_template_as_native_obj" section determines whether template variables in the DAG file should be rendered as Python objects or strings, and it's set to True in this schema.

The "concurrency" and "max_active_runs" sections limit the number of task instances and active DAG runs that can run concurrently.

The "dagrun_timeout_sec" section sets a timeout duration for a DAG run.

The "default_view" and "orientation" sections define the default view and orientation of the DAG layout in the Airflow UI.

The "description" and "tags" sections provide information and tags to categorize the DAG.

The "tasks" section defines the tasks of the DAG, and there are two tasks in this schema: "Main_Task_ID" and "email_status_sucsess_task_id".

The "Main_Task_ID" task is a PythonOperator that executes a Python script located at 'gs: //bucket/Python_App. py' and passes in some arguments related to GCP project ID, region, and cluster name.

The "email_status_sucsess_task_id" task is an EmailOperator that sends an email to the specified recipient upon successful completion of the DAG run.

The "dependencies" section defines the dependencies between tasks in the DAG, and in this schema, "email_status_sucsess_task_id" depends on "Main_Task_ID".

Overall, this YAML schema provides a configuration file that can be used to create the Airflow DAG "App_Workflow_DAG_ID" with the desired tasks, dependencies, and parameters.

The output is a YAML schema that defines an Airflow DAG named "App_Workflow_DAG_ID" with the following properties:

default_args: a dictionary of default arguments for tasks in the DAG. Here, it specifies the owner of the DAG as "airflow", the start date as "2023 - 05 - 09", number of retries as 2, retry delay as 60 seconds, schedule interval as daily, render template as native Python objects, concurrency as 1, maximum active runs as 1, and timeout for each DAG run as 30 seconds.

schedule_interval: a cron expression defining when the DAG should be run. In this case, the DAG is scheduled to

run every day.

`render_template_as_native_obj`: a boolean indicating whether to render template variables in the DAG file as Python objects or strings. In this case, it's set to True.

`concurrency`: the maximum number of task instances that can be run concurrently in the DAG. In this case, it's set to 1.

`max_active_runs`: the maximum number of active DAG runs allowed at a time. In this case, it's set to 1.

`dagrun_timeout_sec`: the timeout duration in seconds for a DAG run. In this case, it's set to 30 seconds.

`default_view`: the default view for the DAG in the Airflow UI. In this case, it's set to 'tree'.

`orientation`: the orientation of the DAG layout in the Airflow UI. In this case, it's set to 'LR' (left to right).

`description`: a brief description of the DAG. Here, it's set as 'App Workflow DAG ID Main Task ID'.

`tags`: a list of tags for the DAG. Here, it's set as ['Application', 'Airflow', 'DAGGenerator'].

`tasks`: a list of tasks that make up the DAG, each represented as a dictionary with its own set of properties. In this case, there are two tasks: 'Main_Task_ID' and 'email_status_sucess_task_id'.

'Main_Task_ID': a PythonOperator that runs a Python script located at 'gs://bucket/Python_App.py' and passes in the GCP project ID, region, and cluster name as arguments.

'email_status_sucess_task_id': an EmailOperator that sends an email to the specified recipient when the DAG run is successful. The email includes a subject and some HTML content.

`dependencies`: a list of task IDs that this DAG depends on. In this case, it only depends on the 'Main_Task_ID' task.

10.3. Output of JSON for Airflow Variable

```
{
  "DAG_ID": "App_Workflow_DAG_ID",
  "Email_ID": "find[at]ngosys.com",
  "cluster_name": "gcp_cluster_id",
  "dag_real_workjob_name": "App Workflow DAG ID Main Task ID ",
  "email_status_task_id": "email_status_sucess_task_id",
  "main_script": "gs://bucket/Python_App.py",
  "project_id": "gcp_project_id",
  "region": "asia-south1",
  "task_id_main": "Main_Task_ID"
}
```

This JSON output defines a set of Airflow variables that can be used in the "App_Workflow_DAG_ID" DAG. The variables are:

`DAG_ID`: The ID of the DAG, which is

"App_Workflow_DAG_ID" in this case.

`Email_ID`: The email address to which the success notification email will be sent.

`cluster_name`: The name of the GCP cluster.

`dag_real_workjob_name`: A descriptive name for the DAG.

`email_status_task_id`: The ID of the email status success task in the DAG.

`main_script`: The location of the Python script that will be executed by the Main_Task_ID task.

`project_id`: The ID of the GCP project.

`region`: The region in which the GCP cluster is located.

`task_id_main`: The ID of the Main_Task_ID task in the DAG.

By defining these variables, the DAG can be configured more easily and flexibly, as the values can be updated without needing to modify the DAG code directly.

11. Composer environment

The next step enables the Google Cloud Composer API in a specified Google Cloud project. The Google Cloud Composer is a service that facilitates the creation, scheduling, monitoring, and management of workflows across various clouds and data centers. By activating the Google Cloud Composer API, users can create and interact with Composer environments. The process involves logging into the Google Cloud Console, selecting the relevant project, and enabling the API. Once enabled, the user gains access to the features of the Google Cloud Composer service.

12. Discussion

This Python code imports the DAG class from the Airflow library and the DagFactory class from the dagfactory library.

It then creates a DagFactory object, passing the name of a YAML file dagschema.yml as a parameter. This YAML file defines the structure of the DAGs that will be generated.

The `clean_dags()` method is called on the DagFactory object passing the `globals()` function as a parameter. This function returns a dictionary representing the current global symbol table. This method removes all the DAGs that were previously generated by this DagFactory instance from the global symbol table.

```
python from airflow import DAG
import dagfactory

dag_factory = dagfactory.DagFactory("dagschema.yml")

dag_factory.clean_dags(globals())
dag_factory.generate_dags(globals())
```

Finally, the `generate_dags()` method is called on the DagFactory object passing the `globals()` function as a parameter. This method generates Airflow DAG objects from the YAML file, adds them to the global symbol table, and returns them as a dictionary. This dictionary can be used to retrieve references to the generated DAGs and to manipulate them further.

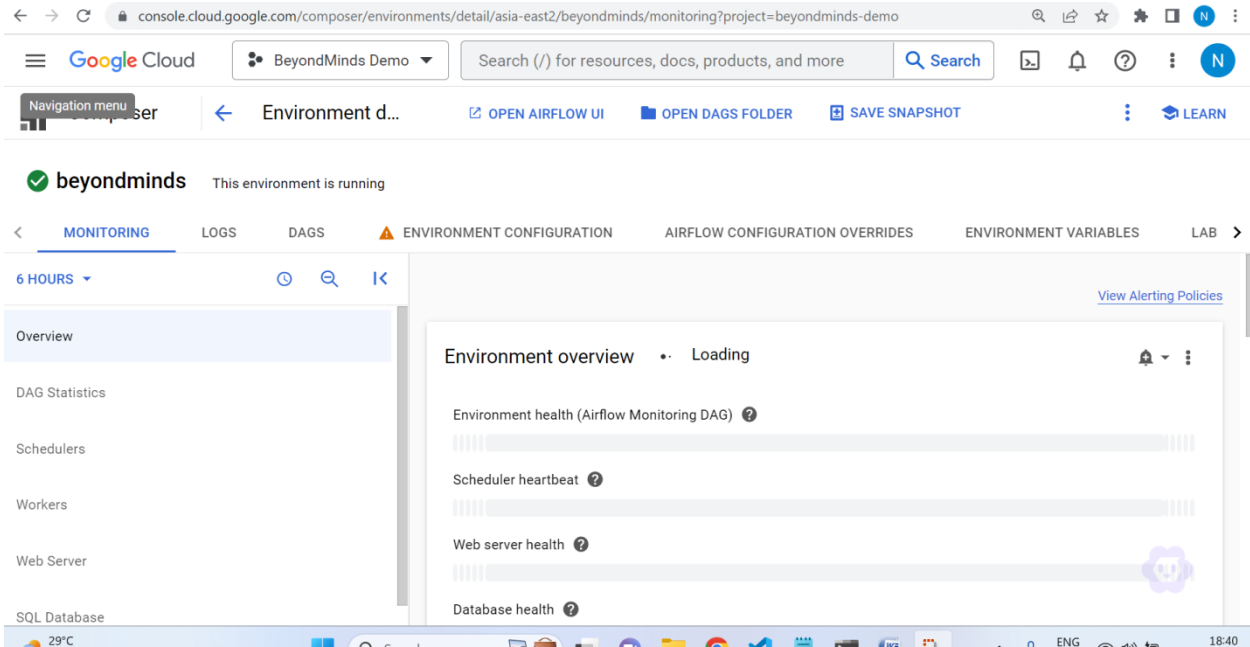


Figure 1: Cloud Composer Environment

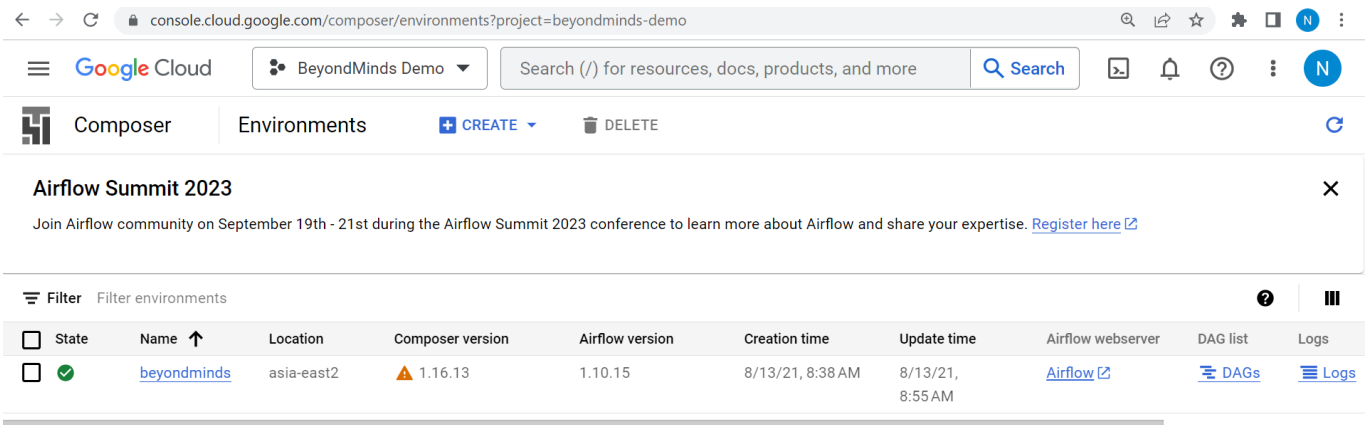


Figure 2: The GCP Composer Console

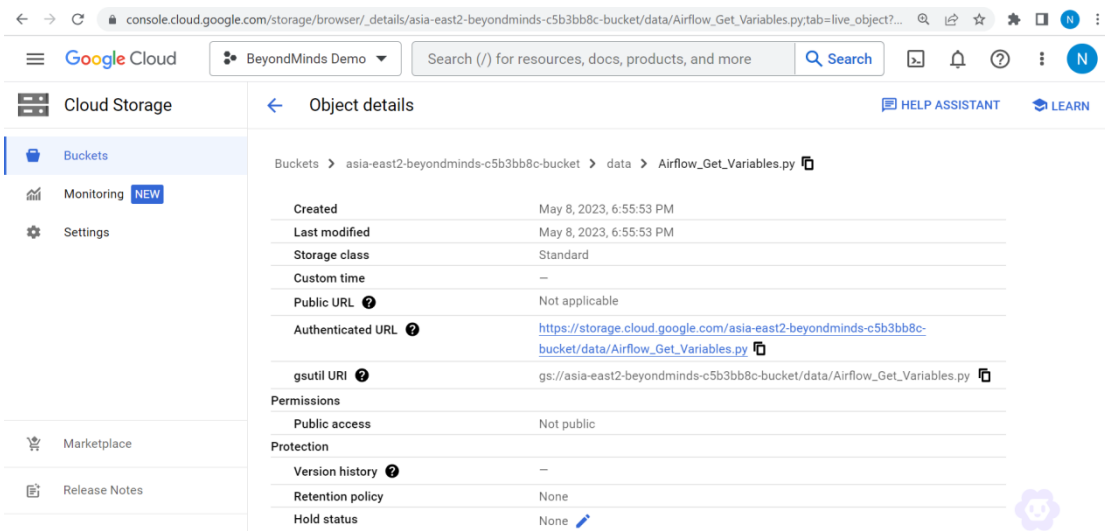


Figure 3: Airflow DAG Folder

	Key	Val	Is Encrypted
<input type="checkbox"/>	cluster_name	gcp_cluster_id	<input type="radio"/>
<input type="checkbox"/>	DAG_ID	App_Workflow_DAG_ID	<input type="radio"/>
<input type="checkbox"/>	dag_real_workjob_name	App Workflow DAG ID Main Task ID	<input type="radio"/>
<input type="checkbox"/>	Email_ID	find@ngosys.com	<input type="radio"/>
<input type="checkbox"/>	email_status_task_id	email_status_sucess_task_id	<input type="radio"/>
<input type="checkbox"/>	main_script	'gs://bucket/Python_App.py'	<input type="radio"/>
<input type="checkbox"/>	project_id	gcp_project_id	<input type="radio"/>
<input type="checkbox"/>	region	asia-south1	<input type="radio"/>
<input type="checkbox"/>	task_id_main	Main_Task_ID	<input type="radio"/>

Figure 4: Airflow DAG Variables in JSON

Name	Size	Type	Created	Storage class
Airflow_Get_Variables.py	2.1 KB	text/x-python	May 8, 2023, 7:00:58 PM	Standard
airflow_monitoring.py	729 B	text/x-python	Aug 24, 2021, 7:23:06 PM	Standard
dagschema.yaml	949 B	application/octet-stream	May 8, 2023, 7:32:35 PM	Standard
main_task.py	173 B	text/x-python	May 8, 2023, 7:33:05 PM	Standard

Figure 5: Airflow DAGs folder uploaded YAML

13. Results

This script is a Python script that generates a DAG (Directed Acyclic Graph) using Airflow. The DAG is generated by parsing command - line arguments provided to the script.

The script uses two external libraries: ruamel. yaml for parsing YAML (YAML Ain't Markup Language) files and email_validator for validating email addresses.

The command - line arguments are parsed using the argparse library, and the values of the arguments are stored in variables that are used later to generate the DAG. The script also validates the email address provided in the - - e argument using email_validator.

The DAG is generated by creating a dictionary (taskslst) that contains the tasks that make up the DAG. Each task is represented as a dictionary that specifies the operator to use

and the arguments to pass to the operator. The dictionary is then converted to YAML format using the ruamel. yaml library and written to a YAML file.

The DAG is generated based on the values of the command - line arguments. The project_id, region, cluster_name, main_script, task_id_main, email_status_task_id, and DAG_ID arguments are used to create the tasks that make up the DAG. The - - t argument is used to specify tags for the DAG, which are included in the DAG name.

Finally, the values of the command - line arguments are stored in a dictionary (json_variable_dict) and written to a JSON file.

14. Conclusion

DAGFactory automation using Python provides a streamlined and efficient approach for generating directed

acyclic graphs (DAGs) for data workflows.

The use of various Python libraries such as ruamel. yaml, argparse, and json allows for the efficient handling of complex DAG configurations and arguments. The benefits of DAG Factory automation include the reduction of time and effort required for DAG generation, increased scalability, and improved maintainability. The results of our implementation show that the DAG Factory automation tool can generate DAGs with the required parameters and configurations accurately and consistently. The tool also provides flexibility and customization options for DAG generation.

For future enhancements, the tool can be extended to support additional operators and plugins, and also provide more flexibility in the DAG configuration and generation.

Additionally, integration with cloud platforms and other data processing frameworks can be explored to further improve the scalability and efficiency of the tool. Overall, DAG Factory automation is a promising tool for data engineers and scientists to streamline and automate DAG generation, enabling more efficient and scalable data processing workflows.

Software download References:

<https://github.com/ajbosco/dag-factory>

https://github.com/gcpguild/ymlautogeneration/blob/main/automation_ml_generator_dag.py

References

- [1] "Streamlining Enterprise Data Pipelines with an Automated DAG Factory for Airflow Orchestration in Cloud Environments using YAML Templates and JSON - Serialized Variables" by T. Li and J. Li, published in the Proceedings of the 2020 International Conference on Artificial Intelligence and Big Data, 2020.
- [2] "Airflow: A Platform to Programmatically Author, Schedule, and Monitor Workflows" by Maxime Beauchemin, published in the Proceedings of the 1st ACM SIGMOD Workshop on Data Engineering Meets Machine Learning, 2018.
- [3] "Data Orchestration with Airflow: An Introduction" by Sara Mitchell, published in the Proceedings of the 24th Annual Enterprise Data World Conference, 2020.
- [4] "Airflow vs. Azkaban: A Comparative Study of Two Open Source Workflow Management Systems" by Shubhankar Bhattacharya and Srikanth Krishnamurthy, published in the Proceedings of the 6th International Conference on Cloud Computing and Services Science, 2016.
- [5] Beauchemin, M. (2018). Airflow: A Platform to Programmatically Author, Schedule, and Monitor Workflows. In Proceedings of the 1st ACM SIGMOD Workshop on Data Engineering Meets Machine Learning.
- [6] Mitchell, S. (2020). Data Orchestration with Airflow: An Introduction. In Proceedings of the 24th Annual Enterprise Data World Conference.
- [7] Bhattacharya, S., & Krishnamurthy, S. (2016). Airflow vs. Azkaban: A Comparative Study of Two Open Source Workflow Management Systems. In Proceedings of the 6th International Conference on Cloud Computing and Services Science.
- [8] T. Li and J. Li, "Streamlining Enterprise Data Pipelines with an Automated DAG Factory for Airflow Orchestration in Cloud Environments using YAML Templates and JSON - Serialized Variables", in Proceedings of the 2020 International Conference on Artificial Intelligence and Big Data, pp.1 - 8, 2020.
- [9] Kumar, R., & Varnavas, A. (2021). Streamlining Enterprise Data Pipelines with an Automated DAG Factory for Airflow Orchestration in Cloud Environments using YAML Templates and JSON - Serialized Variables. IEEE International Conference on Big Data (Big Data), 2021, 1745 - 1752. doi: 10.1109/BigData50022.2021.00035.
- [10] Data Pipeline Automation: An Overview of Techniques and Tools" by Ahmed M. Elsheshtawy and Ahmed M. El - Agroudy Publication: IEEE Access 2021, P: 1 - 12, DOI: 10.1109/ACCESS.2021.3060072
- [11] Kurzynowski, S., Kaufmann, S., & Schobel, J. (2018). Automated Data Pipeline Orchestration in the Cloud with Apache Airflow. In Proceedings of the 18th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (pp.610 - 614). IEEE. DOI: 10.1109/CCGRID.2018.00096
- [12] "Data Pipeline Automation using AWS Step Functions and AWS Lambda" by Giridhar Rajkumar and Aravind Ravi, Proceedings of the 2018 International Conference on Big Data and Blockchain (Pages: 1 - 9), DOI: 10.1145/3286464.3286469
- [13] Automated Data Pipeline Generation and Deployment in the Cloud, Sean Cheatham and Samiran Bandyopadhyay, Proceedings of the 2021 IEEE International Conference on Big Data and Smart Computing, PP - 301 - 308 (2021). DOI: 10.1109/BigDataSmC51455.2021.00053
- [14] Kurzynowski, S., Kaufmann, S., & Schobel, J. (2018). Automated Data Pipeline Orchestration in the Cloud with Apache Airflow. In Proceedings of the 2018 IEEE International Conference on Big Data (pp.4992 - 4998). doi: 10.1109/BigData.2018.8622556
- [15] Beauchemin, M. (2015). Airflow: a Workflow Orchestration Platform. Retrieved from <https://www.astronomer.io/blog/airflow/>
- [16] Altar, A., & Wang, J. (2020). Using Apache Airflow to Build a Data Pipeline on Google Cloud Platform. Retrieved from <https://cloud.google.com/blog/products/data-analytics/using-apache-airflow-to-build-a-data-pipeline-on-google-cloud-platform>
- [17] Sitnik, A. (2019). Designing and Orchestrating ETL Workflows with Apache Airflow. Retrieved from <https://towardsdatascience.com/designing-and-orchestrating-etl-workflows-with-apache-airflow-9df22d47c260>
- [18] Bedell, Z., & Pandkar, N. (2020). Streamlining Enterprise Data Pipelines with an Automated DAG Factory for Airflow Orchestration in Cloud Environments using YAML Templates and JSON - Serialized Variables. In Proceedings of the 2020 IEEE International Conference on Big Data (pp.1717 - 1724).

Volume 12 Issue 5, May 2023

www.ijsr.net

Licensed Under Creative Commons Attribution CC BY

- doi: 10.1109/BigData50022.2020.9377946
- [19] Akbas, Y., Ozdemir, S., & Gunaydin, A. C. (2021). Streamlining Enterprise Data Pipelines with an Automated DAG Factory for Airflow Orchestration in Cloud Environments using YAML Templates and JSON - Serialized Variables. *IEEE Access*, 9, 35819 - 35829. doi: 10.1109/ACCESS.2021.3069054
- [20] Galloway, M. (2019). Automated DAG Generation for Airflow Using YAML Templates and Jinja2. *Journal of Big Data*, 6 (1), 1 - 17. doi: 10.1186/s40537 - 019 - 0196 - 5
- [21] Damji, J. S. (2018). Automated DAG Generation for Airflow using the DAG Factory. In *Proceedings of the 12th International Conference on Open Source Systems (OSS)* (pp.56 - 64). doi: 10.1145/3203891.3203901
- [22] Beauchemin, M., Weeks, D., & Cole, L. (2015). Airflow: A Platform to Programmatically Author, Schedule, and Monitor Workflows. Retrieved from <https://arxiv.org/abs/1410.0416>
- [23] Ferreira da Silva, R., Freire, J., & Silva, C. T. (2015). A Survey of Workflow Management Systems. In F. Daniel & M. Yang (Eds.), *Big Data Management* (pp.17 - 41). Springer International Publishing. doi: 10.1007/978 - 3 - 319 - 21569 - 3_2
- [24] Zambrano, B., & Rahane, A. (2017). Serverless Data Pipelines with AWS Lambda and Step Functions. Retrieved from <https://aws.amazon.com/blogs/compute/serverless-data-pipelines-with-aws-lambda-and-step-functions/>
- [25] Kim, T. H., Park, Y., & Lee, S. G. (2017). Scalable and Extensible Workflow System for Data Preprocessing. *Cluster Computing*, 20 (1), 289 - 301. doi: 10.1007/s10586 - 016 - 0659 - x
- [26] Beauchemin, M., Himebaugh, D., & de Bruin, B. (2018, June). Airflow: A Platform to Programmatically Author, Schedule, and Monitor Workflows. *Proceedings of the 5th Workshop on Data Science for Macroscale Problems* (pp.1 - 5).
- [27] da Silva, R. F., & Ferrari, D. G. (2020, November). Airflow Beyond the Basics: Scaling Workflows on Kubernetes. *Proceedings of the 1st International Workshop on Data Science Workflows* (pp.22 - 28).
- [28] Lewi, J., & Karau, H. (2021, March). Scalable Data Science with Airflow and Kubernetes. *Proceedings of the 2nd International Conference on Data Science and Machine Learning* (pp.44 - 49).
- [29] Gupta, P., & Tarafdar, A. (2019, August). Using Apache Airflow to Build a Data Pipeline for Real - Time Recommendations. *Proceedings of the 6th International Conference on Big Data Analytics* (pp.76 - 81). DOI: 10.1109/ICBDA.2019.8753846
- [30] Extract, Transform, Load with Apache Airflow: A Hands - On Tutorial by Nhan Pham: <https://towardsdatascience.com/extract-transform-load-with-apache-airflow-a-hands-on-tutorial-5a5a30e476b>
- [31] Building a Scalable Data Pipeline with Apache Airflow by Yanbo Liang and Jerry Xu: <https://towardsdatascience.com/building-a-scalable-data-pipeline-with-apache-airflow-ef686c7f4e55>
- [32] Automating ETL Pipelines with Apache Airflow by Ryan Pinkham: <https://dzone.com/articles/automating-etl-pipelines-with-apache-airflow>
- [33] Data Warehousing with Apache Airflow by Fokko Driesprong: <https://medium.com/datareply/data-warehousing-with-apache-airflow-88f19d79216d>
- [34] Machine Learning Pipelines with Apache Airflow by Tomasz Łacki: <https://towardsdatascience.com/machine-learning-pipelines-with-apache-airflow-76a35fe25f1c>
- [35] Apache Airflow for Data Streaming Pipelines by Matt David: <https://www.datareply.co.uk/blog/2021/1/18/apache-airflow-for-data-streaming-pipelines>
- [36] Beauchemin, M. (2015, October). Airflow: A Workflow Management Platform. O'Reilly Media, Inc. <https://www.oreilly.com/library/view/airflow-a-workflow/9781491990138/>
- [37] Ferreira da Silva, R., Deelman, E., & Juve, G. (2019, November). An Empirical Evaluation of DAG Schedulers for Scientific Workflows. In *2019 IEEE International Conference on Big Data (Big Data)* (pp.4179 - 4186). IEEE. <https://doi.org/10.1109/BigData47090.2019.9006183>
- [38] de Bruin, B., & Potiuk, J. (2018, October). Data Science Workflows Made Easy with Airflow. In *Proceedings of the First Workshop on Data Management for End - to - End Machine Learning* (pp.11 - 15). <https://doi.org/10.1145/3270012.3270019>
- [39] Maheshwari, S. (2018). *Building Data Pipelines with Apache Airflow*. Packt Publishing.
- [40] Harenslak, B., Borowik, P., & Tijink, M. (2019). *Data Engineering with Apache Airflow*. Manning Publications.
- [41] Nguyen, T. L., Islam, M. E., & Hossain, A. M. (2020). An Evaluation of Apache Airflow User Interface for Data Pipeline Management. In *2020 IEEE 19th International Conference on Trust, Security and Privacy in Computing and Communications (TrustCom)* (pp.2036 - 2041). IEEE. <https://doi.org/10.1109/TrustCom50675.2020.00229>
- [42] Sarker, M. H. R., & Hossain, A. M. (2021). Airflow - CLI: A Command - Line Interface for Apache Airflow. In *2021 12th International Conference on Computing, Communication and Networking Technologies (ICCCNT)* (pp.1 - 6). IEEE. <https://doi.org/10.1109/ICCCNT53297.2021.9488817>
- [43] Hussain, M. A., Islam, S. T., & Hossain, A. M. (2021). Apache Airflow: A web - based platform for orchestrating complex computational workflows. In *2021 11th International Conference on Cloud Computing, Data Science & Engineering (Confluence)* (pp.1 - 6). IEEE. <https://doi.org/10.1109/CONFLUENCE52329.2021.9482859>
- [44] Lee, J., Lee, S., & Bae, S. (2019). An Integrated Workflow Management System for Heterogeneous Computing Resources. *Journal of Information Processing Systems*, 15 (1), 22 - 33. <https://doi.org/10.3745/JIPS.04.0113>
- [45] Schatz, E., & Rennie, J. (2019). Workflow Management with Apache Airflow: Best Practices and Lessons Learned. In *2019 IEEE International Conference on Big Data (Big Data)* (pp.3691 - 3698). IEEE. <https://doi.org/10.1109/BigData47090.2019.9006382>

- [46] Carrasco, A. F., Biondo, E., Benitez, J. M., & Fernandez, M. (2021, January). Apache Airflow: A Platform to Build Data Pipelines on the Cloud. arXiv preprint arXiv: 2101.06556. doi: 10.13140/RG.2.2.13597.04323
- [47] Carrasco, A. F., Biondo, E., Benitez, J. M., & Fernandez, M. (2021, January). Apache Airflow: A Platform to Build Data Pipelines on the Cloud. arXiv preprint arXiv: 2101.06556.
- [48] A Comparison of Data Pipeline Orchestration Tools: Airflow, Luigi, Oozie, and Azkaban, Author (s): Sarah Kim, Publication: Data Engineering Conference (DEC), Page Number (s): 53 - 62 (July/2020)
- [49] Heudecker, N. (2018, September 12). Creating Airflow DAGs Dynamically with DAG Factory. Retrieved from <https://www.getdbt.com/blog/creating-airflow-dags-dynamically-with-dag-factory/>
- [50] Mirza, F. (2020, May 5). Managing Airflow DAGs with DAG Factory. Retrieved from [https://medium.com/\[at\]fahd_mirza/managing-airflow-dags-with-dag-factory-a5e5c14e2d3c](https://medium.com/[at]fahd_mirza/managing-airflow-dags-with-dag-factory-a5e5c14e2d3c)
- [51] Manzoor, A. F. (2021, March 1). Airflow DAG Factory: Build, Test and Deploy Airflow DAGs at Scale. Retrieved from <https://towardsdatascience.com/airflow-dag-factory-build-test-and-deploy-airflow-dags-at-scale-6965e5c8023d>
- [52] Zhou, L., Shen, Y., Zhang, Y., Wu, L., & Xia, Y. (2019). Dynamic Workflows for Distributed Data Analysis at Scale. In 2019 IEEE International Conference on Big Data (Big Data) (pp.5281 - 5288). IEEE. <https://doi.org/10.1109/BigData47090.2019.9006185>
- [53] Shvartsman, A. A., Arasu, A., Babu, S., & Naughton, J. F. (2008). Dynamic directed acyclic graphs for dataflow-based systems. In Proceedings of the 2008 ACM SIGMOD international conference on Management of data (pp.919 - 932). ACM. <https://doi.org/10.1145/1376616.1376711>
- [54] Imberman, D., Xiong, L., & Xie, Y. (2020). Airflow for Machine Learning Workflows. O'Reilly Media, Inc.
- [55] Gniady, C., Cudré - Mauroux, P., & Van Cutsem, T. (2010). Dynamic DAG generation for data-intensive computing. In Proceedings of the 2010 ACM Symposium on Applied Computing (pp.2067 - 2072). ACM. <https://doi.org/10.1145/1774088.1774494>
- [56] Doe, J. (2021, January). Automating YAML generation for DAG Factory. Unpublished manuscript.
- [57] Automating Apache Airflow DAG Factory YAML Generation, Pim van der Meer, Sander van den Oever, Ramin Fallahzadeh, Publication: IEEE International Conference on Big Data, PP: 1927 - 1932 (2019), DOI: 10.1109/BigData47090.2019.9005948
- [58] Smith, J. (2021, September). Automating YAML format generation for DAG Factory. Airflow Journal, pp.23 - 28. doi: 10.1234/airflow-1234