

Data Orchestration: Modernizing Legacy Data Orchestration with Cloud Composer

Gautami Nadkarni

Department of Information Systems, State University of New York, University at Buffalo, USA

Abstract: *This paper discusses the migration of data orchestration workflows from a legacy tool like Autosys to a modern, cloud - based solution, Google Cloud Composer. It explores the transition from traditional job scheduling to Directed Acyclic Graph (DAG) - based workflows using Apache Airflow, culminating in the deployment and management of these workflows in Cloud Composer. The benefits and challenges of this migration are examined, highlighting the advantages of scalability, flexibility, and cloud integration offered by Cloud Composer.*

Keywords: Data Orchestration, Autosys, Apache Airflow, Cloud Composer, DAG, Migration, Cloud Computing, open source

1. Introduction

Data orchestration is a critical component of modern data engineering, ensuring that data pipelines run reliably and efficiently. Traditionally, tools like Autosys have been used to manage these workflows. However, the rise of cloud computing and the need for more flexible and scalable solutions have driven the adoption of modern orchestration platforms like Apache Airflow and Google Cloud Composer [1].

This paper explores the journey of migrating data orchestration from Autosys to Cloud Composer. It outlines the challenges of legacy systems, the benefits of Airflow's DAG - based approach, and the advantages of deploying Airflow on Cloud Composer.

2. Legacy Orchestration with Autosys

Autosys is a widely used job scheduling system that has served many organizations for decades. It is known for its reliability and robust scheduling capabilities. However, it also presents several limitations:

- **Rigid Structure:** Autosys often requires complex scripting and configuration, making it difficult to adapt to changing requirements.
- **Limited Scalability:** Scaling Autosys can be challenging and often requires significant infrastructure upgrades.
- **Lack of Cloud Integration:** Autosys is typically deployed on - premises, making it difficult to integrate with cloud - based services.
- **Maintenance Overhead:** Managing and maintaining Autosys can be time - consuming and require specialized expertise.

These limitations often hinder agility and innovation, prompting organizations to seek more modern solutions.

3. Transition to Apache Airflow

Apache Airflow, an open - source platform, is instrumental in orchestrating intricate workflows. These workflows are represented using Directed Acyclic Graphs (DAGs), which provide a visual representation of the workflow's structure.

This visualization allows for clear and intuitive definition of dependencies between tasks, ensuring that tasks are executed in the correct order. Airflow's flexibility and scalability make it a popular choice for managing data pipelines and other complex workflows in a variety of industries.

1) DAG - Based Workflows

Airflow's DAG - based approach offers several advantages:

- a) **Clarity and Visibility:** DAGs provide a visual representation of workflows, making it easier to understand and manage complex processes.
- b) **Flexibility:** Airflow supports various operators and integrations, allowing for diverse tasks to be included in workflows.
- c) **Scalability:** Airflow can be scaled horizontally by adding more worker nodes.

2) Advantages of Airflow

- a) **Open Source:** Airflow is an open - source project with a large and active community.
- b) **Extensibility:** Airflow can be extended with custom operators and integrations.
- c) **Monitoring and Logging:** Airflow provides robust monitoring and logging capabilities.

4. Moving to Google Cloud Composer

Google Cloud Composer is a managed orchestration service built on Apache Airflow. It simplifies the deployment and management of Airflow environments, allowing users to focus on building and running workflows.

1) Benefits of Cloud Composer

- a) **Managed Service:** Cloud Composer is a fully managed service, reducing the operational overhead of managing Airflow infrastructure.
- b) **Scalability and Reliability:** Cloud Composer provides automatic scaling and high availability.
- c) **Cloud Integration:** Cloud Composer seamlessly integrates with other Google Cloud services, such as BigQuery, Dataflow, and Cloud Storage [2].

2) Migration Process

The migration from Autosys to Cloud Composer typically involves the following steps:

Volume 14 Issue 4, April 2025

Fully Refereed | Open Access | Double Blind Peer Reviewed Journal

www.ijsr.net

a) Assessment:

- Conduct a thorough analysis of current Autosys workflows, meticulously documenting all job dependencies, schedules, and notification triggers.
- Identify any potential bottlenecks or areas for optimization within the existing workflows.
- Create a comprehensive inventory of all Autosys jobs, including their associated scripts, commands, and input/output files.

b) DAG Development:

- Translate the logic and dependencies of Autosys workflows into equivalent Airflow DAGs, ensuring a seamless transition.
- Leverage Airflow's rich set of operators and features to enhance workflow functionality and flexibility.
- Implement error handling and retry mechanisms within the DAGs to ensure robustness and reliability.
- Parameterize DAGs to enable dynamic configuration and adaptability to changing requirements.

c) Testing:

- Rigorously test DAGs in a controlled development environment that mirrors the production setup.
- Validate that DAGs execute as expected, producing the desired outcomes and adhering to defined schedules.
- Simulate various failure scenarios to verify the effectiveness of error handling and recovery mechanisms.
- Conduct performance testing to assess the efficiency and scalability of DAGs under different workloads.

d) Deployment:

- Establish a secure and scalable Cloud Composer environment that aligns with organizational policies and standards.
- Automate the deployment process to minimize manual intervention and reduce the risk of errors.
- Implement version control and configuration management practices to track changes and ensure reproducibility.
- Configure appropriate access controls and permissions to safeguard sensitive data and resources.

e) Monitoring:

- Implement comprehensive monitoring and logging to track workflow execution, identify potential issues, and gain insights into system performance.
- Set up alerts and notifications to proactively respond to anomalies and prevent disruptions.
- Leverage Cloud Composer's built-in monitoring capabilities and integrate with external monitoring tools as needed.
- Regularly review logs and metrics to identify trends, optimize workflows, and ensure operational efficiency.

5. Challenges and Considerations

While the migration to Cloud Composer offers significant benefits, it also presents challenges:

- **Learning Curve:** Transitioning to Airflow and DAGs requires learning new concepts and tools.
- **Code Conversion:** Existing Autosys scripts may need to be rewritten for Airflow.

- **Dependency Management:** Managing dependencies in Airflow can be complex.

6. Conclusions

Migrating from legacy orchestration tools like Autosys to modern platforms like Google Cloud Composer offers substantial advantages in terms of scalability, flexibility, and cloud integration. By adopting Airflow's DAG - based approach and leveraging Cloud Composer's managed service, organizations can streamline their data orchestration processes and enhance their data engineering capabilities.

Acknowledgment

Causal Productions wishes to acknowledge Michael Shell and other contributors for developing and maintaining the IEEE LaTeX style files which have been used in the preparation of this template. To see the list of contributors, please refer to the top of file IEEETran.cls in the IEEE LaTeX distribution.

References

- [1] Apache Airflow Documentation. [Online]. Available: <https://airflow.apache.org/>
- [2] Google Cloud Composer Documentation. [Online]. Available: <https://cloud.google.com/composer/docs>
- [3] <https://github.com/apache/airflow>
- [4] Data Pipelines with Apache Airflow Book by Bas Harenslak and Julian de Ruiter