

Automation of Data Pipeline Creation in Apache Airflow

Prakarsh Anand¹, Anuj Jain², Suma B³

¹Student, R V College of Engineering, Bengaluru, India

²Employee, Bounceshare Pvt Ltd, Bengaluru, India

³Professor, R V College of Engineering, Bengaluru, India

Abstract - A data pipeline is a series of data processing steps. Data pipelines consist of three key elements: a source, a processing step or steps, and a destination. They are scheduled at regular intervals using job executors. They are organised as a directed acyclic graph(DAG) containing a starting point and an endpoint in the DAG. When tasks are executed in sequential, parallel or any desired manner it gives out a meaningful outcome which is used for data analysis or data sciences. Today, data pipelines are being created using programming python scripts from scratch. Airflow provides a very useful user interface for monitoring of execution of the data pipelines as per schedule and requirement.

The paper proposes an end to end UI solution from data pipeline creation to monitoring. It proposes to provide a state of art user experience. This allows non-technical users to leverage the advantages of data orchestration tools. The paper uses Airflow, an open-source data orchestration tool and a python module named "DAG-factory" to implement the proposed solution. Approximately 65 % of all the users of Airflow are data engineers and 97% of the data pipelines executed on Airflow are used for Data Processing(ETL). Therefore, the paper focuses on data processing users as its primary users to provide the best user experience.

Key Words: ETL, Airflow, data engineering, DAG, DAG - factory, data-pipelines

1. INTRODUCTION

To deal with scheduled job execution many data pipeline orchestration tools like Airflow have been developed. Their primary work is to create robust task inter-linkages in a data-pipeline and perform scheduled runs at a regular interval. Out of several orchestration tools, Airflow provides a lot of monitoring and customisation functionalities as GUI.

Quoting Apache Airflow documentation - "Airflow is a platform created by community to programmatically author, schedule and monitor workflows". Airflow provides features to monitor workflow progress and execute DAGs. It allows visualisations to check how a DAG execution time varies with an increase in data to be processed, to check if a task execution can be performed in parallel without any data dependency. It allows back-filling of data through scheduled DAG execution using DAG execution date and not the start date of DAG. These features have undoubtedly made it massively useful worldwide.

1.1 Existing Systems

Airflow became Top-Level Apache Software Foundation project in January 2019 and has been used by a lot of companies catering their need to perform different data orchestration easily and efficiently. Airflow also supports a lot of different formats of input allowing it to be used universally to process any kind of data stored or task to be executed. Creation of tasks needs an operator in general. Airflow supports different kinds of operators, hooks or sensors where each of them can be used to create a task to be executed in a workflow. Different data pipelines are created using the customised cross-linking of different tasks which may run in parallel or sequence as per specified task dependencies created to serve the purpose of DAG creation in a meaningful manner. Airflow also supports python operator which opens a whole lot of possibility of tasks to be created ranging from machine learning module execution to simple control statement execution. The usefulness of python operator in Airflow environment is because any data pipeline created is a python based configuration file which is read by Airflow's python interpreters. Therefore, a python based configuration file can be used for any DAG creation.

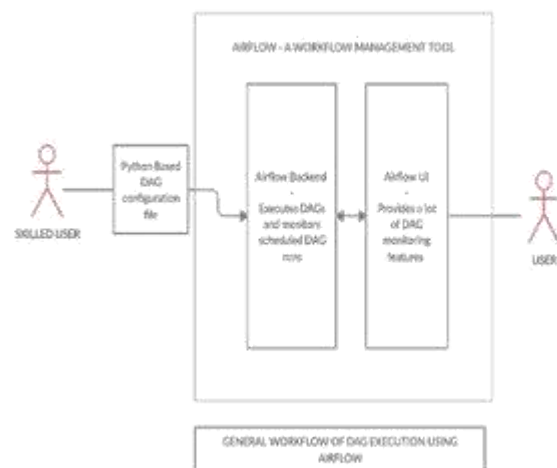


Figure 1 General Workflow of Airflow

Airflow's execution environment and the standard steps to process DAG creation and execution comprise the whole framework used to perform data orchestration. It has been widely in use by a lot of organisations but with a very

obvious lagging step of programmatical data pipeline creation. It leads to different complications like

1. One needs to know python programming to relish the full potential of the variation of data pipelines that can be created and executed.
2. The limited set of operators, sensors or hooks in contrast to rich functionalities provided by python has led to a lot of redundant and non-standardized codes across organisations. They generally are observed as different combinations of same task structures driven by mere syntax usage of operators defined by Airflow.
3. Creation of different ways of solving the problem of repetitive code-based data pipeline creation using pre-processing of different configuration files parsing. For example “.ini” extensions or “.yaml” extensions.
4. If DAGs are being created by not looking after the code repeatability issue, it allows errors like no standardisation of codebase at least across organisation. Human error during implementation increases as the dependability on the way an individual gets his/her work done creeps in every time. Machine-based task creation allows code standardisation and removes human intervention reducing the errors of DAG creation to be negligible. But this doesn't solve the code repeatability issue. This impacts the scalability of code and the organisation as a whole.
5. The steps or creativity applied using the tools available to reduce code repeatability by organisations are generally organisation specific which caters to the needs of an organisation and solid documentation is needed to scale the organisation and code base efficiently. It also means any new user of the non-standard tool created by the organisation needs to learn it from scratch using the organisation resources and it can't be used as a standard tool which doesn't serve the purpose as per the analysis of this paper.

1.2 Scope and Motivation

With every addition of module in the python library, DAG creation becomes more efficient and its userbase scales massively increasing the use of Airflow in particular multiple folds. One such python module addition is “ DAG-factory ” which if used in full potential will open DAG creation and execution for non-technical people. This will allow them to work with data orchestration as efficiently as people who know python programming. The proposed methodology leverages the full potential of DAG factory python module which reads a YAML based configuration file and converts it into a python-based configuration file internally. This python-based configuration file is used by the Airflow execution environment to create and execute data pipelines.

Always a YAML based configuration file readability and writability are far more generalised and easier than a python-based configuration file creation. It allows non-python developers, common users of the data pipelining environment or any person who can understand configuration files to understand and write a configuration file which involves very less complexity as compared to python-based configuration files and thus has fewer chances of human error and increase in the user base of airflow as a data orchestration platform. This is the whole purpose of implementation of “DAG factory” module as an open-source module. It serves the purpose of an increase in the user base of open-source data pipeline execution environment Airflow. It also creates a potential extension of the usage of the module to completely change the way DAG creation has been looked and assumed by the world so far. It has given the potential to create an environment which may create the desired DAG on click of a button or a mere command-line execution. The paper describes an approach which may potentially change the world's notion of data pipeline creation from being a skilled person in required technology to merely using the UI to specify what you want to perform using the data pipeline and the resources to be used.

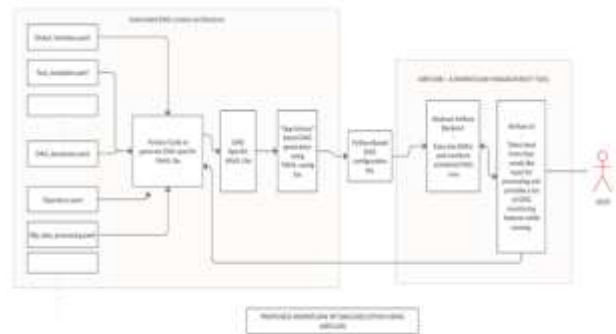


Figure 2 Proposed Airflow Architecture Diagram

2. METHODOLOGY

Airflow defines DAG as a Python script, which represents the DAGs structure (tasks and their dependencies) as code in their official documentation. This has been interpreted as skilled users may generally create a python script to create data-pipelines. Airflow uses a python-based configuration file to interpret and create data pipelines accordingly to be scheduled to execute at regular interval. These python-based configuration files can be created through coding the python configuration file manually or using the DAG-factory module which parses a configuration file and converts it into a python based configuration file to be interpreted for DAG creation. The DAG-factory module only needs the source path of the configuration file to fetch the file content and parses it to create a python based configuration file. A configuration file can contain many DAG definitions where each section in the configuration file may be fetched as a different DAG definition. Here, the configuration file is the standard content

- clubs the required modules in specific order to create a completely defined configuration file.
- 4. The configuration file is processed by “DAG-factory” module to convert it into a python based configuration file.
- 5. The python script is an executable used by Airflow interpreter to execute DAGs.

The whole process of DAG creation requires user to provide the information about what to perform using the DAG to create, knowing the resources and links to be used to accomplish use of DAG. This approach may potentially change the world’s notion of data pipeline creation needs a skilled python programmer.

5. RESULTS

	Existing Implementation	Proposed Implementation
GUI based Automated DAG creation	No	Yes
Prior knowledge of any programming language	No	Yes
Knowledge of Airflow primitives	No	Yes
Duplicacy of Code	Yes	No
Human Error Possibility	Yes	No
Reduces effort required to create customised DAGs	No	Yes
Reduces time taken to create multiple DAGs	No	Yes
Basic Understanding of Airflow required to create data pipelines	Yes	Yes
Transformation from low level programming to abstracted code reusability through templates	No	Yes
Abstraction — Interoperability and platform independence	No	Yes
YAML configuration file used than a python based configuration file	No	Yes
Increases flexibility of DAG generation types	No	Yes
Modularity - Read once write multiple times	No	Yes

Figure 5 Comparison of proposed vs existing implementations

The criteria chosen to compare the 2 implementations define the advancements possible in Airflow user experience. The system has been effectively tested in a rapidly growing software company, and showed satisfactory results. The implementation correctly ensures the reliability and robustness of the system testing process. The comparisons were done between the existing python-based DAG configuration files in the company and their equivalent conversions into the template-based implementation. The equivalent template-based implementation where used to automatically create python based DAG configuration files on issue of desired arguments to the python code generator .

6. CONCLUSION

Template-based automation of DAG creation is an experimental project work taken in order to prove that it’s efficient and effective implementation will definitely cause a shift in mindset from the way, data pipelines are generally written in today’s world to an automated DAG creation platform where specifying the expectation from the data pipeline to be created is enough .This template based implementation works quite well because of the simple reason that a task is the fundamental unit of any data pipeline and any task is created using an operator definition from a set of limited number of operators. Therefore, a proper definition of usage of every operator in a manner such that the definitions can be used by any task implementation is enough. Data pipeline implementation makes the proposed template-based data pipeline creation work. It also allows a very robust implementational abstraction of data pipeline creation from the end user who doesn’t need to know the implementational part ensuring the benefits of abstraction to be leveraged later.

7. FUTURE WORK

It can be started as an open source project to create an ecosystem where different and new tasks created, Different and new DAGs created can be pushed in a general config file and they can be used as template to generate similar DAGs multiple times all over the world on a single search on the file with specific parameters assigned. Further, the future work for the project would be to make it robust and enough user friendly so that it can mainstreamed by the open source community to use it for data pipeline creation rather than writing repeated codes. Making the User interface as easy and user friendly as possible so that people can specify search through the DAG templates available, create a non-existing DAG template through clicks of buttons, select the template for DAG creation and specify parameter changes if required to create data pipeline. Link the DAG creation to direct execution of data pipeline on AIRFLOW. Thus completing the user experience of creation and execution of data pipelines through clicks of button than coding all of it together.

8. ACKNOWLEDGEMENT

We would like to acknowledge the support provided by employees of Data Engineering Team, Bounceshare, Bengaluru, India and teachers of Department of Computer Science & Engineering & Department of Biotechnology, RV

College of Engineering, Bengaluru, India through their assistance during the research work.

M., et al.: Pegasus, a workflow management system for science automation. *Future Generation Computer Systems* 46, 17–35 (2015)

9. REFERENCES

1. H. Bal, M. Haines, "Approaches for Integrating Task and Data Parallelism", *IEEE Concurrency*, vol. 6, no. 3, pp. 74-84, July-Aug. 1998.
2. O. Delannoy, N. Emad, S. Petiton. Workflow Global Computing with YML. In: *Proceedings of the 7th IEEE/ACM International Conference on Grid Computing, GRID'06*. ACM Press, New York, 2006. 25-32.
3. F. P. Guimarães and A. C. M. Melo, "User-Defined Adaptive Fault-Tolerant Execution of Workflows in the Grid," in *Proceedings of the IEEE CIT*, Sep 2011, pp. 356-362.
4. L. Li, Z. Miao, L. Yuqing, Q. Liangjuan, "A Survey on Workflow Management and Scheduling in Cloud Computing", *Cluster Cloud and Grid Computing (CCGrid) 2014 14th IEEE/ACM International Symposium on*, pp. 837-846, 2014.
5. M. Kotliar et al., "CWL-Airflow: a lightweight pipeline manager supporting common workflow language", *bioRxiv*, 2018.
6. M. Beauchemin, (2014) Apache Airflow Project.
7. A. Barker, J. Van Hemert, "Scientific workflow: a survey and research directions", *International Conference on Parallel Processing and Applied Mathematics*, pp. 746-753, 2007.
8. E. Deelman, T. Peterka et al., "The future of scientific workflows", *The International Journal of High Performance Computing Applications*, 2018.
9. R. Sumbaly, J. Kreps, S. Shah, "The big data ecosystem at linkedin", *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, pp. 1125-1134, 2013.
10. W. Hummer, V. Muthusamy, T. Rausch, P. Dube, K. El Maghraoui, "Modelops: Cloud-based lifecycle management for reliable and trusted ai", *2019 IEEE International Conference on Cloud Engineering (IC2E'19)*, Jun 2019.
11. G. Alonso, D. Agrawal, A. El Abbadi, C. Mohan, "Functionality and Limitations of Current Workflow Management Systems", *IEEE Expert*, vol. 1, no. 9, 1997.
12. M. Rosemann, "Evaluation of Workflow Management Systems - a Meta Model Approach", In: *Proceedings of the 2nd International CAiSE/IFIP 8.1. Workshop 'Evaluation of Modeling Methods in Systems Analysis and Design*, 1997.
13. M. Berger et al., *An Evaluation of Workflow Management System*, Austria: Institute for Applied Computer Science and Information Systems, University of Vienna, 1997.
14. Deelman, E., Vahi, K., Juve, G., Rynge, M., Callaghan, S., Maechling, P.J., Mayani, R. Chen, W., Da Silva, R.F., Livny,