

Apache Airflow – Conceptual Write-Up

What is Apache Airflow and how does it work?

Apache Airflow is an open-source tool used for creating, scheduling, and monitoring workflows. It was originally developed at Airbnb and is now widely used in the data engineering field. The main idea of Airflow is that workflows can be written as Directed Acyclic Graphs (DAGs) using Python code. Each workflow is broken down into tasks, and every task is represented by an operator such as a Python function, a shell script, or a SQL query. The scheduler then takes care of running these tasks in the correct order, while the executor is responsible for actually running them, either on the same machine or across distributed systems. Airflow also provides a very user-friendly web interface, where we can monitor running tasks, check logs, retry failed ones, and see the progress of pipelines.

Where does Airflow fit in modern data engineering workflows?

In the current world of data engineering, there are many tools for data storage, processing, and analytics. Airflow plays the role of an orchestration tool, which means it does not process or move data itself but coordinates other tools to do so. For example, Airflow can schedule a job to extract data from an API, trigger a Spark job to transform it, and finally load it into a data warehouse like Snowflake or BigQuery. This makes it very useful for building ETL and ELT pipelines, machine learning workflows, and reporting systems. Since companies rely on multiple technologies in their data ecosystem, Airflow acts as the central controller that ensures everything runs in the right sequence and on time.

How is Airflow different from traditional schedulers or other tools?

Traditional schedulers like cron are very limited. They can run tasks at specific times but cannot handle dependencies, retries, or monitoring. Airflow, on the other hand, is built exactly for these needs. It knows which task should run first, what depends on what, and it can retry tasks automatically if something fails. Airflow is more mature,

widely adopted, and has a very strong ecosystem of plugins and integrations. This makes Airflow a preferred choice in many enterprises.

What are the key components and how do they interact?

Airflow is made up of several important components.

DAGs are the pipelines themselves, written in Python.

Operators are the building blocks of tasks inside a DAG.

Scheduler looks at DAGs and decides which tasks to run and when.

Executor then runs those tasks either locally or in distributed environments.

All the information about past runs, task status, and logs are stored in the **metadata database**.

Finally, the **web UI** provides a friendly way to visualize DAGs, check status, and control workflows. All these parts work together to provide a complete orchestration system.

Where is Airflow useful in real-world scenarios?

Airflow has many use cases in real-world enterprises. In the finance sector, it can schedule pipelines for fraud detection and daily reporting. In e-commerce, it can update recommendation models, sync inventory, and prepare sales reports. In healthcare, it can automate the ingestion of patient data into analytics systems. SaaS companies often use Airflow to process user activity logs, run churn prediction models, and feed data lakes. Basically, whenever there are multiple steps in a data process that depend on each other, Airflow is very helpful in automating and managing them.