

# SciDAP

Airflow and CWL-powered bioinformatics platform

Nicholas Luckey  
Michael Kotliar

Human readable workflow representation

Workflow provenance support

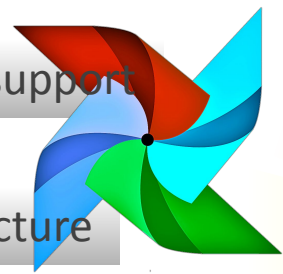
Modular workflow structure

**Apache Airflow**

Interoperability and portability

independent of vendor  
task-driven workflow management  
system developed by Airbnb

Rich ecosystem of available workflows  
widely supported by scientific community



**CWL-Airflow**

python package that adds support  
for CWL to the Apache Airflow

**Common Workflow Language**

open standard  
workflows are  
makes them easy

Arvados Project

Curii

Seven Bridges Genomics

Galaxy Project

Apache Taverna

Institut Pasteur

Wellcome Trust Sanger Institute

University of California Santa Cruz

Harvard T.H. Chan School of Public Health

Cincinnati Children's Hospital Medical Center

Broad Institute

University of Melbourne Center for Cancer Research

Netherlands eScience Center

Agave Platform

CyVerse

Institute for Systems Biology

ELIXIR Europe

BioExcel

BD2K

EMBL Australia Bioinformatics Resource

IBM Spectrum Computing

DNAnexus

CERN

Inputs:

Outputs:

steps:

alternative step:

second step:

second step:

input file:

output:

output\_file

run: second.cwl

input file: first\_step/output\_file

out:

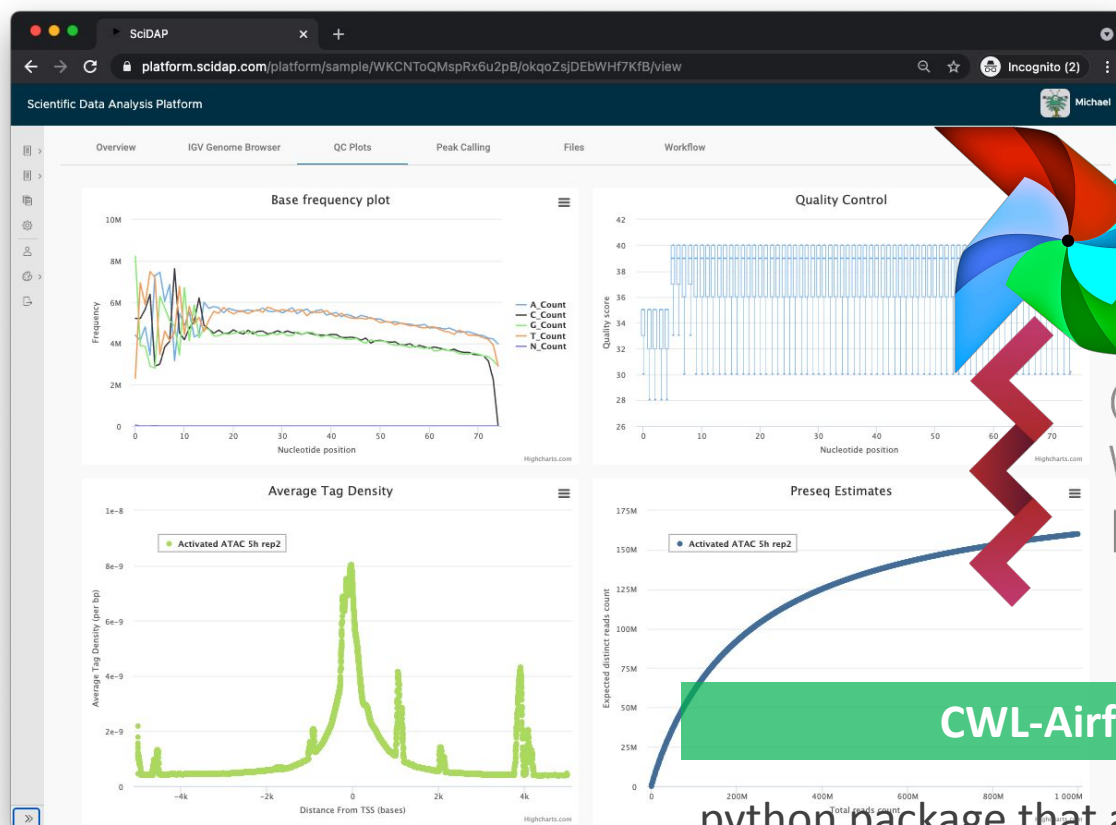
output\_file

Workflow

Research Object

- Findability
- Accessibility
- Interoperability
- Reuse

# SciDAP

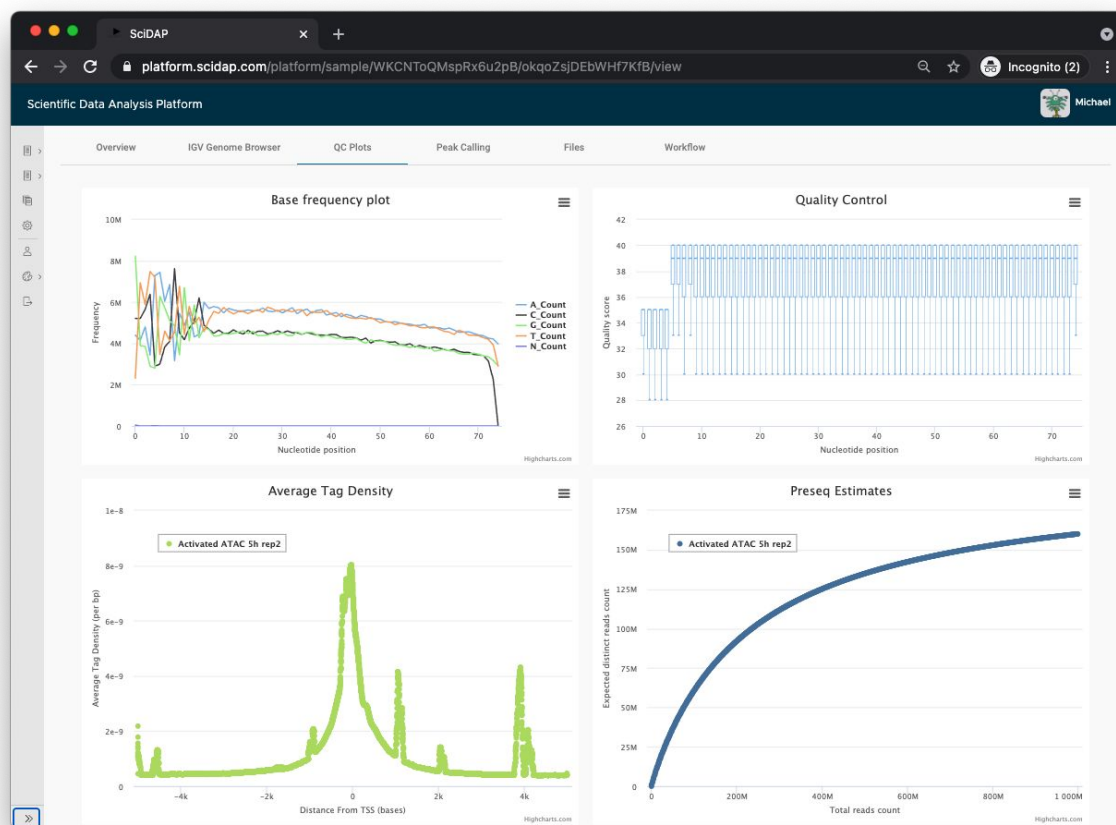


COMMON  
WORKFLOW  
LANGUAGE

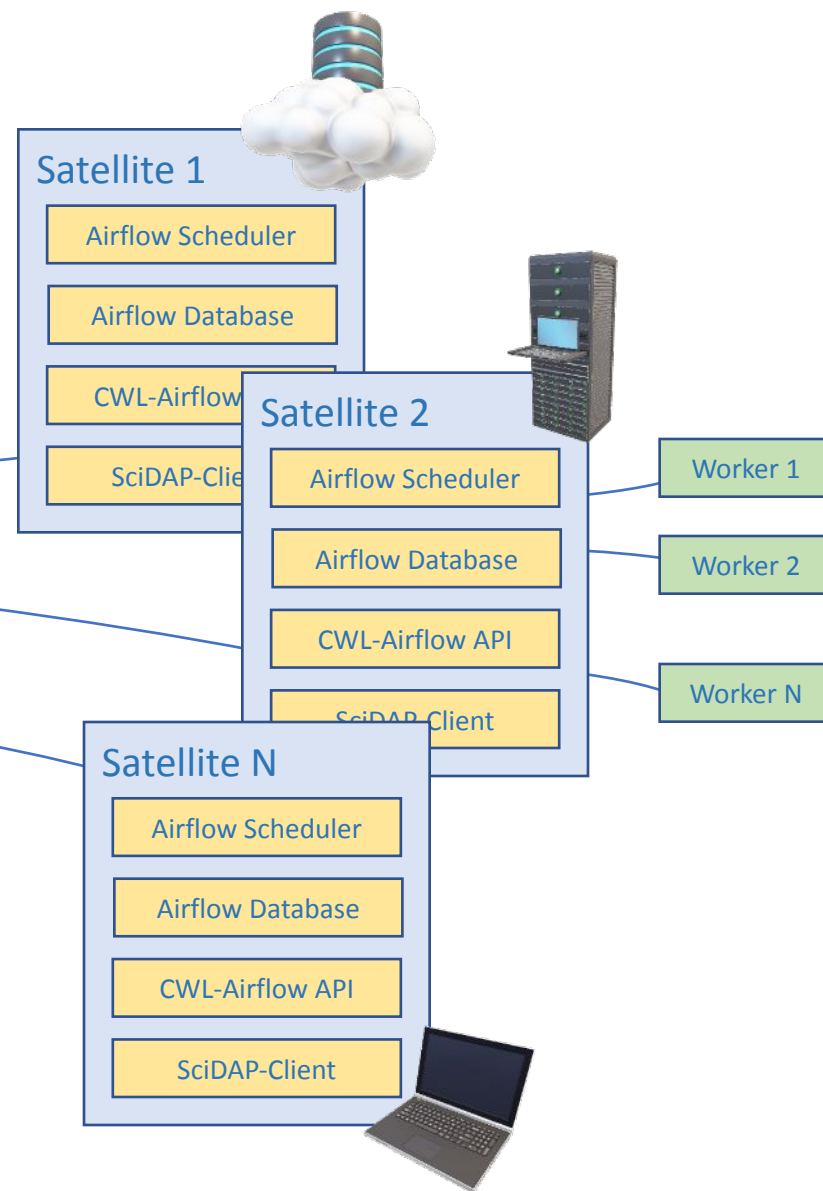
CWL-Airflow

python package that adds support  
for CWL to the Apache Airflow  
user-friendly scientific data analysis platform

# SciDAP

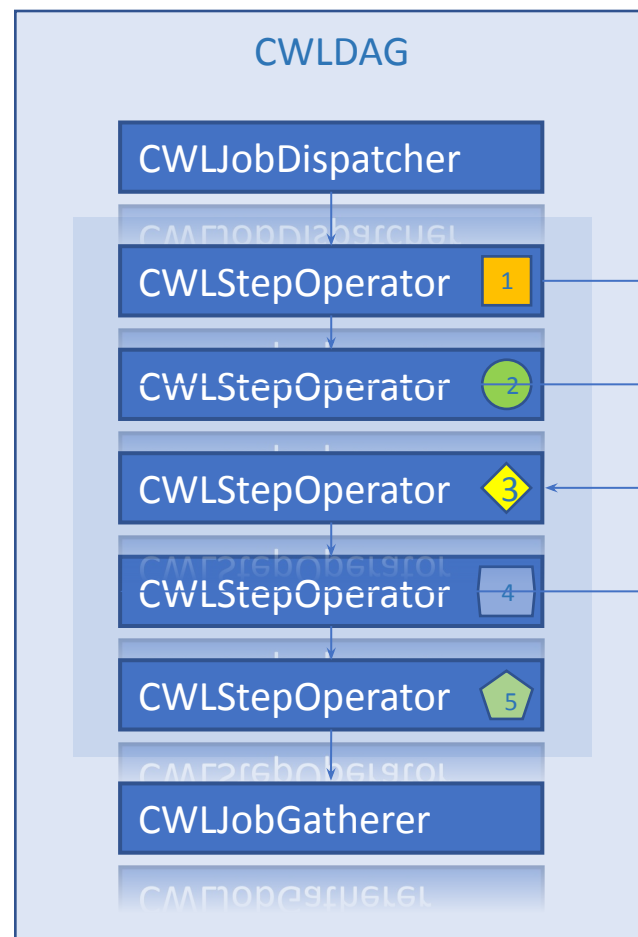
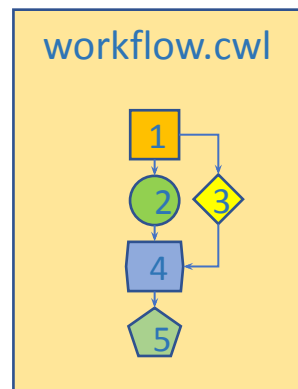


user-friendly scientific data analysis platform



# Challenge 1

## Converting CWL workflow to Airflow DAG



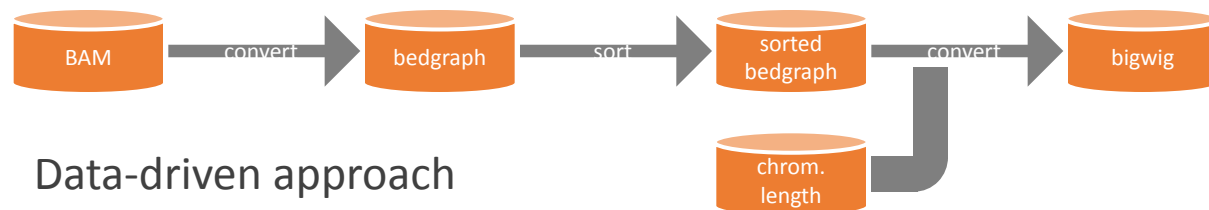
### my\_dag.py

```
#!/usr/bin/env python3
from cwl_airflow.extensions.cwldag import CWLDAG
dag = CWLDAG(
    workflow='workflow.cwl',
    dag_id='bam-bedgraph-bigwig'
)
```

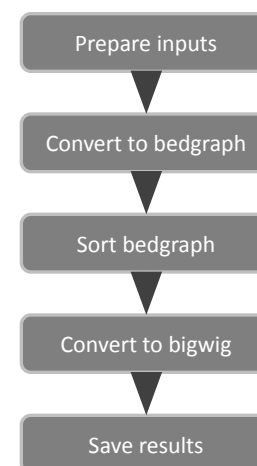
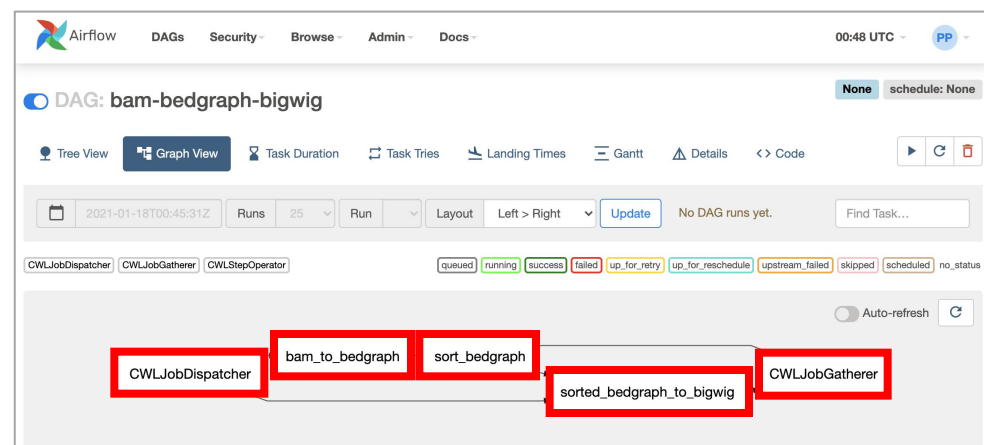
### my\_dag\_with\_embedded\_workflow.py

```
#!/usr/bin/env python3
from cwl_airflow.extensions.cwldag import CWLDAG
dag = CWLDAG(
    workflow='H4slAAAn3/F8/8 ... TikAAA==',
    dag_id='bam-bedgraph-bigwig'
)
```

## Challenge 2 CWL is a data-driven workflow standard



Data-driven approach



Combine tasks based on the CWL workflow step inputs and outputs

Define a mechanism for transferring data between tasks

Traditional for Airflow DAG-based task-driven approach

Task starts its execution if all its predecessors have successfully completed

Task is not supposed to exchange data with other task

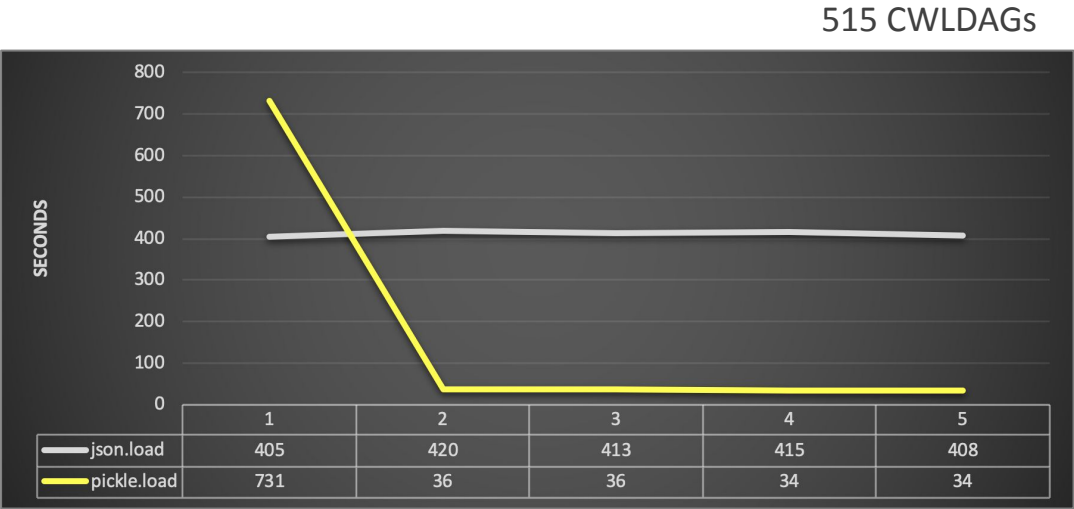
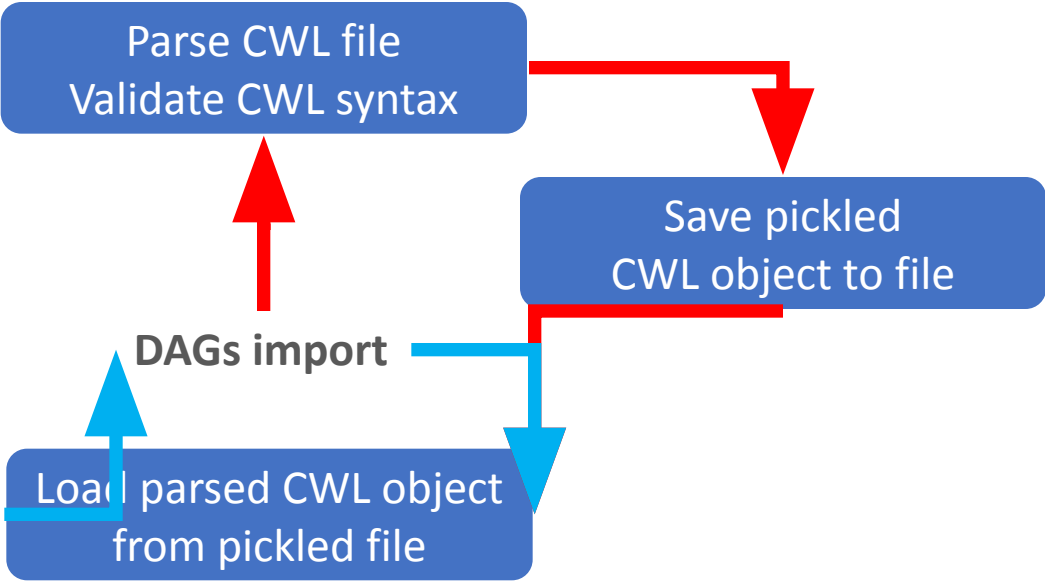
# Challenge 3 CWL slows down Airflow DAG import

- 1. Parse CWL file
- 2. Validate CWL syntax**
- 3. Create DAG

my\_dag.py

```
#!/usr/bin/env python3
from cwl_airflow.extensions.cwldag import CWLDAG
dag = CWLDAG(
    workflow='workflow.cwl',
    dag_id='bam-bedgraph-bigwig'
)
```

How long before timing out a python file import?  
**dagbag\_import\_timeout = 30.0**



## Advantages:

Cross-vendor portability

Shallow learning curve for newcomers

Designed with FAIRness in mind

Rich ecosystem of available and tested workflows

Workflow visualization tools  
for easy pipeline building and viewing

Broad list of participating organizations

## Disadvantages:

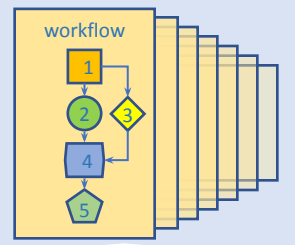
Limited functionality implied by  
using CWL specification

Extra level of complexity, thus  
more places to look for errors



# CWL-Airflow and its role in scientific research

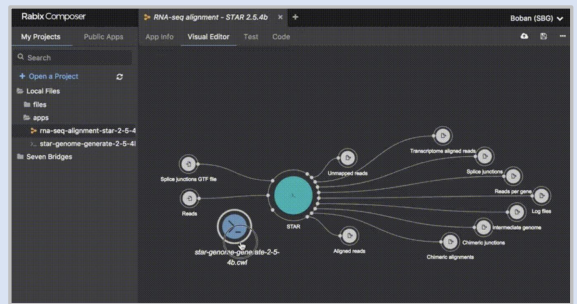
## CWL workflows



## Innovations



## Visual workflow editors



## Workflow management systems

SevenBridges



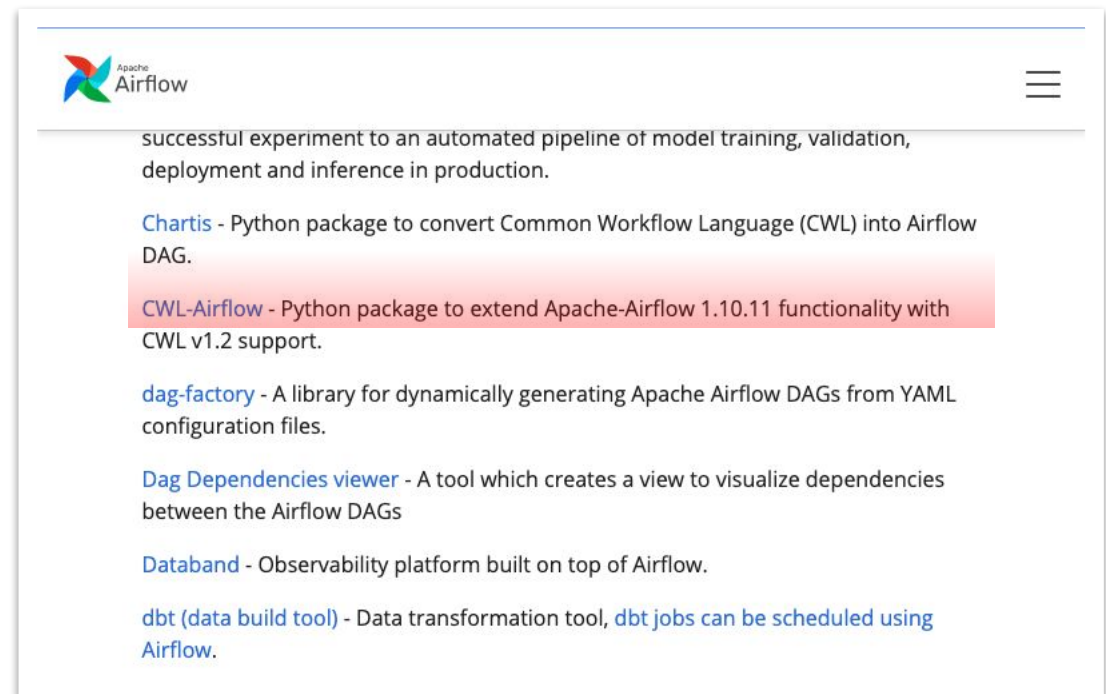
Arvados



reana

- Complete CWL specification support
  - Conditional step execution
  - Scattered step execution
  - Abstract operations support
- Support other workflow specifications
  - WDL
  - Nextflow
  - Snakemake
- Complete Cloud support

## Airflow Ecosystem



The screenshot shows the Apache Airflow Ecosystem page. At the top left is the Apache Airflow logo, and at the top right is a hamburger menu icon. Below the header, there is a paragraph describing a successful experiment in automating model training, validation, deployment, and inference. This is followed by a list of ecosystem tools, each with a brief description. The 'CWL-Airflow' entry is highlighted with a red background. The tools listed are: Chartis, CWL-Airflow, dag-factory, Dag Dependencies viewer, Databand, and dbt (data build tool).

successful experiment to an automated pipeline of model training, validation, deployment and inference in production.

[Chartis](#) - Python package to convert Common Workflow Language (CWL) into Airflow DAG.

[CWL-Airflow](#) - Python package to extend Apache-Airflow 1.10.11 functionality with CWL v1.2 support.

[dag-factory](#) - A library for dynamically generating Apache Airflow DAGs from YAML configuration files.

[Dag Dependencies viewer](#) - A tool which creates a view to visualize dependencies between the Airflow DAGs

[Databand](#) - Observability platform built on top of Airflow.

[dbt \(data build tool\)](#) - Data transformation tool, [dbt jobs can be scheduled using Airflow](#).