



Apache Airflow

@ Jagex – Anum Sheraz

# Managing DAGs at scale

Deployment, versioning & package management



# Agenda

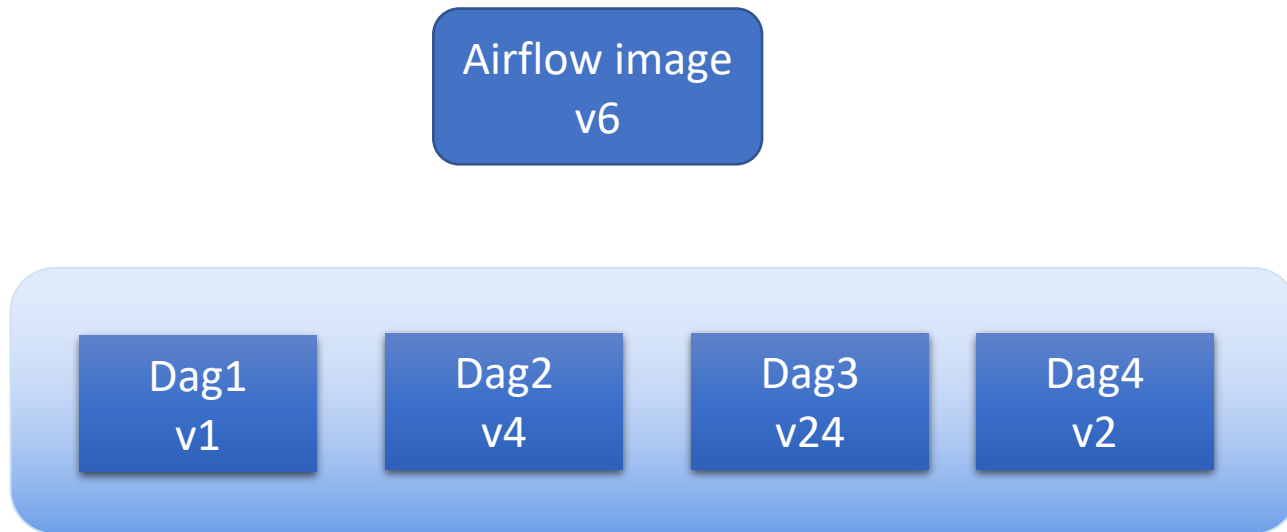
---

- Keep track of airflow state
- keep track of DAGs states
- Revert to previous state
- Manging DAGs code base
- Releasing next airflow version
- Automate airflow deployment process
- How to manage DAGs package dependencies



# Airflow/Dags States

---



# Airflow Repository

---

## Contains

- Dags code
- Dockerfile
  - Inherit from airflow base image
  - install custom\_dags/ into airflow dags/ directory
  - Install additional py packages
  - generate new custom image version x

custom\_dags/

- dag\_1 (v1)
- dag\_2 (v3)
- dag\_3 (v8)

Dockerfile

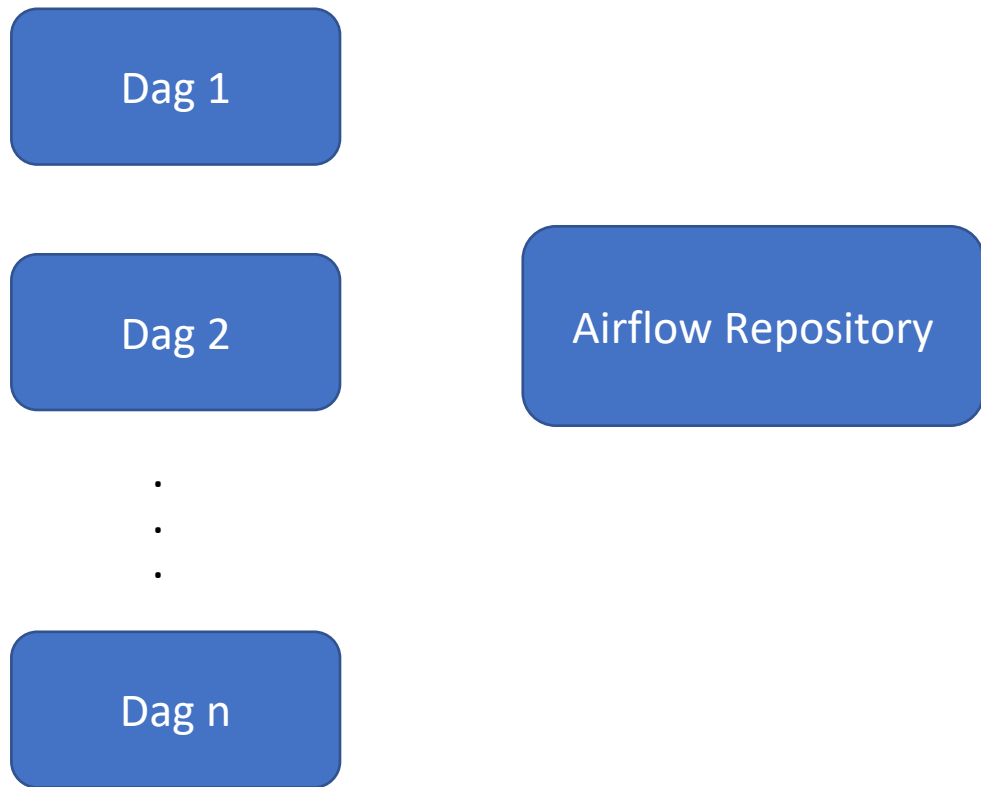
## Problems:

- All dags code in single repository
- hard to manage



# Split Dags repositories

---



# Split Dags repositories

---

Dag 1

Dag 2

⋮

Dag n

Airflow Repository

Dag\_1 v?

Dag\_2 v?

Dag\_3 v?

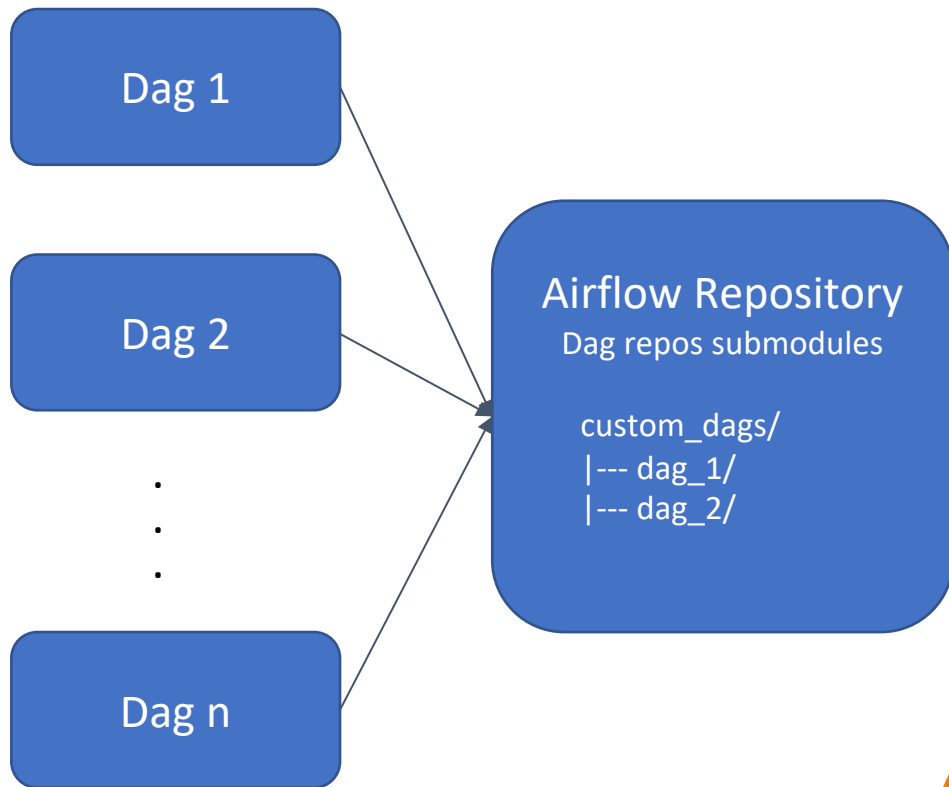


# Dags | Git submodules

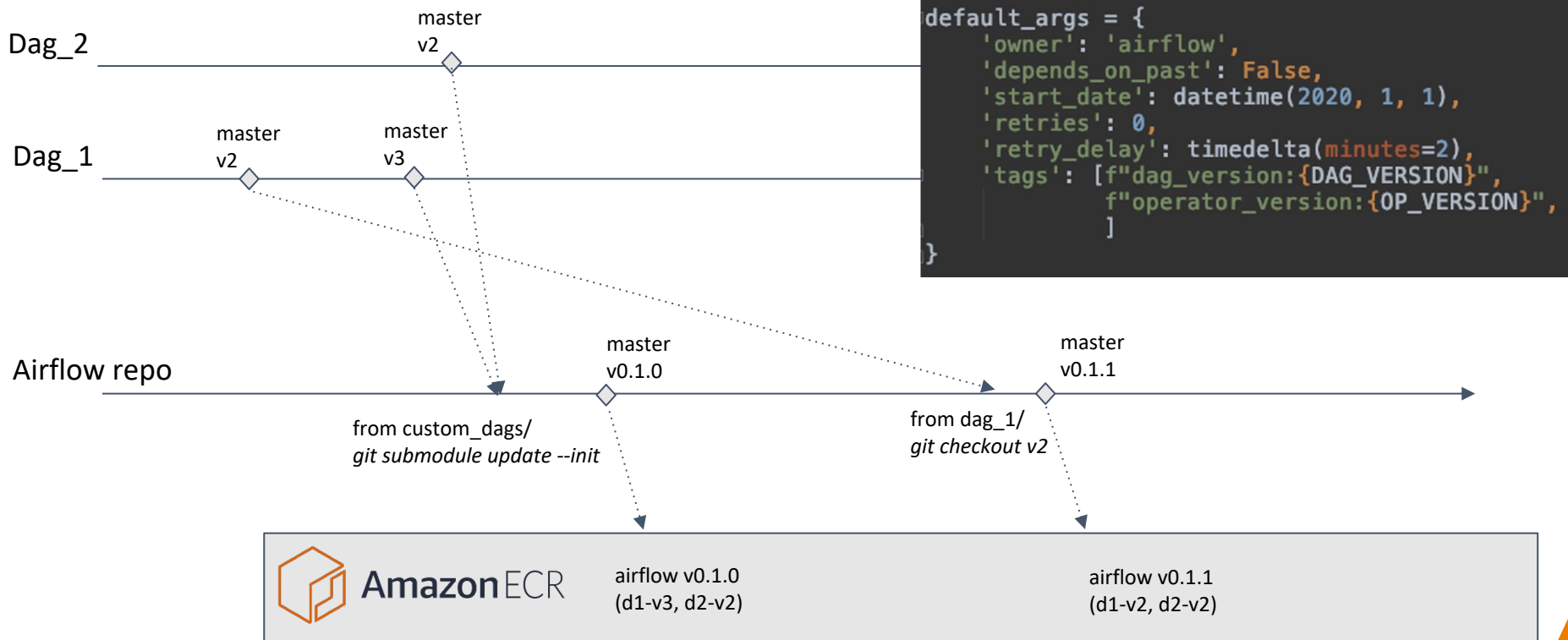
---

## Airflow Repository

- `Git submodule add https://<dag_1_repo>`
- `Git submodule add https://<dag_2_repo>`
- `Git submodule add https://<dag_n_repo>`



# Dags | Git submodules



Amazon ECR

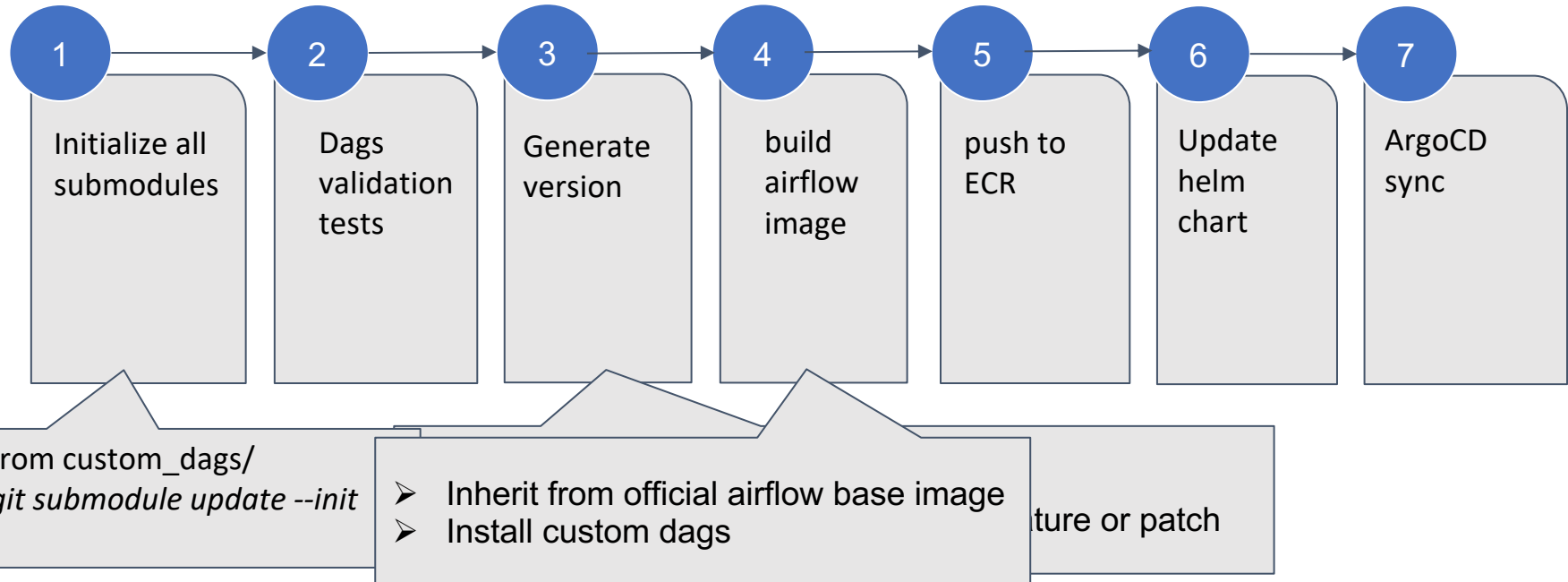
airflow v0.1.0  
(d1-v3, d2-v2)

airflow v0.1.1  
(d1-v2, d2-v2)

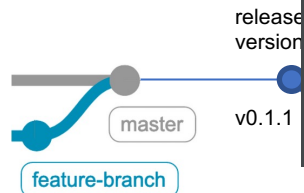


# CICD Pipeline Steps

Pipeline triggers on Airflow repository master branch

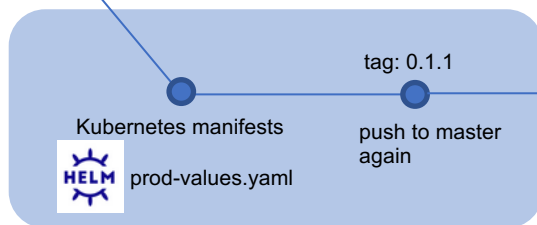


# Airflow De



```
airflow-dags ~/jagex/code/airflow-dags
├── .dev
├── custom-dags
├── deployment
├── templates
├── tests
├── Chart.yaml
├── dev-values.yaml
├── prod-values.yaml
├── README.md
├── values.yaml
├── images
├── plugins
├── tests
├── .bumpversion.cfg
├── .gitignore
├── .gitmodules
├── bitbucket-pipelines.yml
├── check-commit-standards.sh
├── README.md
├── version.py
└── xcom_sidecar.py

scheduler_parsing_processes: 4
statsd:
  enabled: true
  image: prom/statsd-exporter
  image_pull_policy: IfNotPresent
  image_tag: v0.21.0
  ingest_port: 9125
  prefix: airflow-de
  scrape_port: 9102
dags_image:
  pull_policy: IfNotPresent
  repository: [REDACTED].dkr.ecr.eu-west-1.amazonaws.com/airflow_dags
  tag: 0.1.1
secret: jagex-ads-airflow-prod-configuration-credentials
```



commit from cicd pipeline

CD sync



ArgoCD



PROD

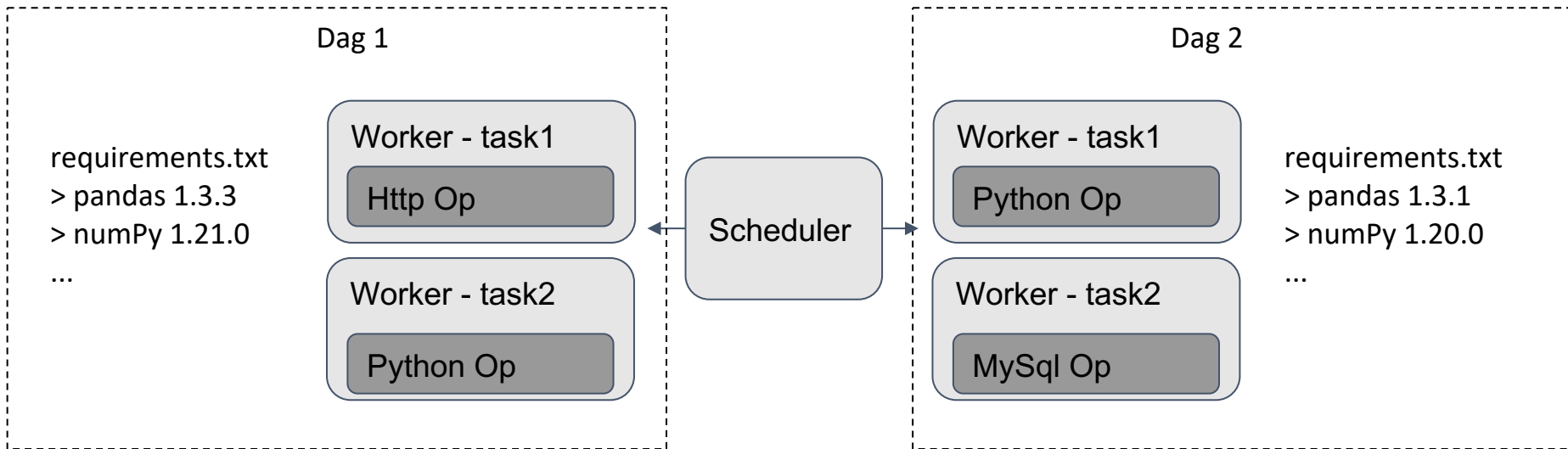


DEV

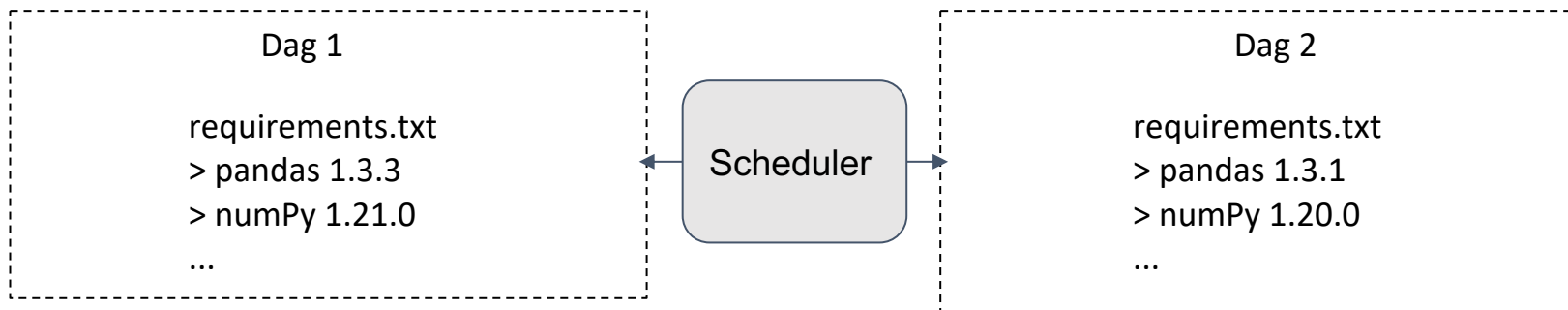


# Package Dependency management

---



# Package Dependency management | Problems



- Airflow workers also use same airflow image (v0.1.1)
- install packages required by all DAGs into airflow image?
  - Hard to track which package is used by which DAG
  - Hard to cleanup packages while DAG removal/upgrade
- Create virtual environment for each DAG ?
  - ramp-up and ramp-down venvs while DAG starts/stops
  - Slows down the workflow
  - needs extra work to maintain venvs

# Package Dependency management | Solution

---

- Workflow execution: business logics performed inside DAG's tasks
- Workflow management: Dags and tasks triggering, alerts, monitoring etc

Need Isolation between these two layers... idea ?

1. Build docker container to perform business logics
2. Let airflow run those containers

How ?

- Docker Operator
- Kubernetes Pod Operator



# Docker Operator

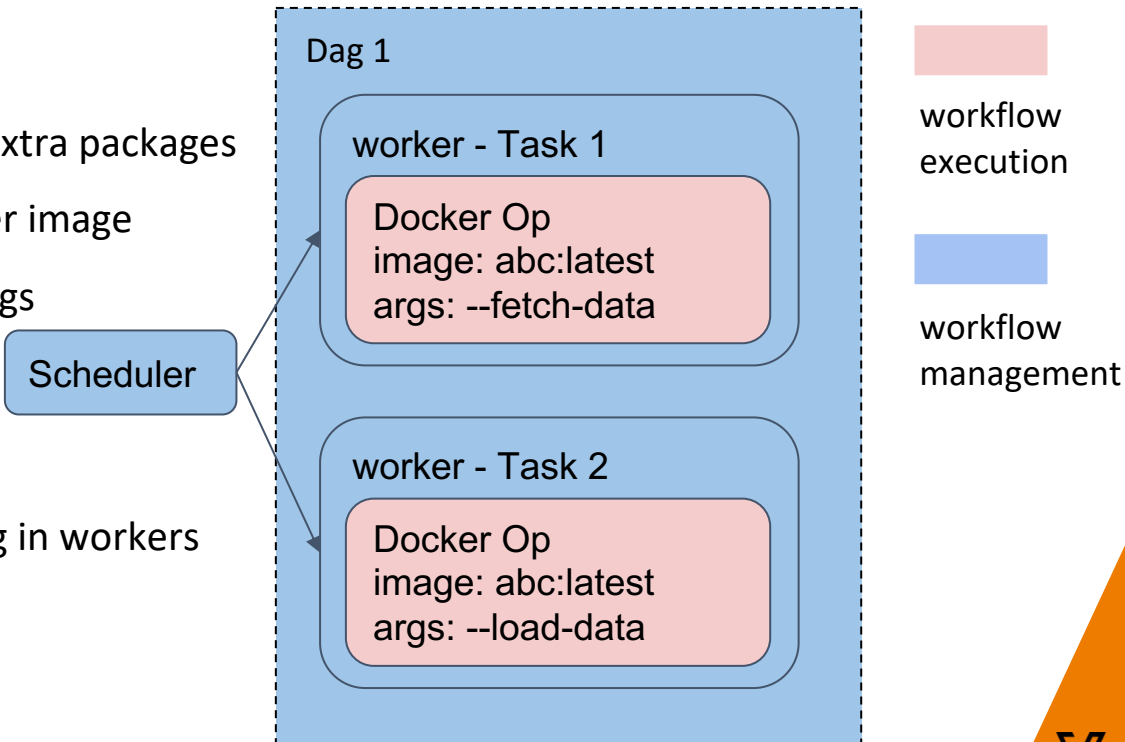
Lets you run docker container within same airflow worker node

## Pros:

- not populating airflow image with extra packages
- Tasks packages installed in container image
- Reuse same image with different args
- Not bound to python anymore

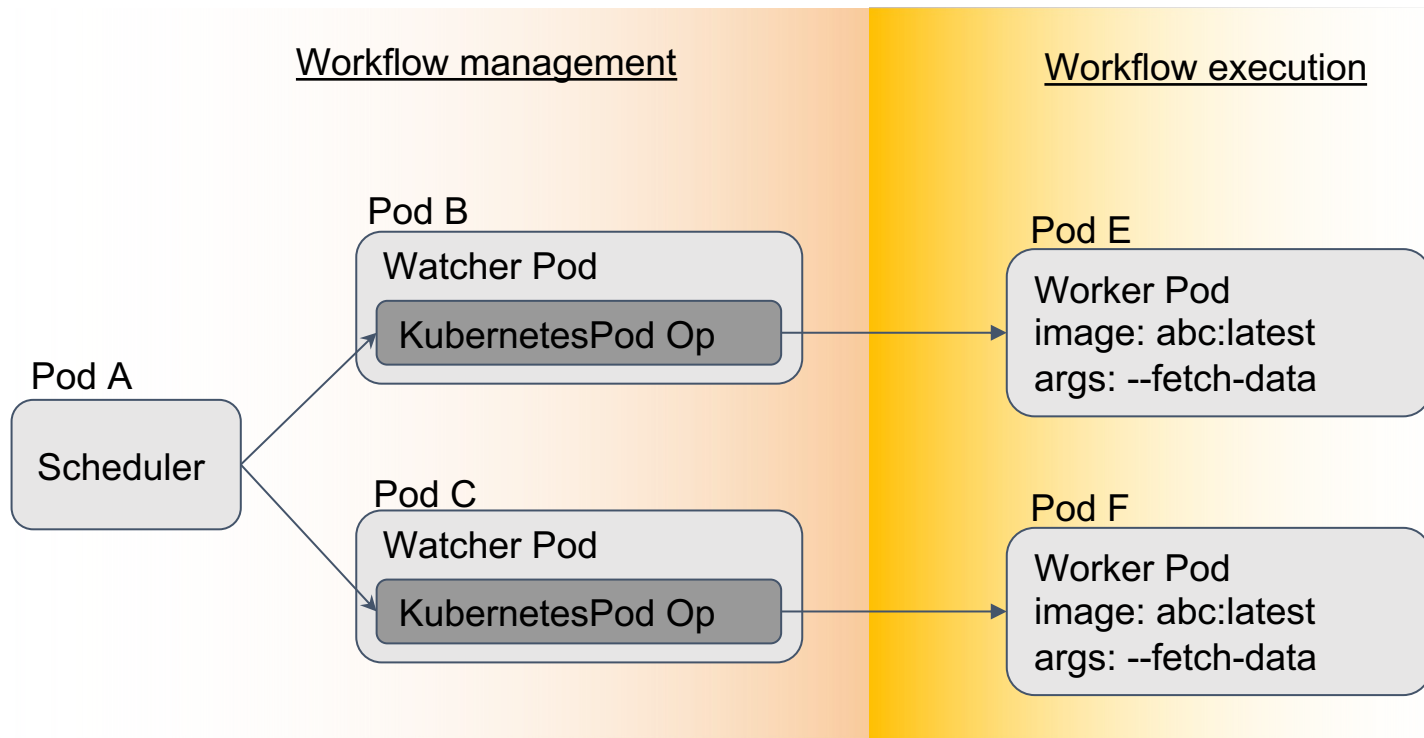
## Cons:

- resource limited - container running in workers



# Kubernetes Pod Operator


Lets you run docker container on dedicated worker node (pod) in kubernetes cluster






# Pros & Cons



---

## Docker Operator

 resource limited - container running inside workers

-  not populating airflow image with extra packages
-  Tasks packages installed in container image
-  Not bound to python anymore
-  Reuse same image with different args
-  Easy to debug

## Kubernetes Pod Operator

-  Kubernetes knowledge required
-  slightly slow - Each task require pair of nodes

