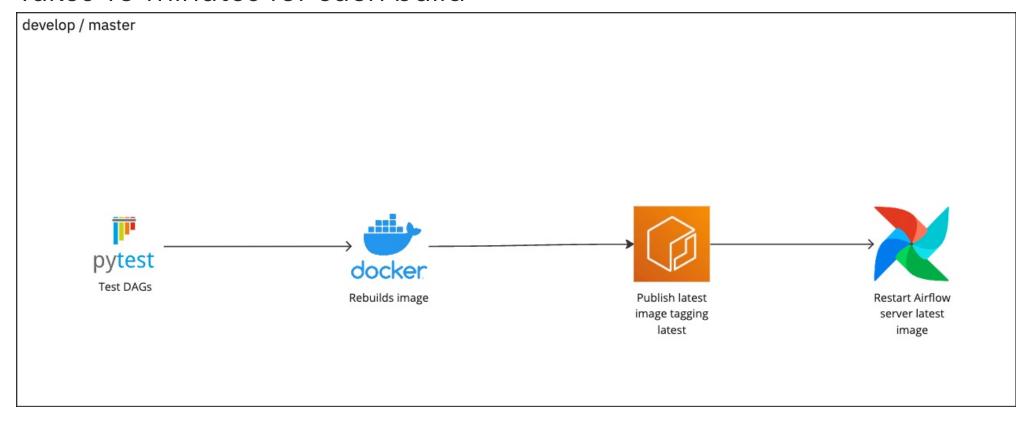


Agenda

- # CI/CD
 - Present
 - Planned
- Future of Airflow
- ? Questions

** Present CI/CD

- Builds a new image for any changes
- Publishes this image ~500 MB each time on updated changes
- Takes 10 minutes for each build

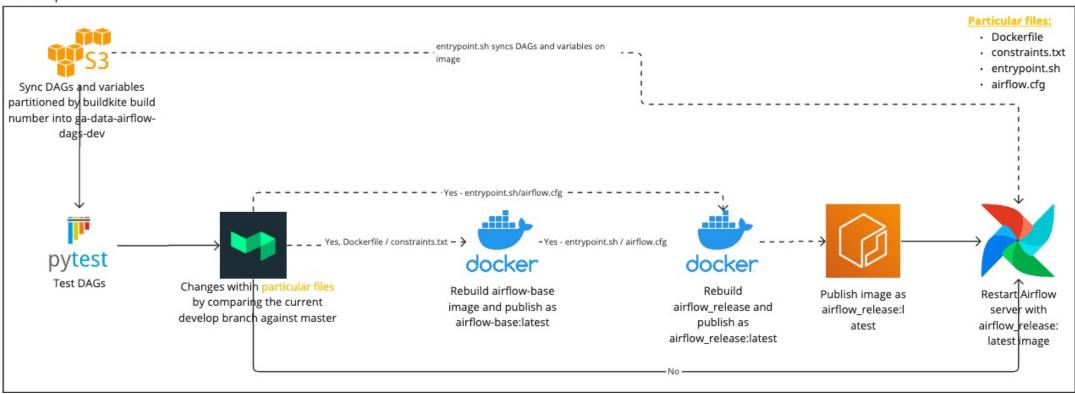


© Planned CI/CD

- Goals:
 - Reduce the time to build by 50%
- Achieved by:
 - i. Sync DAGs / variables into ga-data-airflow-dags-{env} S3 bucket.
 - ii. Splitting the Dockerfile into a multi stage build
 - airflow_base : for requirements
 - airflow_release: for entrypoint.sh / airflow.cfg changes
 - iii. Update entrypoint.sh using aws s3 sync to sync DAGs back
- Other changes:
 - Using SEEK's open source software plugins (seek-oss/docker-ecr-publish)

@ Planned CI/CD - develop

develop



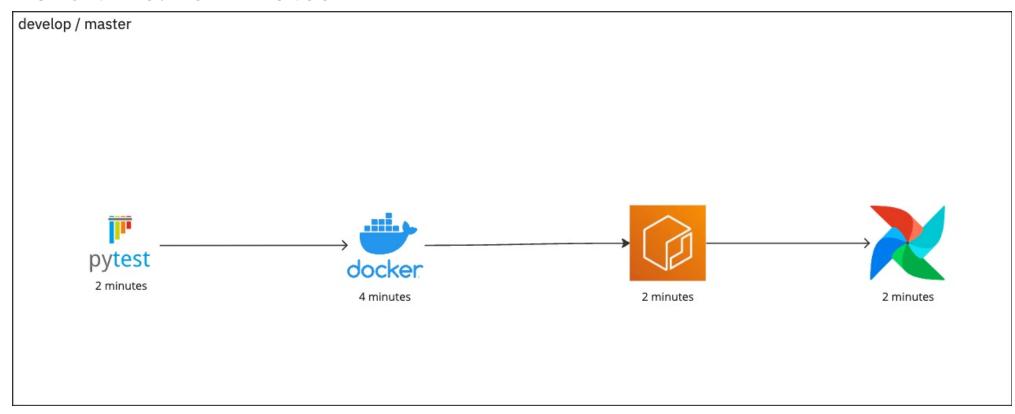
Planned CI/CD - master

Difference between develop and master is that each PR is a squash and merge.

master Particular files: Dockerfile entrypoint.sh syncs DAGs and variables on image · constraints.txt entrypoint.sh · airflow.cfg Sync DAGs and variables partitioned by buildkite build number into ga-data-airflowdags-prod - Yes - entrypoint.sh / airflow.cfg Yes, Dockerfile / constraints.txt - > - Yes - entrypoint.sh / airflow.cfg · -> pytest docker docker Test DAGs Changes within particular files Rebuild airflow-base Rebuild Restart Airflow Publish image as by comparing the current image and publish as airflow_release and airflow release:l server with commit and previous commit airflow-base:latest publish as airflow release: atest airflow release:latest latest image Assumes squash and merge PR

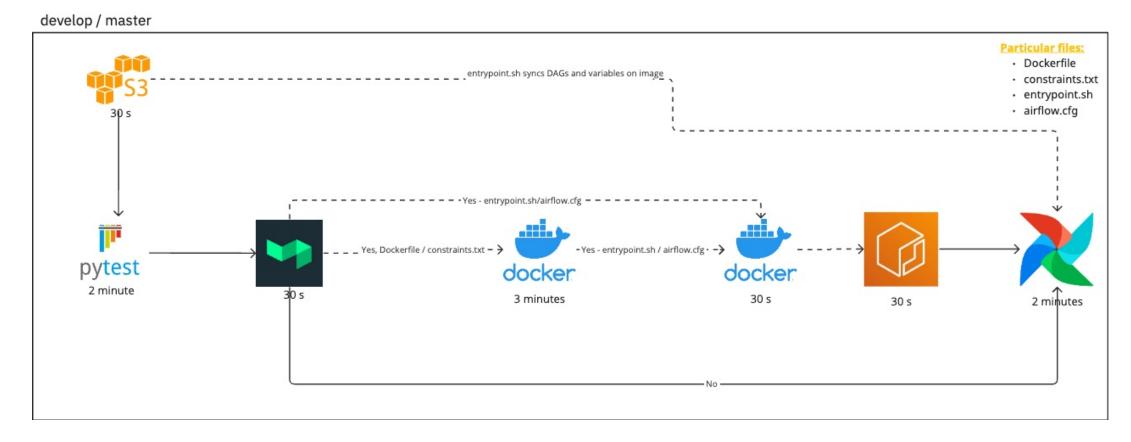
Timing differences - Present CI/CD

• Build time: 10 minutes



Timing differences - Planned CI/CD

- Build times:
 - 5 minutes for DAG / variable changes (most of the time ~95%)
 - 10 minutes for rebuilding image



Future Airflow

- The use of variables
- Dynamic DAGs
- Airflow upgrades

Use of variables

- Migration of all variables from dev.yml and prod.yml
- Use common_airflow.read_airflow_env() function to pull variables.

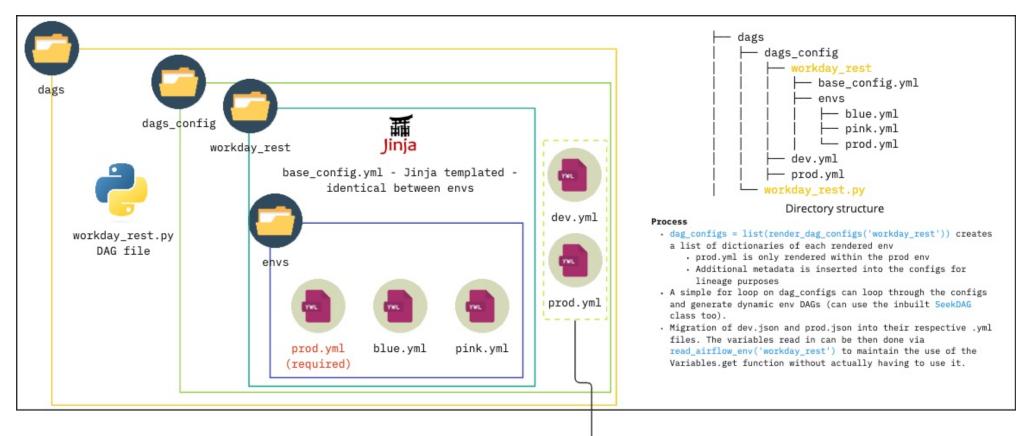
Motivation

Reduce manual effort maintaing multiple DAGs for different envs

Planned Implementation

- A base_config.yml in Jinja2 template to have a consistently render each
- Different dags_config/<dag_name>/envs files
- Use common_airflow.render_dag_configs function; adds additional metadata on rendering.

**Putting it all together



Replicates of how Enterprise DataOps currently use variables

base_config.yml

```
pipeline_suffix: "{{ pipeline_suffix }}"
concurrency: {{ concurrency }}
default_args:
  retries: 1
  retry_delay: {{ default_args.retry_delay }}
list_schedule_interval: {{ list_schedule_interval }}
optimise_schedule_interval: {{ optimise_schedule_interval }}
optimise_glue_databases: {{ optimise_glue_databases }}
optimise_instance_upsize_objects:
 - raw/ecc/gl_items/gl_line_items/
 - raw/ecc/copa_items/copa_items/
  - transformed/semantic.fact_talent_search_activity.delta/
```

prod.yml

```
pipeline_suffix: ""
concurrency: 40
default_args:
    retry_delay: 300
list_schedule_interval: 30 21 * * SAT
optimise_schedule_interval: 0 22 * * SAT
optimise_glue_databases:
    - spectrum_general
    - spectrum_base
    - spectrum_raw
    - spectrum_salesforce
```

Sample Code

```
dags_configs = list(render_dag_configs())
for dag config in dags configs:
    # additional metadata `render dag configs()`
    dag config file name = dag config.get("file name")
    base dag name = dag config.get("base dag name")
    seek_env = dag_config.get("seek_env")
    config_absolute_path = dag_config.get('config_absolute_path')
    dag id = dag config.get('dag id')
    with SeekDAG(
        dag id=dag config.get("dag id"),
        catchup=False,
        default args=default args.
        concurrency=dag_config.get("concurrency"),
        doc md=airflow md env,
        max active runs=1,
        schedule_interval=dag_config.get("optimise_schedule_interval"),
        start date=datetime(2023, 2, 10, tzinfo=local tz),
        seek env config=dag config,
        tags=tags,
    ) as dag:
        logging.info("Rendering DAG %s", dag.dag_id)
        globals()[dag.dag id] = dag
```

Airflow upgrades

- Planned to complete by end of July 2023
- Migration from 1.10.15 to 2.6.2 (with support for Python 3.11)
- Planned changes
 - Backup Postgres database
 - Update all operators
 - Run upgrade_check CLI tool
- Post upgrade change migrate all JSON variables into YML files

? Questions

References

• Migration to Airflow 2