



Bujar Bakiu · Oct 14, 2022 · 5 min read

# Dockerizing dbt Transformations for Managed Airflow: Docker, dbt, and GCP Cloud Composer



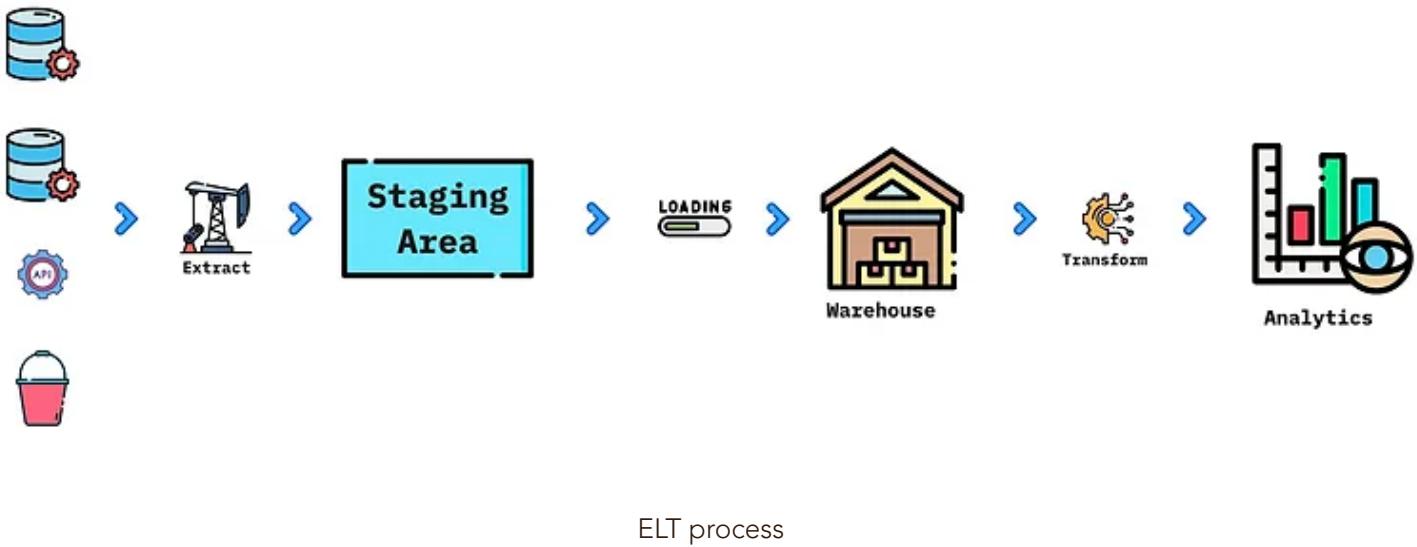
Airflow is one of the most popular pipeline orchestration tools out there. It has been around for more than 8 years, and it is used extensively in the data engineering world.

Popular cloud providers offer Airflow as a managed service e.g: GCP offers Cloud Composer and AWS offers Amazon Managed Workflows for Apache Airflow (MWAA). It is a very fast way to start an ETL (Extract, Transform and Load) pipeline and use it in a production environment.

Meanwhile, dbt (Data Build Tool) is exclusively focused on the Transformation step (the T in ETL). Its popularity has grown significantly and is now defacto the default tool for building powerful and flexible SQL transformations.

Often ETL process is implemented as an ELT (Extract, Load, and Transform) process. First, the data is extracted (or collected) from various sources e.g: databases, APIs, or files of different formats, and stored in a staging area. This

could be a data lake in the form of an S3 bucket. As a second step, the data is loaded in a data warehouse (DWH). Typical data warehouses are Redshift from AWS, BigQuery from GCP, or Snowflake. With the rapid adoption of cloud technologies, storage has become cheaper. Therefore loading the data in the warehouse early is no longer as costly. The last step is Transformations, where the data is aggregated, and adapted to the requirements of the Analytics teams or other stakeholders.



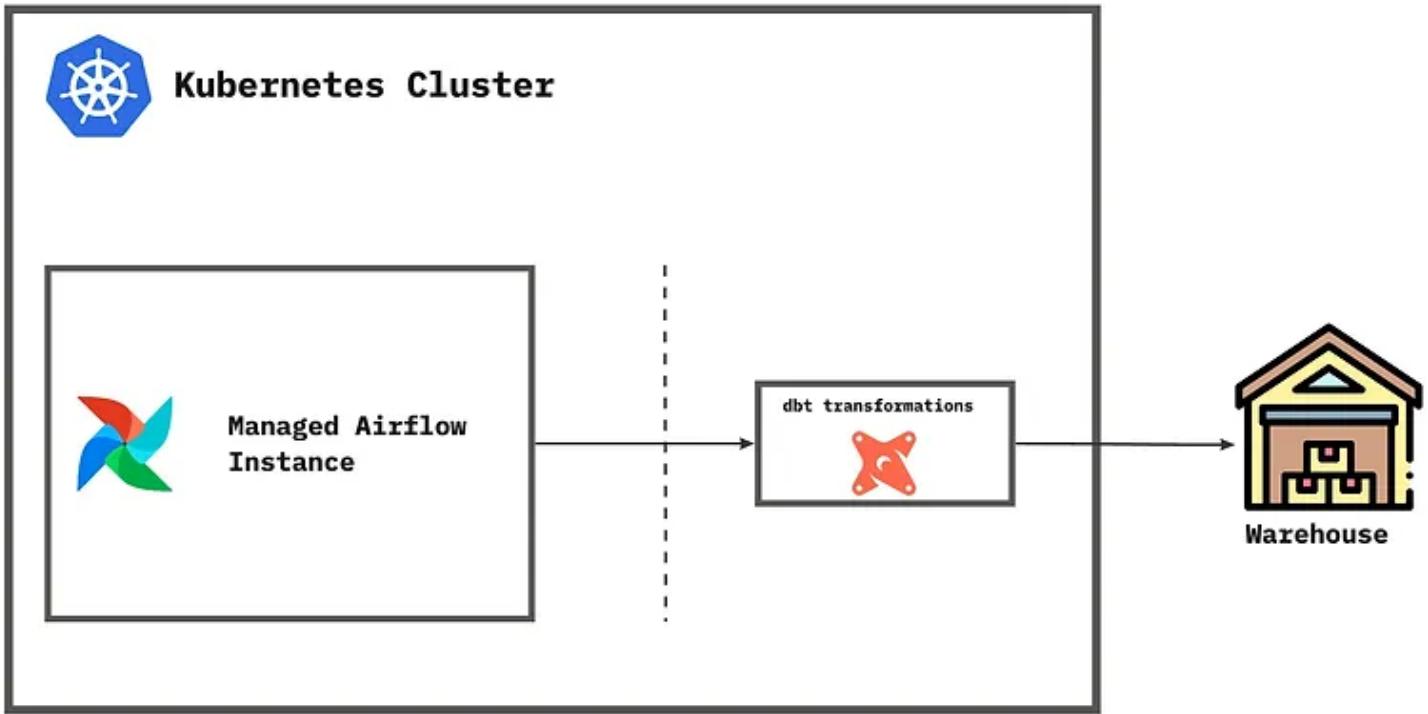
The goal of this post is to show how *dbt* transformations can be integrated into a *managed Airflow instance*. The main challenge we have faced in our projects with managed Airflow instances and dbt is the process of *resolving dependency conflicts*. Airflow dependencies, dbt dependencies, and the default packages pre-installed in the managed Cloud Composer instance do not always work well together. Finding the right setup for them to work is time-consuming and can be obsolete when one dependency changes.

Another limitation of managed Airflow instances is the versions of Python (and Airflow), that are available. They are managed by the cloud provider and are limited in what they support. It takes a while from the moment one version is released until it is available to be used. Similar issues are present in AWS MWAA as well.

**Source Code:** GitHub repo: <https://github.com/data-max-hq/dbt-docker>

## Solution shortly

The idea that we came up with is to run the dbt transformations in an isolated environment e.g: a Kubernetes pod, and run them as independent workloads.



Isolating dbt transformations in a separate pod in Kubernetes

This seems to be a popular idea from other teams that faced similar challenges.

In addition to the obvious environment isolation advantage, this solution provides the added benefit of being agnostic to the underlying infrastructure, be it a managed Kubernetes cluster by a cloud provider or a cluster in on-premise infrastructure.

## But first, a bit of dbt

dbt is a tool focused on the Transformations in ELT. Before dbt, SQL scripts were stored in random files that were not very flexible. dbt brings best practices to SQL transformations. It allows for powerful SQL data modeling, that is flexible and version controlled. In addition, it brings testing and quality check for the data.

Major benefits of integrating dbt in the ELT pipeline include:

- a. Reusability - it is no different than calling a function with different parameters.
- b. History - with source code version control, it allows for keeping track of the changes over time.
- c. Easily document the sources, sinks, and lineage of the data.
- d. Include testing and quality check for the data.
- e. Keeps track of the volume of data inserted, updated, or deleted.

## Why containerize dbt

In order to run the transformations in an isolated environment, e.g: a pod in a Kubernetes cluster, we need to containerize them. This way the transformations are modularised and independent from the other components in the system's architecture.

## How to build and test transformations

Starting a dbt project is pretty easy:

```
$ dbt init <your_project_name>
```

This command will create the directory `<your\_project\_name>` and also create the file `~/.dbt/profiles.yml` somewhere in your local system. To generalize this a bit we came up with the idea of having the `profiles` in the project directory. To achieve this:

```
$ cd <your_project_name>
$ mkdir profiles
$ cd profiles
$ touch profiles.yml
$ cd ..
```

After creating the `profiles.yml` file, fill it with the correct data based on dbt's [documentation](#).

What are dbt profiles?

In short, dbt profiles define how dbt connects to the data warehouse. Using profiles, dbt knows how to connect to a Postgres, BigQuery, Redshift, or Snowflake warehouse.

A `profiles.yml` file looks something like this:

```
1 profile-name:
2   target: dev
3   outputs:
4     dev:
5       type: redshift
6       host: dev-hostname.region.redshift.amazonaws.com
7       user: [dev-username]
8       password: [dev-password]
9       port: 5439
10      dbname: [dev-database name]
11      schema: [dev-database schema]
12      threads: 4
13
14      prod:
15        type: redshift
16        host: prod-hostname.region.redshift.amazonaws.com
17        user: [prod-username]
18        password: [prod-password]
19        port: 5439
20        dbname: [prod-database name]
21        schema: [prod-database schema]
22        threads: 4
```

profiles.yml hosted with ❤ by GitHub

[view raw](#)

In this case, there are two targets, meaning two different warehouses. One is the dev warehouse and the other is the production one. Credentials for each are provided respectively. This is the case for a Redshift warehouse. For BigQuery or Snowflake the profiles file will look similar, but some other attributes might be required.

## How to run a dbt project?

After creating the models, sources, seeds, and other entities testing it is rather easy. One need only run the correct dbt command. Here are some examples:

```
$ dbt compile --profiles-dir ./profiles  
$ dbt seed --profiles-dir ./profiles  
$ dbt run --profiles-dir ./profiles  
$ dbt test --profiles-dir ./profiles
```

## How to wrap it in a Docker container

The docker image will look something similar to the picture below. The first layer will be a version of python, then some system dependencies are added. The third layer contains the dbt python packages required. Finally, layered on top of each other will be the dbt transformations grouped by project.



Make sure to add your project in the docker image by including the directory in this docker file. Lines 24-25 or 29-30 show how to include a dbt project in the docker image.

```

1  FROM python:3.9.13
2
3  # Update and install system packages
4  RUN apt-get update -y && \
5      apt-get install --no-install-recommends -y -q \
6      git libpq-dev && \
7      apt-get clean && \
8      rm -rf /var/lib/apt/lists/* /tmp/* /var/tmp/*
9
10 # Set environment variables
11 ENV DBT_DIR /dbt
12
13 # Set working directory
14 WORKDIR $DBT_DIR
15
16 # Copy requirements
17 COPY requirements.txt .
18
19 # Install DBT
20 RUN pip install -U pip
21 RUN pip install -r requirements.txt
22
23 # Add dbt_project_1 to the docker image
24 COPY dbt_project_1 ./dbt_project_1
25 RUN ["dbt", "deps", "--project-dir", "./dbt_project_1"]
26
27 # Add dbt_project_2 to the docker image
28 # Repeat these line for every project in the repository
29 COPY dbt_project_2 ./dbt_project_2
30 RUN ["dbt", "deps", "--project-dir", "./dbt_project_2"]

```

Dockerfile hosted with ❤ by GitHub

[view raw](#)

## How to run dbt transformations from the docker image

```

$ docker build -t dbt-transformations:latest .
$ docker run dbt-transformations:latest dbt run --project-dir
./<your_project_name> --profiles-dir ./<your_project_name>/profiles

```

First, the docker image is built from the Dockerfile. Afterward, run the image followed by the dbt command. In the example above it will run the dbt transformations for <your\_project\_name>.

## How to push dbt transformations to the cloud

To add the image to the cloud will need to push it to a container registry. It could be Docker Hub, AWS ECR, GCP Artifact Registry, Azure Container Registry, or another self-hosted option.

Using GitHub Actions is rather easy and usually, lots of prebuilt templates exist.

Here are examples of AWS ECR and GCP Artifact Registry.

```

1  name: Build and Push to AWS ECR
2  on:
3    push:
4      branches:
5        - main
6    workflow_dispatch:
7
8  env:
9    REPOSITORY_NAME: transformations-repository
10   IMAGE_NAME: dbt-transformations
11   AWS_REGION: eu-central-1
12
13 jobs:
14   build-push:
15     name: Build and push image to ECR
16     runs-on: ubuntu-latest
17     if: "!contains(github.event.head_commit.message, 'ci skip')"
18     environment: production
19
20   steps:
21     - name: Checkout
22       uses: actions/checkout@v3
23
24     - name: Configure AWS credentials
25       uses: aws-actions/configure-aws-credentials@v1
26       with:
27         aws-access-key-id: ${{ secrets.AWS_ACCESS_KEY_ID }}
28         aws-secret-access-key: ${{ secrets.AWS_SECRET_ACCESS_KEY }}
29         aws-region: ${{ env.AWS_REGION }}
30
31     - name: Login to Amazon ECR
32       id: login-ecr
33       uses: aws-actions/amazon-ecr-login@v1
34
35     - name: Build Docker Image
36       run: docker build -t $IMAGE_NAME:latest .
37
38     - name: Tag, and push image to Amazon ECR
39       env:
40         ECR_REGISTRY: ${{ steps.login-ecr.outputs.registry }}
41       run: |
42         docker tag $IMAGE_NAME:latest $ECR_REGISTRY/$REPOSITORY_NAME
43         docker tag $IMAGE_NAME:latest $ECR_REGISTRY/$REPOSITORY_NAME:$AWS_REGION
44         docker push $ECR_REGISTRY/$REPOSITORY_NAME/$IMAGE_NAME:latest
45         docker push $ECR_REGISTRY/$REPOSITORY_NAME/$IMAGE_NAME:$AWS_REGION

```

build-push-dbt-aws.yaml hosted with ❤ by GitHub [view raw](#)

```

1  name: Build and Push to GCP Artifact registry
2  on:
3    push:
4      branches:
5        - main
6      workflow_dispatch:
7
8  env:
9    PROJECT_ID: ${{ secrets.GCP_PROJECT_ID }}
10   REPOSITORY_NAME: transformations-repository
11   IMAGE_NAME: dbt-transformations
12   GCP_REGION: europe-west1
13
14  jobs:
15    build-push:
16      runs-on: ubuntu-latest
17      steps:
18        - uses: actions/checkout@v2
19
20        - uses: google-github-actions/auth@v0
21        with:
22          credentials_json: ${{ secrets.GCP_CREDENTIALS }}
23
24        - name: Set up Cloud SDK
25          uses: google-github-actions/setup-gcloud@v0
26
27        - name: Configure Docker Client
28          run: gcloud auth configure-docker ${GCP_REGION}-docker.pkg.dev
29
30        - name: Push Docker Image to GCP Artifact Registry
31          run: |
32            gcloud builds submit --tag ${GCP_REGION}-docker.pkg.dev/${REPO_NAME}/main:latest
33            gcloud container images add-tag ${GCP_REGION}-docker.pkg.dev/${REPO_NAME}/main:latest ${IMAGE_NAME}:latest

```

build-push-dbt-gcp.yaml hosted with ❤ by GitHub [view raw](#)

## How to use the transformations in Airflow



KubernetesPodOperator will pull the image of the transformations from the container registry and run it in a Kubernetes Pod (DEH!). Transformations are run in the Pod using the docker image in the desired environment, connected to the DWH. Here is an example of what the code looks like:

```
1  migrate_data = KubernetesPodOperator(  
2      namespace='default',  
3      image='europe-west1-docker.pkg.dev/<project-id>/transformations',  
4      cmds=["dbt", "run"],  
5      arguments=[  
6          '--project-dir", "./<project_dir>", "--profiles-dir", "  
7      ],  
8      name="dbt_transformations",  
9      task_id="dbt_transformations",  
10     get_logs=True  
11 )
```

kubernetes-pod-operator.py hosted with ❤ by GitHub

[view raw](#)

## Summary

This post gave a solution to running dbt transformations in a managed Airflow service from popular cloud providers. However, this solution can be used also in self-managed Airflow deployments.

Using KubernetesPodOperator to run dockerized transformations helps in resolving dependency conflicts, isolating the code, and modularizing the components of the infrastructure.

We would love to hear your feedback. Drop us a line at [hi@data-max.io](mailto:hi@data-max.io).