Google

From S3 to BigQuery - How A First-Time Airflow User Successfully Implemented a Data Pipeline

Leah Cole (with huge thanks to Emily Darrow!)

July 13th, 2020

Intro to Leah





Today's Story

- Prologue
- Chapter 1: BigQuery Public Datasets
- Chapter 2: Growing Pains
- Chapter 3: The Goal
- Chapter 4: The DAG
- Epilogue
- Q&A











Intro to Composer

@leahecole



Google

Google BigQuery

(II)

Google Cloud's **enterprise data warehouse** for analytics

Gigabyte to **petabyte scale** storage and SQL queries

Encrypted, durable, And highly available UNIQUE

Fully managed and **serverless** for maximum agility and scale

UNIQUE

Real-time insights from streaming data

JNIQUE

Built-in **ML** for out-of-the-box predictive insights

UNIQUE

High-speed, in-memory **BI Engine** for faster reporting and analysis



BigQuery: architecture

Serverless. Decoupled storage and compute for maximum flexibility.





Chapter 1: BigQuery Public Datasets



You need to...

Discover the dataset and where to access it.

Negotiate access to the dataset.

Understand the dataset, how it can be joined with your data, and its changes.

The "Data Science" method

Load the data into your systems.

Update, maintain, and secure your data and database.

Manage access and keep the data updated.

Link public data with private data.

Analyze, Visualize and communicate your results.



What if you only did this?

You need to...

Discover the dataset and where to access it.

Negotiate access to the dataset.

Understand the dataset, how it can be joined with your data, and its changes.

Load the data into your systems.

Update, maintain, and secure your data and database.

Manage access and keep the data updated.

Link public data with private data.

Analyze, Visualize and communicate your results.



<u>Current catalog</u> >180 datasets

Onboarded and maintained by Googler(s) with data provider input/guidance

Data providers:















g.co/cloud/marketplace-datasets



Chapter 2: Growing Pains





Growing Pains in the Public Datasets Program

Understand the dataset, how it can be joined with your data, and its changes.

Load the data into your systems.



Late 2019: Onboarding a New Dataset

- New dataset comes in
- Temporarily stored
- Perform transformations
- Ends up in BQ





Late 2019: Problems with Current Process

- Disparate data sources + formats
- Internal/external resource communication
- Access control inconsistent
- Tooling
- Transformations
- Manual















Chapter 3: The Goal





The Goals

- Unified, repeatable process
- Utilize GCP products designed for this
- Hopefully open source process
- See process through eyes of first-time Airflow user (Leah + Emily)





Early 2020-Present: Proposed solution



Chapter 4: The DAG





The DAG Development Process

- Shared repo
- Shared GCP project
 - Leah + Emily both owners
- Shared notes
- Meetings
 - Pairing as needed
 - Regular team meetings





DAG version 0.0

79

42	<pre>move_file_from_s3 = s3_to_gcs_operator.S3ToGoogleCloudStorageOperator(</pre>
43	<pre>task_id='move_file_from_s3',</pre>
44	<pre>bucket=config['source_bucket'],</pre>
45	<pre>prefix='new_dataset/hourly/2019/10/2019-10-01.parquet',</pre>
46	<pre>aws_conn_id="aws_default",</pre>
47	<pre>dest_gcs_conn_id='google_cloud_default',</pre>
48	<pre>dest_gcs='gcs://us-central1-leah-emily-bucket/dags/datasets/',</pre>
49	replace=False,
50	gzip=True
51)
52	make_bq_dataset_for_this = bash_operator.BashOperator(
53	<pre>task_id='make_bq_dataset_for_this', bash_command='bq mk ' + config['target_dataset']</pre>
54)
55	
56	move_parquet_from_gcs_to_bq = gcs_to_bq.GoogleCloudStorageToBigQueryOperator(
57	<pre>task_id='move_parquet_from_gcs_to_bq',</pre>
58	<pre>bucket='us-central1-leah-emily-bucket',</pre>
59	<pre>source_objects=['dags/datasets/new_dataset/hourly/2019/10/2019-10-01.parquet'],</pre>
60	autodetect='true',
61	<pre>source_format=config['source_format'],</pre>
62	destination_project_dataset_table=config['target_dataset_table'], Pro
63	write_disposition='WRITE_TRUNCATE',
64	<pre>trigger_rule='all_done'</pre>
65)
66	
67	
68	delete_bq_dataset_for_this = bash_operator.BashOperator(
69	<pre>task_id='delete_bq_dataset_for_this',</pre>
70	<pre>bash_command='bq rm -rf ' + config['target_dataset'],</pre>
71	<pre>trigger_rule='all_done'</pre>
72)
73	
74	print_dag_finished_message = bash_operator.BashOperator(
75	<pre>task_id='print_dag_finished_message',</pre>
76	<pre>bash_command='echo \"Operation Complete\"',</pre>
77	<pre>trigger_rule='all_done'</pre>
78)

- Get data from S3, store in GCS
- Make target dataset
- Put data into BigQuery

Problem:

Leftover GCS bucket

DAG version 0.1

```
1
2
3
4
5
6
7
8
9
```

```
from gcs_delete_operator import GCSDeleteObjectsOperator
```

```
# Since we have imported the file, we no longer need it. Lets delete it.
delete_parquet_from_gcs = GCSDeleteObjectsOperator(
    task_id="delete_parquet_from_gcs",
    bucket_name=GCS_BUCKET,
    objects=[DEST_FOLDER + SOURCE_PREFIX_DATED]
```

- Get data from S3, store in GCS
- Make target dataset
- Put data into BigQuery
- Delete staging bucket



```
DAG version 1.x - Schema
Definition, Resource
Creation
```



```
{
  "description": "distribution channel this came from",
  "mode": "NULLABLE",
  "name": "data_source",
  "type": "STRING"
},
{
  "description": "download time by hour - UTC",
  "mode": "REQUIRED",
  "name": "time",
  "type": "TIMESTAMP"
},
  "description": "package name (ex: pandas)",
  "mode": "REQUIRED",
  "name": "name",
  "type": "STRING"
},
  "description": "package version (ex: 0.23.0)",
  "mode": "NULLABLE",
  "name": "version",
  "type": "STRING"
},
```

Google

DAG version 1.x - YAML config

16	# s3 bucket template
17	# Uses template syntax to indicate day/month/year.
18	<pre>source_bucket: 'new_public_dataset-package-data'</pre>
19	<pre>source_prefix: 'new_public_dataset/hourly/{year}/{month}/{year}-{month}-{day}.parquet</pre>
20	source_format: 'PARQUET'
21	
22	# BigQuery table
23	<pre>target_dataset: 'new_public_dataset'</pre>
24	<pre>target_dataset_table: 'new_public_dataset.hourly_downloads'</pre>
25	
26	# GCS Config
27	<pre>gcs_bucket: 'us-central1-leah-emily-bucket'</pre>
28	gcs_dest_folder: 'dags/datasets/'
29	
30	# General config
31	lag: 45
32	
33	# Column renames, defined as a key-value mapping - Must be exhaustive for now,
34	# including all columns and column names desired.
35	# This is a mask as well as a dictionary.
36	columns_with_aliases:
37	<pre>data_source: data_source</pre>
38	time: time
39	pkg_name: name
40	pkg_version: version
41	<pre>pkg_platform: platform</pre>
42	<pre>pkg_python: python_version</pre>
43	<pre>counts: total_downloads</pre>
44	<pre>target_time_field_name: 'time'</pre>
45	source_time_field_name: 'time'



DAG version 1.x - Verify

```
146
            # Make sure the configuration is set up correctly
147
            verify configuration = python operator.PythonOperator(
148
                 task id='verify configuration',
                 python callable=check config,
149
150
                 op_kwargs={'config_string': config_string}
151
     # Make sure we've got configuration variables for everything we need, or else don't run.
37
38
     def check_config(config_string):
39
         required variables = (
             'qcs bucket',
40
41
             'gcs_dest_folder',
             'source format',
42
             'source bucket',
43
44
             'target_dataset',
             'target dataset table',
45
             'source prefix',
46
             'columns_with_aliases',
47
             'lag'
48
49
         config_dict = json.loads(config_string)
50
         for config_key in required_variables:
51
52
             assert config_key in config_dict.keys() and config_dict[config_key], "{key} is undefined
```

DAG version 1.x - Verify

154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169

```
# We are assuming that the dataset has already been created. This is an ingestion job, not a setup job.
# This will fail if the dataset doesn't exist, and should echo a message to that effect.
verify_dataset_requirement = bash_operator.BashOperator(
    task_id='verify_dataset_requirement', bash_command='if bq show ' + TARGET_DATASET + '; then ' +
    'echo "Dataset exists"; else echo "Dataset ' + TARGET_DATASET + ' is required - failing now."; exit 1; fi'
)
#
# As above, verify that the dataset contains the master table we're planning to insert into.
# This will fail if the master table doesn't exist.
# We do not verify that the master table's schema matches any expectation at this time.
verify_target_table_requirement = bash_operator.BashOperator(
    task_id='verify_target_table_requirement', bash_command='if bq show --schema ' + TARGET_DATASET_TABLE + '; then ' +
    'echo "Target table exists"; else echo "Target table ' + TARGET_DATASET_TABLE + ' is required - failing now."; exit 1;
)
```



DAG version 1.x - Extract

```
171
          # The S3 to GCS operator copies files from s3 to GCS, but it can copy multiple files.
          # We will be using macros to specify prefixes and elements.
172
173
          # It is worth noting that this operator copies the full folder structure.
174
          #
175
          # Requirements:
176
          # - Boto3 package in python environment
177
          # - S3 Credentials with access to the bucket defined in the airflow configuration.
178
          move_file_from_s3 = s3_to_gcs_operator.S3ToGoogleCloudStorageOperator(
179
              task id='move file from s3',
180
              bucket=SOURCE BUCKET,
              prefix= SOURCE PREFIX DATED.
181
182
              aws_conn_id="aws_default",
              dest_gcs_conn_id='google_cloud_default',
183
              dest_gcs='gcs://' + GCS_BUCKET + '/' + DEST_FOLDER,
184
185
              replace=False,
186
              gzip=True
187
188
189
190
          #
191
          # This is where we actually put things into BigOuery. Since this is a parguet file, we can skip the
          # schema parameter. Parquet files are self-describing. If this were a csv file, we would need to
192
          # describe the expected schema.
193
194
          #
          # It is worth noting that the autodetect parameter is required for parquet files to work.
195
          move parquet from qcs to bq = qcs to bq.GoogleCloudStorageToBigQueryOperator(
196
197
              task_id='move_parquet_from_gcs_to_bq',
198
              bucket=GCS BUCKET.
              source objects=[DEST FOLDER + SOURCE PREFIX DATED].
199
200
              autodetect='true',
              source_format=SOURCE_FORMAT,
201
202
              destination_project_dataset_table=TEMP_TABLE_DATED,
203
              write disposition='WRITE TRUNCATE'
204
```

DAG version 1.x - Transform + Load

```
216
          # Here we transfer the data from temp to Master, with the column remapping. The actual operator.
217
          transfer_data = bigguery_operator.BigQueryOperator(
218
              task id='transfer data',
219
              sql=MERGE_TRANSFORM_STATEMENT,
220
              write_disposition='WRITE_TRUNCATE', # Append to existing tables
221
              create disposition='CREATE NEVER', # Don't create tables.
222
              use_legacy_sgl=False,
              allow large results=True
223
224
206
           #
207
           # Since we have imported the file, we no longer need it. Lets delete it.
208
           delete parquet from gcs = GCSDeleteObjectsOperator(
209
               task id="delete parguet from gcs",
210
               bucket name=GCS BUCKET,
211
               objects=[DEST_FOLDER + SOURCE_PREFIX_DATED]
212
           ۱
226
           #
227
           # We no longer need the temporary table, so we'll wipe it out.
228
           drop temp table = bash operator.BashOperator(
229
               task_id='drop_temp_table',
230
               bash_command='bg rm -f -t ' + TEMP_TABLE_DATED + '; echo "Deleted the temp table"'
231
```

Epilogue





Lessons Learned

- Double check your Composer and Airflow versions
- Documentation is extremely important
- Changelogs and release notes are extremely important
- Transferring data between cloud providers is REALLY easy with Airflow





Call to Action

- Contribute
- Automate
- Collaborate





Thank you to

- Emily Darrow for their technical legwork with this project
- Tim Swast technical advice + vision, moral support
- Shane Glass technical advice, vision, and presentation content
- Rafał Biegacz + the Composer Team presentation content + tireless engineering work
- Seth Hollyman and my Data Analytics DevRel colleagues presentation content, moral support, and constant inspiration
- Moderators, sponsors, and attendees!





Q&A with Leah, Tim, and Shane



