



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

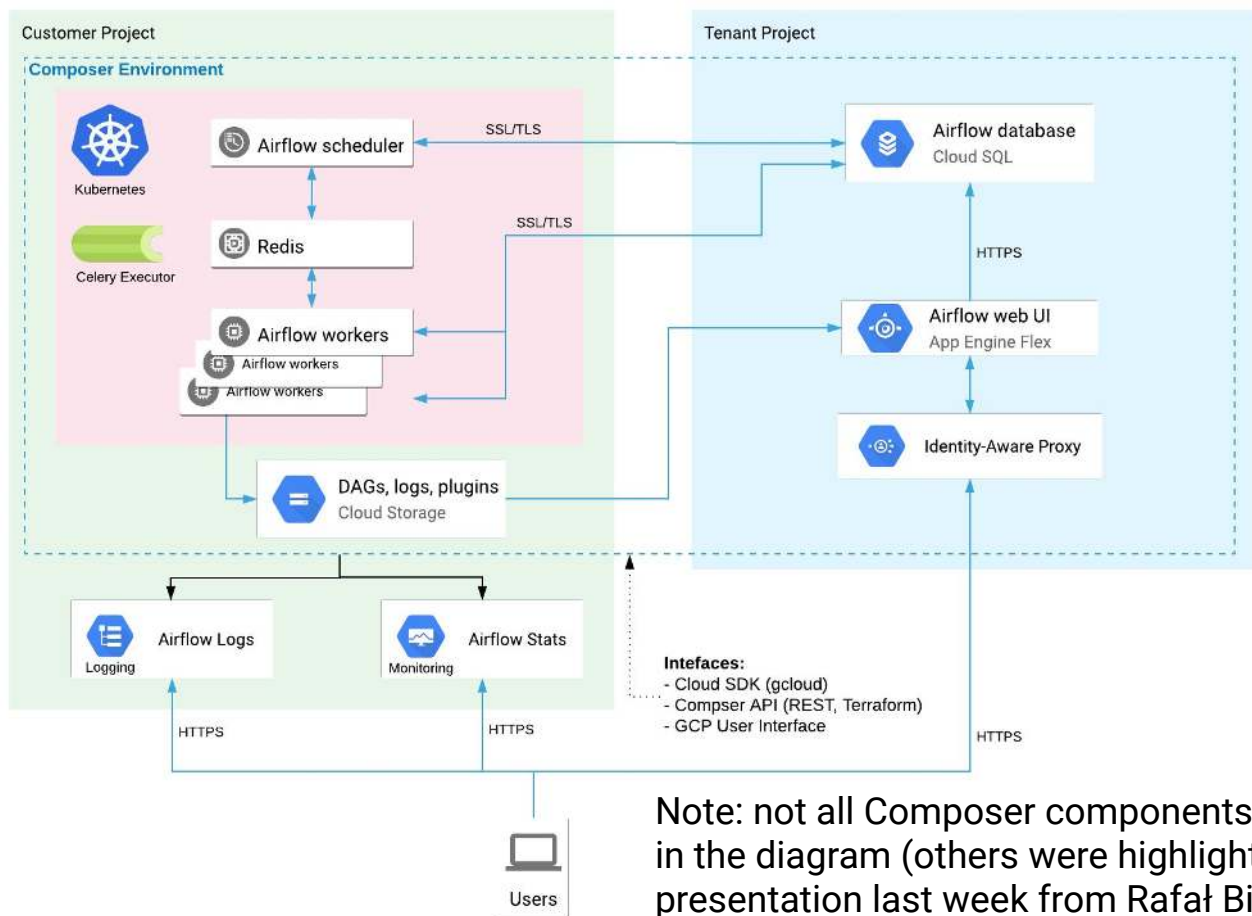
- Prologue
- Chapter 1: BigQuery Public Datasets
- Chapter 2: Growing Pains
- Chapter 3: The Goal
- Chapter 4: The DAG
- Epilogue
- Q&A



@leahecole

Prologue

Intro to Composer



Note: not all Composer components are depicted in the diagram (others were highlighted in a presentation last week from Rafał Biegacz)



@leahecole

Google



Google BigQuery

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

UNIQUE

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

UNIQUE

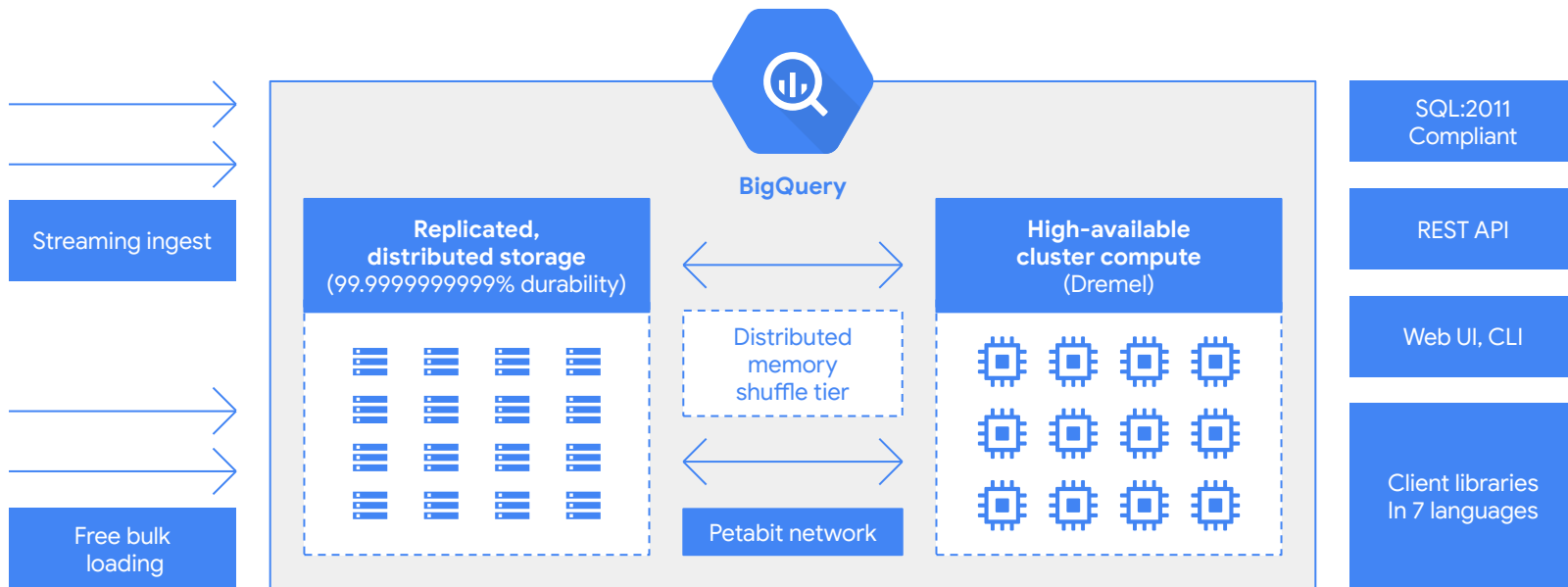
High-speed, in-memory **BI Engine** for faster reporting and analysis



@leahecole

BigQuery: architecture

Serverless. Decoupled storage and compute for maximum flexibility.



@leahecole

Google

Chapter 1: BigQuery Public Datasets



@leahecole



The "Data Science" method

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.



@leahecole

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.



@leahecole

Current catalog

>180 datasets

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

Data providers:



g.co/cloud/marketplace-datasets



@leahecole

Chapter 2: Growing Pains

Growing Pains in the Public Datasets Program

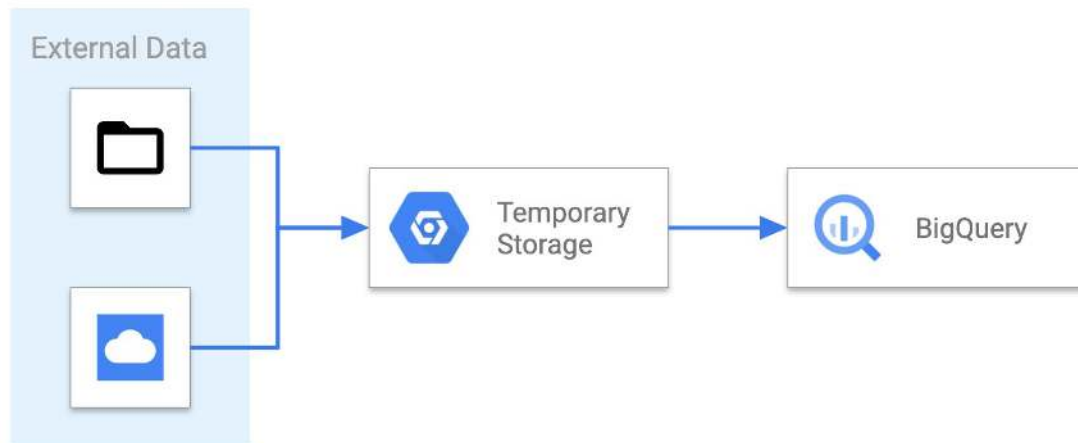
Understand the dataset, how it can be joined with your data, and its changes.

Load the data into your systems.



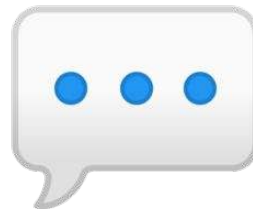
Late 2019: Onboarding a New Dataset

- 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)

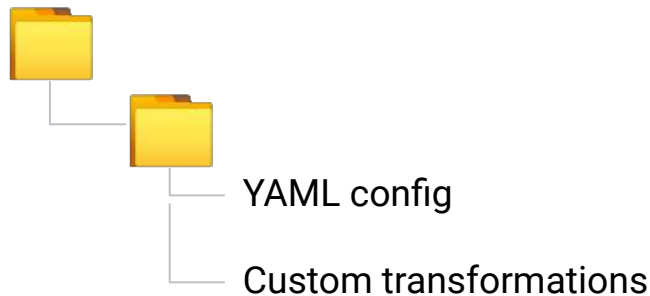


Early 2020-Present: Proposed solution

1. Clone repo,
make branch



2. Add config + transformations



3. Generate DAG + .tf config



4. Create a PR



5. Presubmit checks



6. Human review



7. Deploy



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

```
42 move_file_from_s3 = s3_to_gcs_operator.S3ToGoogleCloudStorageOperator(  
43     task_id='move_file_from_s3',  
44     bucket=config['source_bucket'],  
45     prefix='new_dataset/hourly/2019/10/2019-10-01.parquet',  
46     aws_conn_id="aws_default",  
47     dest_gcs_conn_id='google_cloud_default',  
48     dest_gcs='gcs://us-central1-leah-emily-bucket/dags/datasets/',  
49     replace=False,  
50     gzip=True  
51 )  
52 make_bq_dataset_for_this = bash_operator.BashOperator(  
53     task_id='make_bq_dataset_for_this', bash_command='bq mk ' + config['target_dataset']  
54 )  
55  
56 move_parquet_from_gcs_to_bq = gcs_to_bq.GoogleCloudStorageToBigQueryOperator(  
57     task_id='move_parquet_from_gcs_to_bq',  
58     bucket='us-central1-leah-emily-bucket',  
59     source_objects=['dags/datasets/new_dataset/hourly/2019/10/2019-10-01.parquet'],  
60     autodetect=True,  
61     source_format=config['source_format'],  
62     destination_project_dataset_table=config['target_dataset_table'],  
63     write_disposition='WRITE_TRUNCATE',  
64     trigger_rule='all_done'  
65 )  
66  
67  
68 delete_bq_dataset_for_this = bash_operator.BashOperator(  
69     task_id='delete_bq_dataset_for_this',  
70     bash_command='bq rm -rf ' + config['target_dataset'],  
71     trigger_rule='all_done'  
72 )  
73  
74 print_dag_finished_message = bash_operator.BashOperator(  
75     task_id='print_dag_finished_message',  
76     bash_command='echo \"Operation Complete\"',  
77     trigger_rule='all_done'  
78 )  
79
```

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

Problem:

- Leftover GCS bucket

DAG version 0.1

```
1 from gcs_delete_operator import GCSDeleteObjectsOperator
2
3 # Since we have imported the file, we no longer need it. Lets delete it.
4 delete_parquet_from_gcs = GCSDeleteObjectsOperator(
5     task_id="delete_parquet_from_gcs",
6     bucket_name=GCS_BUCKET,
7     objects=[DEST_FOLDER + SOURCE_PREFIX_DATED]
8 )
9
```

- 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

```
14  ✓ bq mk --table \  
15      --schema hourly_downloads_schema.json \  
16      --time_partitioning_field time \  
17      --clustering_fields name \  
18      --description "New Public Dataset" \  
19      new_public_dataset.hourly_downloads
```

```
{
  "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"
},
{
```

DAG version 1.x – YAML config

```
16 # s3 bucket template
17 # Uses template syntax to indicate day/month/year.
18 source_bucket: 'new_public_dataset-package-data'
19 source_prefix: 'new_public_dataset/hourly/{year}/{month}/{year}--(month)--(day).parquet'
20 source_format: 'PARQUET'
21
22 # BigQuery table
23 target_dataset: 'new_public_dataset'
24 target_dataset_table: 'new_public_dataset.hourly_downloads'
25
26 # GCS Config
27 gcs_bucket: 'us-central1-leah-emily-bucket'
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   data_source: data_source
38   time: time
39   pkg_name: name
40   pkg_version: version
41   pkg_platform: platform
42   pkg_python: python_version
43   counts: total_downloads
44 target_time_field_name: 'time'
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',
149         python_callable=check_config,
150         op_kwargs={'config_string': config_string}
151     )
37 # Make sure we've got configuration variables for everything we need, or else don't run.
38 def check_config(config_string):
39     required_variables = (
40         'gcs_bucket',
41         'gcs_dest_folder',
42         'source_format',
43         'source_bucket',
44         'target_dataset',
45         'target_dataset_table',
46         'source_prefix',
47         'columns_with_aliases',
48         'lag'
49     )
50     config_dict = json.loads(config_string)
51     for config_key in required_variables:
52         assert config_key in config_dict.keys() and config_dict[config_key], "{key} is undefined
```



@leahecole

DAG version 1.x - Verify

```
154 # We are assuming that the dataset has already been created. This is an ingestion job, not a setup job.
155 # This will fail if the dataset doesn't exist, and should echo a message to that effect.
156 verify_dataset_requirement = bash_operator.BashOperator(
157     task_id='verify_dataset_requirement', bash_command='if bq show ' + TARGET_DATASET + '; then ' +
158     'echo "Dataset exists"; else echo "Dataset ' + TARGET_DATASET + ' is required - failing now."; exit 1; fi'
159 )
160
161 #
162 # As above, verify that the dataset contains the master table we're planning to insert into.
163 # This will fail if the master table doesn't exist.
164 # We do not verify that the master table's schema matches any expectation at this time.
165 verify_target_table_requirement = bash_operator.BashOperator(
166     task_id='verify_target_table_requirement', bash_command='if bq show --schema ' + TARGET_DATASET_TABLE + '; then ' +
167     'echo "Target table exists"; else echo "Target table ' + TARGET_DATASET_TABLE + ' is required - failing now."; exit 1;
168 )
169
```

DAG version 1.x – Extract

```
171 # The S3 to GCS operator copies files from s3 to GCS, but it can copy multiple files.
172 # We will be using macros to specify prefixes and elements.
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,
181     prefix= SOURCE_PREFIX_DATED,
182     aws_conn_id="aws_default",
183     dest_gcs_conn_id='google_cloud_default',
184     dest_gcs='gcs://' + GCS_BUCKET + '/' + DEST_FOLDER,
185     replace=False,
186     gzip=True
187 )
188
189
190 #
191 # This is where we actually put things into BigQuery. Since this is a parquet file, we can skip the
192 # schema parameter. Parquet files are self-describing. If this were a csv file, we would need to
193 # describe the expected schema.
194 #
195 # It is worth noting that the autodetect parameter is required for parquet files to work.
196 move_parquet_from_gcs_to_bq = gcs_to_bq.GoogleCloudStorageToBigQueryOperator(
197     task_id='move_parquet_from_gcs_to_bq',
198     bucket=GCS_BUCKET,
199     source_objects=[DEST_FOLDER + SOURCE_PREFIX_DATED],
200     autodetect='true',
201     source_format=SOURCE_FORMAT,
202     destination_project_dataset_table=TEMP_TABLE_DATED,
203     write_disposition='WRITE_TRUNCATE'
204 )
```

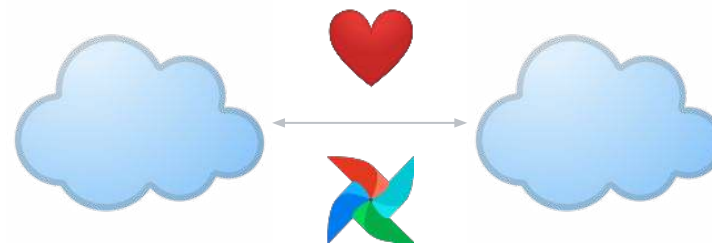
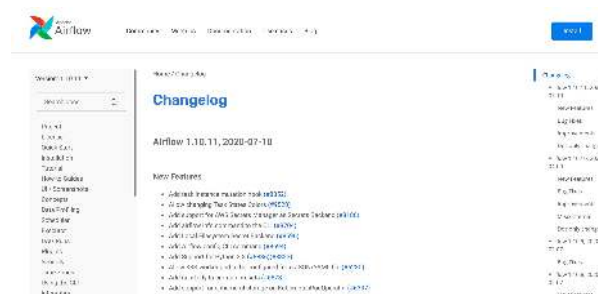
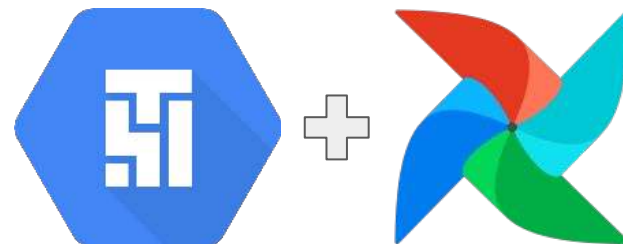
DAG version 1.x - Transform + Load

```
216 # Here we transfer the data from temp to Master, with the column remapping. The actual operator.
217 transfer_data = bigquery_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_sql=False,
223     allow_large_results=True
224 )
225
226 #
227 # Since we have imported the file, we no longer need it. Lets delete it.
228 delete_parquet_from_gcs = GCSDeleteObjectsOperator(
229     task_id="delete_parquet_from_gcs",
230     bucket_name=GCS_BUCKET,
231     objects=[DEST_FOLDER + SOURCE_PREFIX_DATED]
232 )
233 #
234 # We no longer need the temporary table, so we'll wipe it out.
235 drop_temp_table = bash_operator.BashOperator(
236     task_id='drop_temp_table',
237     bash_command='bq rm -f -t ' + TEMP_TABLE_DATED + '; echo "Deleted the temp table"'
238 )
```

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!



@leahecole

Q&A with Leah, Tim, and Shane



@leahecole