



Manage Data Pipelines with Cloud Data Fusion and Cloud Composer

In this module, we will discuss how to manage end-to-end data pipelines using tools in Google Cloud.

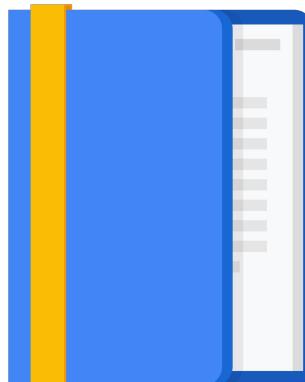
Agenda

Building Batch Data Pipelines visually with Cloud Data Fusion

- Components
- UI Overview
- Building a Pipeline
- Exploring Data using Wrangler

Orchestrating work between GCP services with Cloud Composer

- Apache Airflow Environment
- DAGs and Operators
- Workflow Scheduling
- Monitoring and Logging



The two solutions we will discuss in this module are Cloud Data Fusion and Cloud Composer.

Cloud Data Fusion



Cloud Data Fusion is a **fully-managed, cloud native, enterprise data integration service** for quickly building and managing data pipelines.



Cloud Data Fusion provides a graphical user interface and APIs that increase time efficiency and reduce complexity. It equips business users, developers and data scientists to quickly and easily build, deploy and manage data integration pipelines. Cloud Data Fusion is essentially a graphical no code tool to build data pipelines.



Developer, Data Scientist and Business Analyst

1

Need to cleanse, match, de-dupe, blend, transform, partition, transfer, standardize, automate, and monitor.

2

Use Cloud Data Fusion to visually build integration pipeline, test, debug and deploy.

3

Run it at scale on Google Cloud, operationalize (monitor, report) pipelines, inspect rich integration metadata.



Cloud Data Fusion is used by developers, data scientists, and business analysts alike. For developers, Cloud Data Fusion allows you to cleanse, match, remove duplicates, blend, transform, partition, transfer, standardize, automate, and monitor data. Data scientists can use Cloud Data Fusion to visually build integration pipelines, test, debug, and deploy applications. Business analysts can run Cloud Data Fusion at scale on Google Cloud, operationalized pipelines, and inspect rich integration metadata.



Benefits

Integrate with any data

Increase productivity

Reduce complexity

Increase flexibility

The screenshot shows the Cloud Data Fusion interface with a table titled "sales_clean.txt". The table has 11 rows and 10 columns. The columns are labeled: price, city, zip, type, beds, baths, size, lot_size, and stories. The data represents real estate sales information.

	Double	String	String	String	String	Double	Integer	Long	Integer
1	1000000.0	Santa Clara	95050	Condo	2	2.5	1410	1422	3
2	1000000.0	Santa Clara	95050	Condo	3	2.5	1670	1740	2
3	1000000.0	Santa Clara	95050	Condo	3	2.5	1708	1750	2
4	1000000.0	Santa Clara	95051	Single-Family Home	3	2.0	1068	5600	1
5	1000000.0	Palo Alto	94308	Condo	2	1.5	998	499	2
6	1000000.0	Sunnyvale	94089	Single-Family Home	3	2.0	1108	5824	1
7	1000000.0	Santa Clara	95054	Single-Family Home	3	2.0	1612	6250	1
8	1000000.0	Mountain View	94040	Condo	2	2.0	1206	1880	1
9	1000000.0	Sunnyvale	94085	Condo	2	2.5	1198	1082	3
10	1000000.0	Sunnyvale	94086	Condo	3	3.5	1513	1575	3
11	1000000.0	Santa Clara	95054	Single-Family Home	3	2.0	1097	6200	1



Integrate with any data - through a rich ecosystem of connectors for a variety of legacy and modern systems, relational databases, file systems, cloud services, object stores, NoSQL, EBCDIC, and more.

Increase productivity - If you have to constantly move between numerous systems to gather insight, your productivity is significantly reduced. With Cloud Data Fusion, your data from all the different sources can be pooled into a view like in BigQuery, Spanner or any other Google Cloud technologies, allowing you to be more productive faster.

Reduce Complexity - through a visual interface for building data pipelines, code free transformations, and reusable pipeline templates.

Increase flexibility - through support for on-prem and cloud environments, interoperability with OSS CDAP

Build data pipelines with a friendly UI



Rich graphical interface

100+ plugins - connectors, transforms & actions

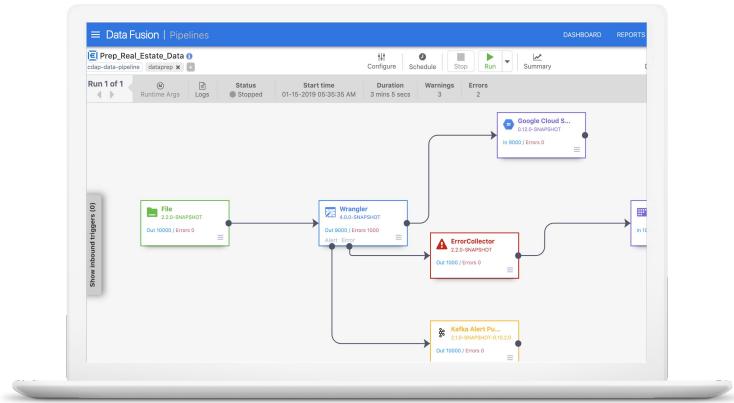
Code free visual transformations

1000+ transforms, data quality

Test and debug pipeline

Pre-built pipelines

Developer SDK



At a high level, Cloud Data Fusion provides you with a graphical user interface to build data pipelines with no code. You can use existing templates, connectors to Google Cloud, and other Cloud services providers and an entire library of transformations to help you get your data in the format and quality you want. Also, you can test and debug the pipeline and follow along with each node as it receives and processes data.



Integration Metadata

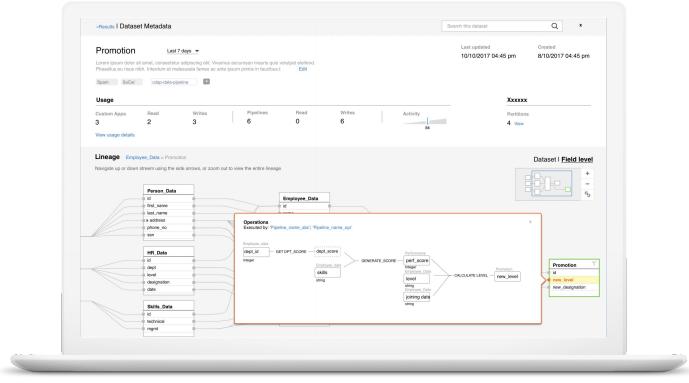
Tags and Properties support

- Pipeline
- Dataset
- Schema

Search integrated entities by

- Keyword
- Schema name and type

Dataset level and Field level Lineage



As you will see in the next section, you can tag pipelines to help organize them more efficiently for your team, and you can use the unified search functionality to quickly find field values or other keywords across your pipelines and schemas. Lastly, we will talk about how Cloud Data Fusion tracks the lineage of transformations that happen before and after any given field on your dataset.

Extensible



Pipeline templatization

Conditional Pipeline Triggers

Plugin Management

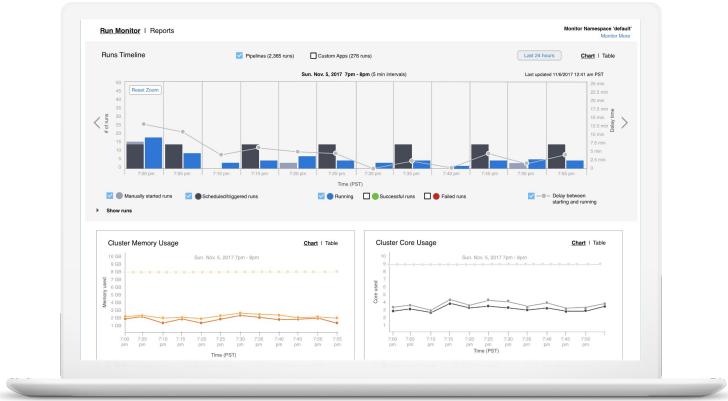
Plugin templatization

Plugin UI Widget

Custom Provisioners

Custom Compute Profiles

Hub Integration



One of the advantages of Cloud Data Fusion is that it's extensible. This includes the ability to template pipelines, create conditional triggers, and manage and template plugins. There is a UI widget plug-in as well as custom provisioners, custom compute profiles, and the ability to integrate to hub.

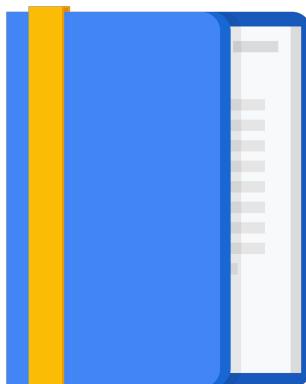
Agenda

Building Batch Data Pipelines visually with Cloud Data Fusion

- [Components](#)
- UI Overview
- Building a Pipeline
- Exploring Data using Wrangler

Orchestrating work between GCP services with Cloud Composer

- Apache Airflow Environment
- DAGs and Operators
- Workflow Scheduling
- Monitoring and Logging





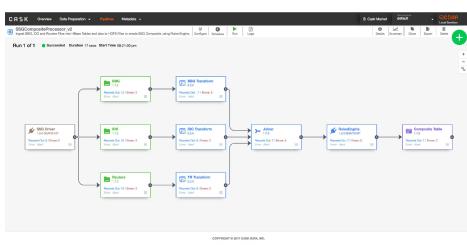
Components of Cloud Data Fusion

The screenshot shows a table with 15 rows of data from a CSV file named 'titanic.csv'. The columns include PassengerId, Survived, Pclass, Name, Sex, Age, SibSp, Parch, Ticket, Cabin, and Embarked. The data includes various passenger details like names, ages, and embarkation points.

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Cabin	Embarked
1	0	3	Miss. Liane Marie Belder	Female	30	1	0	PC1738	B	S
2	1	1	Carly Ms. John Bradley Prentiss Bringe	Male	26	0	0	STON/Oak 321	C	S
3	1	3	Harrison Miss Lena	Female	26	0	0	PC1738	C	S
4	0	1	Kurtis Mr. Jacobus Hart Clio Wigfield	Male	35	1	0	STON/Oak 321	C	S
5	0	3	Leigh Mr. James	Male	35	0	0	PC1738	C	S
6	0	3	Miles Mr. James	Male	35	0	0	AMM7	C	S
7	1	1	McCarthy Mr. Winifred J	Male	34	0	0	PC1738	C	S
8	0	3	Parker Miss. Georgina Woodward	Female	23	1	0	34990	C	S
9	0	3	Perkins Mr. William Gresham	Male	35	0	0	PC1738	C	S
10	0	1	Rosen Mr. Charles H. H. Rosenbaum	Male	34	1	0	34990	C	S
11	1	3	Sawyer Miss. Margaret F	Female	4	1	0	PC1738	C	S
12	0	1	Sennett Miss. Dorothy	Female	50	0	0	113837	C	S
13	0	3	Sparks Mr. George Henry	Male	30	0	0	113837	C	S
14	1	3	Anderson Mr. Andrew John	Male	30	1	0	34990	C	S
15	0	0	Nomura Miss. Kikyo Kameki Melling	Female	14	0	0	34990	C	S

Wrangler — Framework

- Data Preparation for on-boarding new sources and datasets.
- Perform Data Transformations, Data Quality checks with visual feedback.
- Extend the Wrangler by building new user defined directives.
- Integrates with Data Pipeline for operationalizing transformations



Data Pipeline — Framework

- User interface for building complex data workflows
- Join, Lookup, Aggregate, Filtering data in-flight
- Building complex workflows with 100s of connectors
- Extend Data Pipeline using simple APIs
- Integrates with Dataprep, Rule Engine and Metadata Aggregator

The two major user interface components we will focus our attention on in this course, are the Wrangler UI for exploring data sets visually, and building pipelines with no code, and the Data Pipeline UI for drawing pipelines right on to a canvas. You can choose from existing templates for common data processing tasks like Cloud Storage, to BigQuery.



Components of Cloud Data Fusion

The screenshot shows the 'Rules Engine' section of the Cloud Data Fusion interface. It displays a list of rules under a 'Rulebook' named 'Client1CompositeRulebook'. Each rule is defined with a unique ID, name, and a brief description of its purpose. For example, one rule checks if a user's first name is 'John' or 'Mike'. Another rule checks if a user's last name is 'Doe' or 'Smith'. The interface includes a search bar at the top and a detailed view of each rule's configuration below.

Rules Engine — Tool*

- Business Data Transformations and checks codified for business users
- Define Complex rules using intuitive and simple to use user interface
- Logically group Rules in Rulebook and trigger or schedule processing.
- Integrates with Data Pipeline for operationalizing Rules.

The screenshot shows the 'Metadata Aggregator' section of the Cloud Data Fusion interface. It displays a pipeline diagram with nodes labeled 'source', 'transform', and 'sink'. The 'Lineage' tab is selected, showing the flow of data from the source through the transform step to the sink. Below the diagram, there are sections for 'Business', 'Audit Log', 'Provider', 'Usage', and 'Dictionary'. A detailed view of the 'Lineage' section shows the specific fields and their relationships across the pipeline stages.

Metadata Aggregator — Tool

- Aggregate Business, Technical and Operational Metadata
- Track the flow of data (Lineage) for richer data needed for governance
- Create Data Dictionary and Metadata Repository
- Integrate with enterprise MDM solutions.
- Integrates with Data Pipeline, Rules Engine

There are other features of Cloud Data Fusion that you should be aware of too. There's an integrated rules engine where business users can program in their pre-defined checks and transformations, and store them in a single place. Then data engineers can call these rules as part of a rule book or pipeline later. We mentioned data lineage as part of field metadata earlier. You can use the metadata aggregator to access the lineage of each field in a single UI and analyze other rich metadata about your pipelines and schemas as well. For example, you can create and share a data dictionary for your schemas directly within the tool.

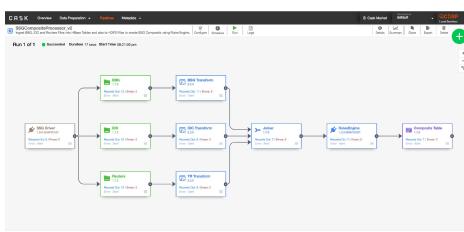


Components of Cloud Data Fusion



Microservice — Framework*

- Build specialized logic for processing data
- Create loosely coupled network for processing events
- Bind processing to varied set of queues



Event Condition Action (ECA) — Application*

- Delivers a specialized solutions for IoT event processing
- Parses any events, triggers conditions and executes Action.
- Real-time notification system, with easy-to-use user interface for configuring event parsing, condition and actions

Other features such as the microservice framework allow you to build specialized logic for processing data. You can also use the Event Condition Action (ECA) Application to parse any event, trigger conditions, and execute an action based on those conditions.

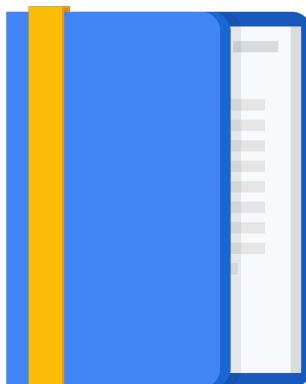
Agenda

Building Batch Data Pipelines visually with Cloud Data Fusion

- Components
- [UI Overview](#)
- Building a Pipeline
- Exploring Data using Wrangler

Orchestrating work between GCP services with Cloud Composer

- Apache Airflow Environment
- DAGs and Operators
- Workflow Scheduling
- Monitoring and Logging





Cloud Data Fusion - UI Overview

- Control Center
- Pipelines
- Wrangler
- Metadata
- Hub
- Entities
- Administration

The screenshot shows the Cloud Data Fusion UI with the title "Cloud Data Fusion | Pipelines". On the left, there's a sidebar with "NAMESPACE" set to "default" and options for "Control Center", "Pipelines" (which is selected), "Studio", "Transform", and "Metadata". The main area is titled "by pipeline name" and contains a table of pipelines. The columns are "Type", "Status", "Last start time", "Next run in", "Total runs", and "Tags". The table lists the following pipelines:

Type	Status	Last start time	Next run in	Total runs	Tags
Batch	Failed	10-18-2018 05:01:35 PM	58 sec	2	Wrangler Real_Estate (cop-data-pipeline)
Batch	Running	10-18-2018 05:01:35 PM	1 hr	10	Wrangler Real_Estate (cop-data-pipeline)
Batch	Succeeded	10-18-2018 05:01:35 PM	1 day	13	Real_Estate SoCal (cop-data-pipeline)
Realtime	Deployed	(cop-data-pipeline)
Realtime	Running	10-18-2018 05:01:35 PM	...	234	Wrangler Real_Estate new (cop-data-pipeline)
Batch	Failed	10-18-2018 05:01:35 PM	...	35	(cop-data-pipeline)
Batch	Succeeded	10-18-2018 05:01:35 PM	1 hr	5,678	(cop-data-pipeline)
Realtime	Succeeded	10-18-2018 05:01:35 PM	1 month	345	(cop-data-pipeline)
Batch	Succeeded	10-18-2018 05:01:35 PM	...	1	(cop-data-pipeline)
Batch	Succeeded	10-18-2018 05:01:35 PM	...	24	(cop-data-pipeline)

At the bottom left of the main area, it says "Namespace Admin".



Here are some of the key user interface elements that you will encounter when using Cloud Data Fusion. Let's look at each of them in turn under Control Center.



Control Center

- Application
- Artifact
- Dataset

Cloud Data Fusion | Control Center

Dashboard Reports Hub System Admin +

Entities in namespace "default"

Displaying Applications, Datasets, sorted by Newest

Type	Name	Programs	Operations	Writes	Running	Failed
Data Pipeline	join-w-customer	1	0	0	0	0
Data Pipeline	sales-with-zip	1	10	3	0	0
Dataset	customers.csv	1	0	0	0	0
Dataset	sales.csv	1	0	0	0	0
Data Pipeline	sales-ingest	2	0	0	0	0
Data Pipeline	customers-ingest	1	12	9	0	0
Dataset	customers	1	0	0	0	0
Dataset	U.S._Chronic_Disease_Indicator...	1	0	0	0	0
Data Pipeline	Mask_data	1	0	0	0	0
Dataset	Mask_data	1	42,119,097	42,119,097	0	0
Data Pipeline	MobileManagementApp	1	0	0	0	0
Dataset	experiment_model_meta	1	0	0	0	0
Dataset	experiment_model_components	1	0	0	0	0
Dataset	experiment_splits	1	0	0	0	0
Data Pipeline	Titanic_test_pipeline_v1_test_v1	2	0	0	0	0
Dataset	titanic.csv	0	0	0	0	0
Data Pipeline	Titanic_test_pipeline_v1_test	2	0	0	0	0
Dataset	Titanic, test	0	0	0	0	0
Dataset	Error_sink_titanic	0	0	0	0	0
Application	dataprep	1	1	0	0	0
Dataset	connections	1	0	0	0	0
Dataset	datapreps	1	0	0	0	0
Dataset	workspace	1	0	239	0	0



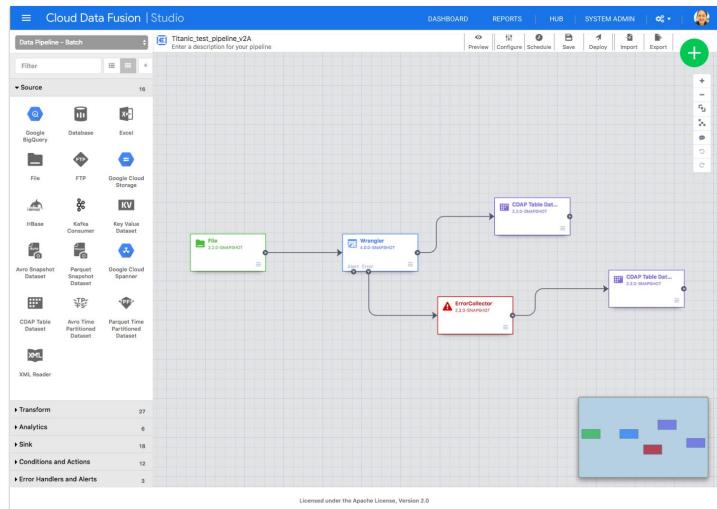
There is a section for applications, artifacts and a dataset. Here you could have multiple pipelines associated with a particular application.

The Control Center gives you the ability to see everything at a glance and search for what you need, whether it's particular dataset, pipeline or other artifact like a data dictionary for example.



Pipelines

- Developer Studio
- Preview
- Export
- Schedule
- Connector and function palette
- Navigation



Under the Pipelines section you have a Developer Studio. You can preview, export, schedule a job or a project. You also have a connector and function palette and a navigation section.



Wrangler

- Connections
 - Transforms
 - Data Quality
 - Insights
 - Functions

Cloud Data Fusion Wrangler													Dashboard	Hub	System Admin	Enterprise Edition	
													Create a Pipeline	More			
Data Insights													Columns (16)				
Integer	String	String	String	String	Integer	String	String	String	String	Double	String	String	#	Name	Completion	Search	Column names ▾
PassengerId	Survived	Pclass	Sex	SibSp	Parch	Ticket	Fare	Cabin	Embarked				1	PassengerId	100%		
1	0	3	male	22	1	A/5 21171	7.25	none	S				2	Survived	100%		
2	1	1	female	38	1	PC 17599	71.3833	CB	C				3	Pclass	100%		
3	1	3	female	26	0	STON/O2. 3101282	7.925	none	S				4	Sex	100%		
4	1	1	female	35	1	313603	53.1	C123	S				5	Age	100%		
5	0	3	male	35	0	373460	8.05	none	S				6	SibSp	100%		
6	0	1	male	54	0	17463	51.9625	E46	S				7	Parch	100%		
7	0	3	male	2	3	349909	21.075	none	S				8	Ticket	100%		
8	1	3	female	27	0	2	347742	11.1333	none	S			9	Embarked	100%		
9	1	2	female	14	1	237736	30.0708	none	C				10	Last_Name	100%		
10	1	3	female	4	1	PP 9549	16.7	G6	S				11	Salutation	100%		
11	1	1	female	58	0	0	113783	26.55	C120	S			12	First_Name	100%		
12	0	3	male	20	0	0	345151	8.05	none	S			13	Title_Freq	100%		
13	0	3	male	39	1	347082	31.275	none	S				14	ID	100%		
14	0	3	female	14	0	0	350468	7.8542	none	S			15				
15	1	2	female	56	0	248706	16.0	none	S				16				
16	0	3	male	2	4	1	382652	29.125	none	Q			17				
17	0												18				
18													19				



Under the Wrangler section you have connections, transforms, data quality, insights, and functions.



Integration Metadata

- Search
- Tags & Properties
- Lineage - Field and Data

The screenshot shows the Cloud Data Fusion Metadata interface. The left sidebar has a 'NAMESPACE' dropdown set to 'default' and a 'Metadata' section selected. The main area displays a search results table with columns for rank, name, type, description, and creation date. The results include:

Rank	Name	Type	Description	Created
18	recipes	Dataset	Recipe store.	Oct 12, 2018
2	workspace	Dataset	Dataspace workspace dataset.	Oct 12, 2018
4		Dataset		Oct 12, 2018
18	dataprefs	Dataset	Store DataPrep index files.	Oct 12, 2018
	connections	Dataset	DataPrep connections store.	Oct 12, 2018
	Error_sink_titanic	Dataset	No description provided for this Dataset.	Oct 17, 2018
	Titanic_test	Dataset	No description provided for this Dataset.	Oct 17, 2018
	Titanic_copy	Dataset	No description provided for this Dataset.	Oct 17, 2018

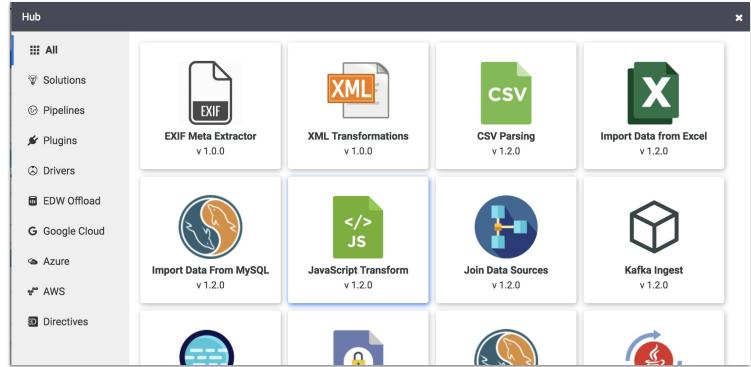


Under the Integration metadata section you can search, add tags and properties, and see the data lineage for field and data.



Hub

- Plugins
- Use cases
- Pre-built pipelines



The Hub allows you to see all the available plugins, sample use cases, and pre-built pipelines.



Entities



- Pipeline
- Application
- Plugin
- Driver
- Library
- Directive

Add entity

 Pipeline A pipeline allows you to create, manage, and operate complex batch and real-time workflows intuitively. Create Import	 Application An application is a collection of datasets and programs that read and write data to datasets. Upload	 Plugin A plugin is an easy way to extend the functionality of an application. Upload
 Driver A driver is a JAR file that contains third-party code to communicate with systems such as MySQL, Oracle, and PostgreSQL using JDBC. Upload	 Library A library is a JAR file that can contain reusable third-party code (e.g. external Spark programs). Upload	 Directive A directive is a data manipulation instruction that can be used to perform data cleansing, transformation and filtering. Upload



Entities include, the ability to create pipelines, upload an application, plug-in, driver, library, and directives.



Administration

- Management
 - Services
 - Metrics
- Configuration
 - Namespace*
 - Compute Profiles
 - Preferences
 - System Artifacts
 - REST Client

Cloud Data Fusion | Administration

DASHBOARD HUB SYSTEM ADMIN

Management | Configuration

Uptime 9 hours 5 mins 32 secs Version: 6.0.0-SNAPSHOT

Services

Status	Name	Provisioned	Requested	Action
Green	App Fabric	1	1	View Logs
Green	Dataset Executor	1	1	View Logs
Green	Log Saver	1	1	View Logs
Green	Messaging Service	1	1	View Logs
Green	Metadata Service	1	1	View Logs
Green	Metrics	1	1	View Logs
Green	Metrics Processor	1	1	View Logs
Green	Transaction	1	1	View Logs

System metrics

Entities	Last hour load
Datasets	0
Programs	0
Namespaces	1
Artifacts	36
Applications	0
Total requests	192,059
Successful	192,057

Transactions

Metric	Value
NumCommittingChangeSets	0
NumInProgressTransactions	0
NumInvalidTransactions	0

Cloud Data Fusion | Administration

DASHBOARD HUB SYSTEM ADMIN

Management | Configuration

Reload system artifacts Make HTTP calls

Namespaces (1) Create, View, and manage namespaces

System Compute Profiles (1) Manage compute profiles available to launch programs in all namespaces

Create New Profile Import

Default	Profile name	Provisioner	Scope	Last 24 hrs runs	Last 24 hrs node hours	Total node hours	Schedules	Triggers	Status
★	Dataproc	Google Cloud DataProc	SYSTEM	—	—	—	0	0	Enabled

System Preferences (1) Manage system preferences available as runtime arguments to programs across all namespaces



There are two components in Administration: Management and Configuration. Under Management you have services and metrics. Under Configuration you have namespace, compute profiles, preferences, system artifacts, and the REST client.

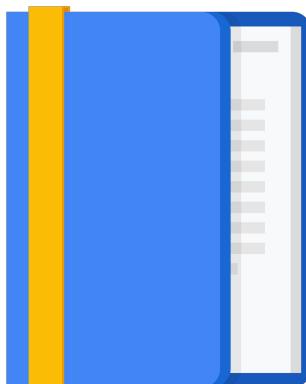
Agenda

Building Batch Data Pipelines visually with Cloud Data Fusion

- Components
- UI Overview
- **Building a Pipeline**
- Exploring Data using Wrangler

Orchestrating work between GCP services with Cloud Composer

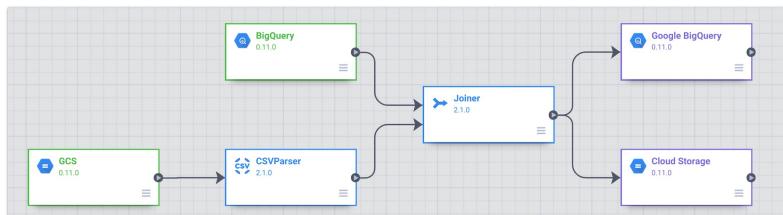
- Apache Airflow Environment
- DAGs and Operators
- Workflow Scheduling
- Monitoring and Logging





Data Pipeline | Directed Acyclic Graph (DAG)

- Represented by a series of stages arranged in a DAG. This forms a **one-way** pipeline.
- Stages, which are the "nodes" in the pipeline graph, can be of different types



A pipeline is represented visually as a series of stages arranged in a graph. These graphs are called DAGs or directed acyclic graphs because they flow from one direction to another and they can not feed into themselves. Acyclic simply means not a circle. Each stage is a node, and as you can see here, it can be of a different type. You may start with a node that pulls data from Google Cloud Storage, then passes it on to a node that parses a CSV. The next node takes multiple nodes, has an input, and joins them together before passing the join data to two separate data sync nodes.

Data Pipeline

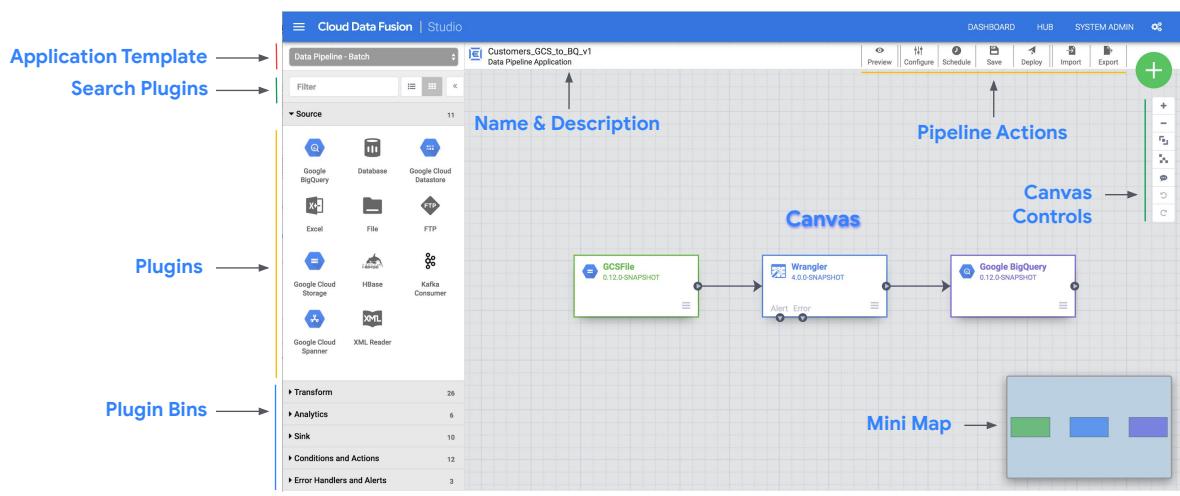


- Allows non-linear pipelines.
- Can fork from a node, where output from a node can be sent to two or more stages.
- Two or more forked nodes can merge at a transform or a sink node.



As you saw in our previous example, you can have multiple nodes fork out from a single parent node. This is useful because you may want to kick off another data processing workstream that should not be blocked by any processing on a separate series of nodes. You can combine data from two or more nodes into a single output in a sync.

Studio is the UI where you create new pipelines

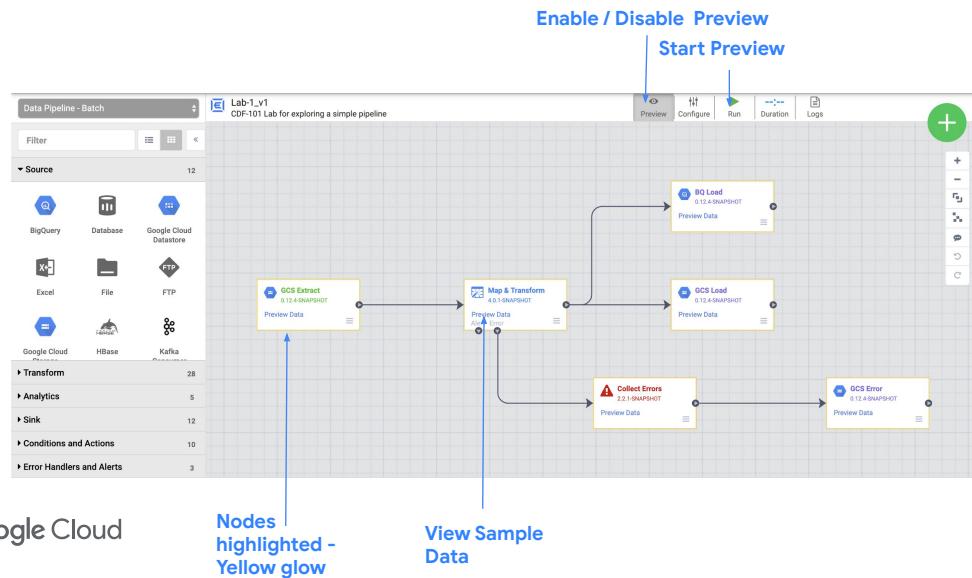


In Cloud Data Fusion, the studio is the user interface where you author and create new pipelines.

The area where you create nodes and chain them together in your pipeline is your canvas. If you have many nodes in a pipeline, the canvas can get visually cluttered, so use the mini map to help navigate around a huge pipeline quickly. You can interact with the canvas and add objects by using the Canvas Control Panel. When you're ready to save and run the entire pipeline, you can do so with the pipeline actions toolbar at the top. Don't forget to give your pipeline a name and description, as well as make use of the many preexisting templates and plugins, so you don't have to write your pipeline from scratch. Here, we've used a template or data pipeline batch which gives us the three nodes you see here to move data from a Cloud Storage file, process it in a wrangler, and output it to BigQuery.



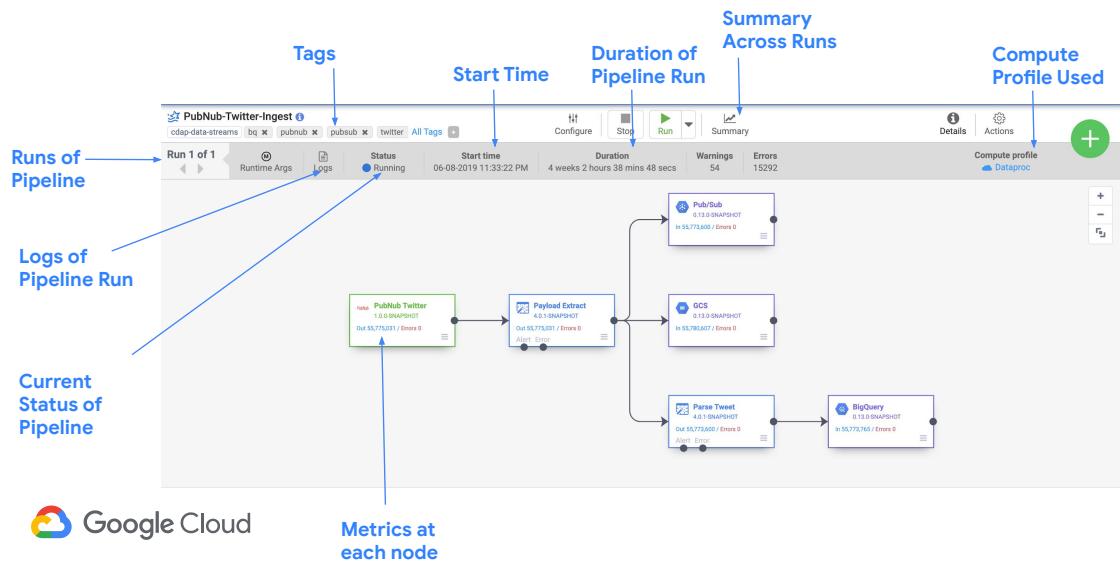
Use Preview Mode to see how the pipeline will run



You should make use of preview mode before you deploy and run your pipeline in production to ensure everything you run will run properly. While a pipeline is in preview, you can click on each node and see any sample data or errors that you will need to correct before deploying.



Monitor the health of the entire pipeline



After deployment, you can monitor the health of your pipeline and collect key summary stats of each execution. Here, we are ingesting data from Twitter and Google Cloud, and parsing each tweet before loading them into a variety of data syncs.

If you have multiple pipelines, it is recommended that you make liberal use of the tags feature to help you quickly find and organize each pipeline for your organization. You can view the start time, the duration of the pipeline run, and the overall summary across runs for each pipeline. You can quickly see the data throughput at each node in the pipeline simply by interacting with the node. Note the Compute profile used in the Cloud. Currently, Cloud Data Fusion supports running on Dataproc, but Cloud Data Source support is on the road map.



Monitor the health of a single node



Clicking on a node gives you detail on the inputs, outputs, and errors for that given node. Here, we are integrating with the Speech-to-Text API to process audio files into searchable text. You can track the individual health of each node and get useful metrics like records out per second, average processing time, and max processing time, which can alert you to any anomalies in your pipeline.

You can schedule batch pipelines



The screenshot shows the Cloud Data Fusion Pipeline interface. On the left, there's a sidebar with tabs for 'No Runs', 'Runtime Args', 'Logs', and 'Status'. The main area is titled 'Configure schedule for pipeline batch_pipeline'. It has a 'Basic' tab selected, showing a dropdown for 'Pipeline run repeats' set to 'Daily', 'Repeats every 1 day(s)', and 'Starting at 1:00 AM'. A summary note says: 'This pipeline is scheduled to run everyday, at 1:00AM. The pipeline cannot have concurrent runs.' Below this are fields for 'Max concurrent runs' (set to 1) and 'Compute profiles' (set to 'Dataproc (Google Cloud Dataproc)'). At the bottom are two buttons: 'Save and Start Schedule' (highlighted in blue) and 'Save Schedule'.



You can set your pipelines to run automatically at certain intervals. If your pipeline normally takes a long time to process the entire dataset, you can also specify a maximum number of concurrent runs to help avoid processing data unnecessarily. Keep in mind that Cloud Data Fusion is designed for batch data pipelines. We will dive into streaming data pipelines in future modules.

After data is transformed, you can track field-level lineage



Relationship between fields of datasets: Provenance, Impact

The screenshot shows the Cloud Data Fusion Field Level Lineage interface. At the top, there's a navigation bar with 'DASHBOARD', 'HUB', 'SYSTEM ADMIN', and 'Enterprise Edition'. Below the navigation is a breadcrumb trail: '« Back | doubleclick | Dataset'. The main area is titled 'Field level lineage' with the subtitle 'Explore root cause and impact for each of the fields of the dataset'. A dropdown menu 'Last 7 days ▾' is open. On the left, under 'Cause for: doubleclick: campaign', there's a table for '1 Dataset' with columns 'Dataset name' and 'Field name'. It lists 'campaign' with 'offset' and 'body' as its fields. In the center, under 'Dataset: doubleclick', there's a table for '8 Fields' with columns 'Field name' and a search bar. It lists 'advertiser_id', 'timestamp', 'referrer_url', 'campaign_id', 'landing_page_url', 'advertiser', and 'landing_page_id'. On the right, under 'Impact for: doubleclick: campaign', there's a table for '1 Dataset' with columns 'Dataset name' and 'Field name'. It lists 'hits' with 'campaign' as its field. A central diagram shows the flow of data from 'campaign' through various fields to 'hits', with arrows indicating 'Incoming operations' pointing to 'campaign' and 'Outgoing operations' pointing away from it.

One of the big features of Cloud Data Fusion is the ability to attract the lineage of a given field value. Let's take this example of a campaign field for double-click dataset and track every transform operation that happened before and after this field.

See every operation that is made on a field



Operations applied to a field

A screenshot of the Cloud Data Fusion interface. The top navigation bar shows 'Cloud Data Fusion | Field Level Lineage'. On the right, there are links for 'DASHBOARD', 'HUB', 'SYSTEM ADMIN', and 'Enterprise Edition'. Below the navigation, a breadcrumb trail shows 'doubleclick > Dataset'. A modal window titled 'Cause operations for field 'campaign'' is open, displaying the lineage information. The modal has a dark header and a light body. It shows 'Operations between 'campaign' and 'doubleclick''. There is a note 'Last executed by '00-BigQuery_SQL_DoubleClick_v1' on 03-27-2019 09:48:28 AM'. A table below lists three operations:

Input	Input fields	Operation	Description	Output fields	Output
1	campaign	campaign.Re...	Read from Google Cloud Storage.	offset, body	--
2	--	[offset], [body]	parse campaigns.P... Data	advertiser_id, campaign_id, campaign	--
3	--	[campaign]	Joiner3.Iden... parse campaigns.c...	campaign	--

Here, you can see the lineage of operations that are applied to the campaign field between the campaign dataset and the double-click dataset. Note the time this field was last changed by a pipeline run and each of the input fields and descriptions that interacted with the field as part of processing it between datasets. Imagine the use cases if you have inherited a set of analytical reports and you want to walk back upstream all of the logic that went into a certain field. Well now, you can.

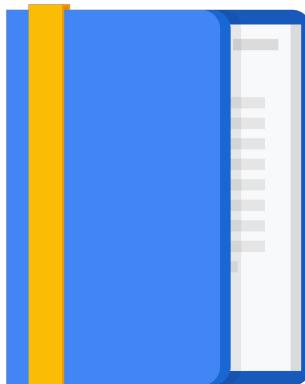
Agenda

Building Batch Data Pipelines visually with Cloud Data Fusion

- Components
- UI Overview
- Building a Pipeline
- Exploring Data using Wrangler

Orchestrating work between GCP services with Cloud Composer

- Apache Airflow Environment
- DAGs and Operators
- Workflow Scheduling
- Monitoring and Logging



We've discussed the core components, tools, and processes of building data pipelines. Now, we'll look at using Wrangler to explore the dataset.



Data analysts can explore datasets in the Wrangler

Wrangler is a code-free, visual environment for transforming data in data pipelines

The screenshot shows the Google Cloud Data Fusion Wrangler interface. On the left, a sidebar lists 'Namespace' (default), 'Control Center', 'Pipeline', 'List', 'Studio', 'Wrangler' (selected), and 'Metadata'. Below this is a 'Google' logo. The main area has two tabs: 'Select data' and 'Pipeline'. Under 'Select data', there's a 'Connections in "default"' section with 'Upload', 'Database (2)' (Cloud SQL MySQL TheSQL, Cloud SQL Postgresql), 'Kafka (0)', 'S3 (1)', and 'Google Cloud Storage (2)' (Cloud Storage Default, Sample Buckets). A '+' button is at the top right. The 'Pipeline' tab shows a pipeline named 'credit_card_data...'. It contains three stages: 'credit_card_data.csv' (Google Cloud Storage), 'customers.csv' (Google Cloud Storage), and 'customers.csv' (Google Cloud Storage). The middle stage is highlighted. Stage details include 'Data' and 'Insights' tabs, and a 'Create a Pipeline' button. Below the stages is a table with 8 rows of data. To the right of the table is a 'Columns (6)' section and a 'Directives (4)' section with the following content:

#	Directives
1	parse-as-csv body;"true'
2	drop body
3	send-to-error !dq!isCreditCard(credit_card)
4	filter-rows-on regex-match credit_card_type

At the bottom right, it says 'Instance ID: cloud-data-fusion-demos/cdf'.

So far in the course, we have focused on building new pipelines for our datasets. That presumes we know what the data is and what transformations need to be made already. Oftentimes, a new dataset still needs to be explored and analyzed for insights. The Wrangler UI is the Cloud Data Fusion environment for exploring new datasets visually for insights. Here, you can inspect the dataset and build a series of transformation steps called directives to stitch together a pipeline.



Wrangler UI Overview for exploring datasets

The screenshot illustrates the Wrangler UI interface. On the left, the 'Connection' sidebar shows 'Connections in "default"' including 'Upload', 'Database (2)' (Cloud SQL MySQL, PostgreSQL), and 'Kafka (0)'. A blue arrow points from the 'Add New Connection Type' button to the 'Add Connection' button in the sidebar. Below the sidebar is the 'Google Cloud' logo. In the center, the 'Workspace' area displays a 'Select data' search bar and a list of files under 'Root / demo-datasets / csv'. A blue arrow labeled 'Live Browsing of Connection' points to the file list. To the right, the 'Sample Insights' section shows a preview of 'customers.csv' with columns 'id', 'first_name', 'last_name', and 'email'. A blue arrow labeled 'Power Mode' points to the bottom of this section. Further right, the 'Turn all transformation into Pipeline' section includes a green '+' button, 'Add Directive', 'Recipe', 'Data Quality', and 'Schema' tabs. A blue arrow labeled 'Schema' points to the schema definition area. At the bottom right is a table preview of the 'customers.csv' data.

	String	String	String	String
1	id	first_name	last_name	email
2	1	Joyce	Taylor	jaylor0@bbc.co.uk
3	2	Katherine	Wallace	kwallace1@aristoteles.com
4	3	Sandra	Mcdonald	smcdonald2@outlook.com
5	4	Henry	Crawford	hcrawford3@ukfast.co.uk
6	5	Lawrence	Lane	llane4@mapquest.com
7	6	Jean	Garza	jgarza5@usatoday.com
8	7	Gerald	Austin	gaustin6@tele2.se
9	8	Lawrence	Adams	ladams7@scribd.com

Here's what the Wrangler UI looks like. Starting from the left, you have your connections to existing datasets. Here, you can add new connections to a variety of data sources like Google Cloud Storage, BigQuery, or even other cloud providers. Once you specify your connection, you can browse all of the files or tables in that source. Here, you see a Cloud Storage bucket of demo data sets and all the CSV files of customer complaints. Once you've found an example dataset like `customers.csv` here, you can explore the rows and columns visually and view sample insights. As you explore the data, you might want to create new calculated fields, drop columns, filter rows, or otherwise wrangle the data. You can do so using the Wrangler UI by adding new directives to form a data transformation recipe. When you're happy with your transformations, you can create a pipeline that you can then run at regular intervals.



Building and executing a pipeline graph in Cloud Data Fusion

Objectives

- Connect Cloud Data Fusion to a couple of data sources
- Apply basic transformations
- Join two data sources
- Write data to a sink

Building and Executing a Pipeline Graph with Data Fusion

<https://gcpstaging.qwiklabs.com/labs/24798/edit>

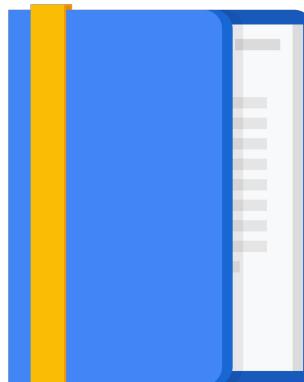
Agenda

Building Batch Data Pipelines visually with Cloud Data Fusion

- Components
- UI Overview
- Building a Pipeline
- Exploring Data using Wrangler

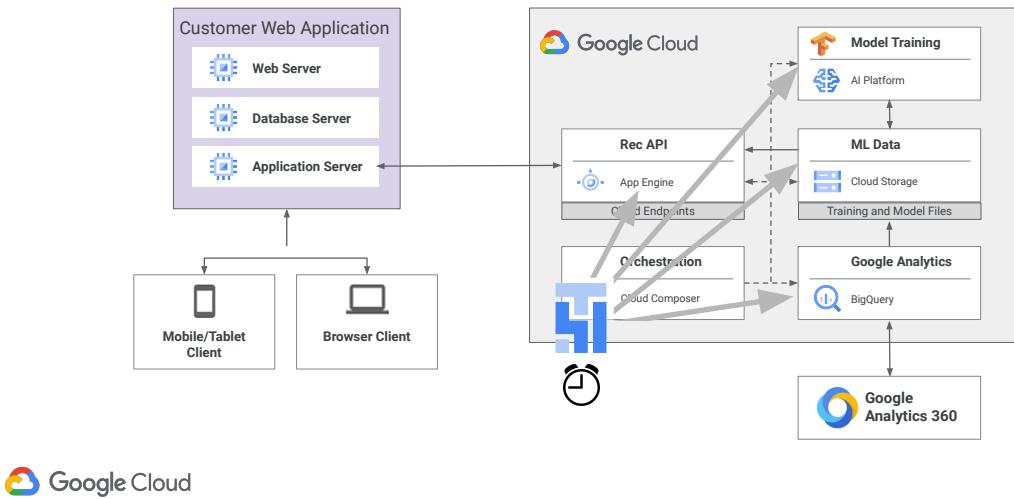
Orchestrating work between GCP services with Cloud Composer

- Apache Airflow Environment
- DAGs and Operators
- Workflow Scheduling
- Monitoring and Logging



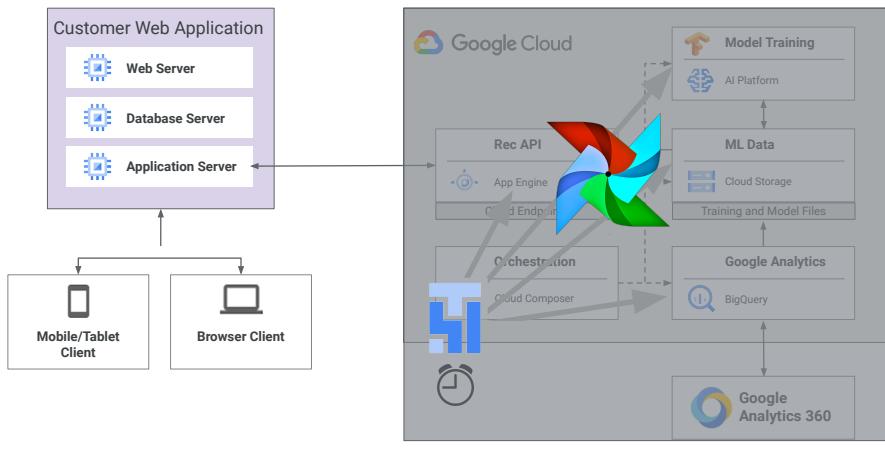
The next big task for managing data pipelines is to orchestrate the work across multiple Google Cloud services. For example, if you had three Cloud Data Fusion Pipelines and two ML models that you want it to run in a certain order, you need an orchestration engine. In this module, we'll look at using Cloud Composer to help out with tasks like that.

Cloud Composer orchestrates automatic workflows



Cloud Composer will command the Google Cloud services that we need to run. But Cloud Composer is simply a serverless environment on which an open source workflow tool runs.

Cloud Composer is managed Apache Airflow



Google Cloud

That workflow tool is called Apache Airflow, which is an open-source orchestration engine.

Use Apache Airflow DAGs to orchestrate GCP services



```
bq_rec_training_data → bq_export_op → ml_engine_training_op → app_engine_deploy_version
```

DAG = Directed Acyclic Graph



The heart of any workflow (workflow pipeline and DAG are interchangeable), is to DAG. As you saw with Cloud Data Fusion, you're also building DAGs with Apache Airflow as you see here. What's happening in this particular DAG are four tasks that update our training data, export it, we train our model, and we deploy it. You can tell your DAG to pretty much do anything you need it to do. Here it's sending tasks to BigQuery, Cloud Storage, and AI Platform, but yours could orchestrate among four completely different services.

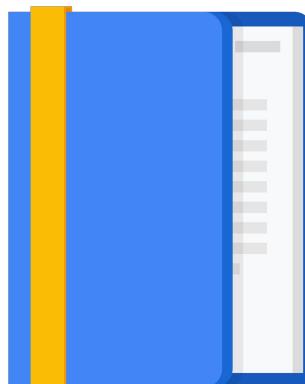
Agenda

Building Batch Data Pipelines visually with Cloud Data Fusion

- Components
- UI Overview
- Building a Pipeline
- Exploring Data using Wrangler

Orchestrating work between GCP services with Cloud Composer

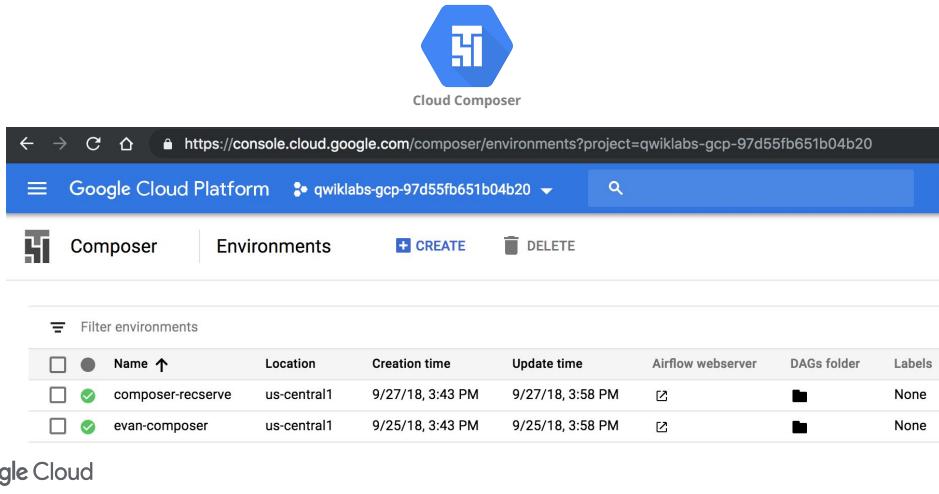
- Apache Airflow Environment
- DAGs and Operators
- Workflow Scheduling
- Monitoring and Logging



Building any workflow in Cloud Composer consists of these four steps and building a ML workflow for our recommendation engine is no different.

We'll start by previewing the actual Cloud Composer environment.

Cloud Composer creates managed Apache Airflow environments



The screenshot shows the Google Cloud Platform interface for Cloud Composer environments. At the top, there's a navigation bar with the URL <https://console.cloud.google.com/composer/environments?project=qwiklabs-gcp-97d55fb651b04b20>. Below it, the main header reads "Google Cloud Platform" and "qwiklabs-gcp-97d55fb651b04b20". The main content area is titled "Composer" and "Environments". It features a "CREATE" button and a "DELETE" button. A table lists two environments:

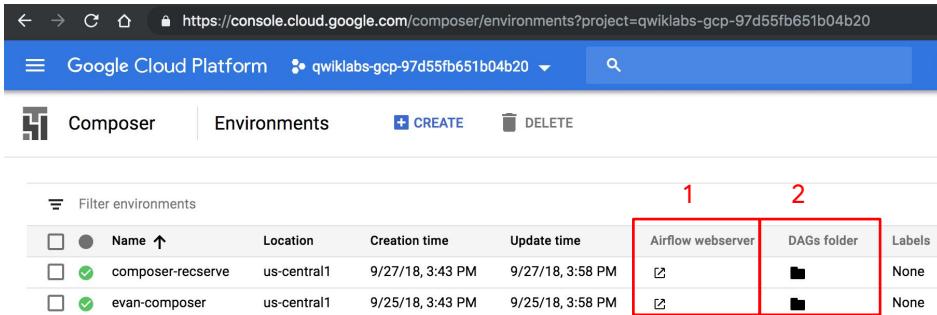
Name	Location	Creation time	Update time	Airflow webserver	DAGs folder	Labels
composer-recserve	us-central1	9/27/18, 3:43 PM	9/27/18, 3:58 PM	[Edit]	[Edit]	None
evan-composer	us-central1	9/25/18, 3:43 PM	9/25/18, 3:58 PM	[Edit]	[Edit]	None

At the bottom left, there's a "Google Cloud" logo.

Once you use the command line or GCP web UI to launch a Cloud Composer instance, you'll be met with a screen like this. Keep in mind that you can have multiple Cloud Composer environments and which each environment you can have a separate Apache Airflow instance which could have zero to many DAGs.

An important note here is that sometimes you'll be required to edit environment variables for your workflows (like specifying your specific GCP project account). Normally you will not do that at the Cloud Composer level but on the actual Apache Airflow instance level which we will show you in the next demo. Again, generally I'm only on the Cloud Composer page here to create new environments before I launch directly into the Airflow webserver.

Each Airflow environment has a separate webserver and folder in GCS for pipeline DAGs



The screenshot shows the Google Cloud Platform Composer environments page. At the top, there's a navigation bar with the URL <https://console.cloud.google.com/composer/environments?project=qwiklabs-gcp-97d55fb651b04b20>. Below it, the Google Cloud Platform header includes 'Google Cloud Platform' and 'qwiklabs-gcp-97d55fb651b04b20'. The main content area has tabs for 'Composer' and 'Environments', with 'Environments' selected. It features a 'CREATE' button and a 'DELETE' button. A 'Filter environments' dropdown is open, showing 'Name ↑' as the current sort order. The table lists two environments:

Name	Location	Creation time	Update time	Airflow webserver	DAGs folder	Labels
composer-recserve	us-central1	9/27/18, 3:43 PM	9/27/18, 3:58 PM	View	View	None
evan-composer	us-central1	9/25/18, 3:43 PM	9/25/18, 3:58 PM	View	View	None

Red numbers 1 and 2 are overlaid on the 'Airflow webserver' and 'DAGs folder' columns respectively, highlighting the specific components being discussed.

To access the Airflow admin UI where you can monitor and interact with your workflows you'll click on the link underneath Airflow webserver. The second box you see is the DAGs folder which is where the code of your actual workflows will be stored.

The DAGs folder is simply a GCS bucket where you will load your pipeline code



Cloud Composer

Buckets / us-central1-evan-composer-0e85530c-bucket / dags						
Name	Size	Type	Storage class	Last modified	Public access	Encryption
dataflow/	—	Folder	—	—	Per object	—
simple_load_dag.py	6.79 KB	text/x-python-script	Multi-Regional	10/1/18, 1:11 PM	Not public	Google-managed key
simple.py	2.51 KB	text/x-python-script	Multi-Regional	10/1/18, 1:10 PM	Not public	Google-managed key



The DAGs folder for each Airflow instance is simply a GCS bucket that is automatically created for you when you launch your Cloud Composer instance. Here is where you upload your DAG files, written in python, and bring your first workflow to life in Airflow. Let's take a look at a quick demo of what the simplest DAG looks like in the UI.

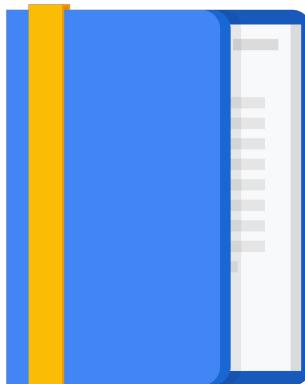
Agenda

Building Batch Data Pipelines visually with Cloud Data Fusion

- Components
- UI Overview
- Building a Pipeline
- Exploring Data using Wrangler

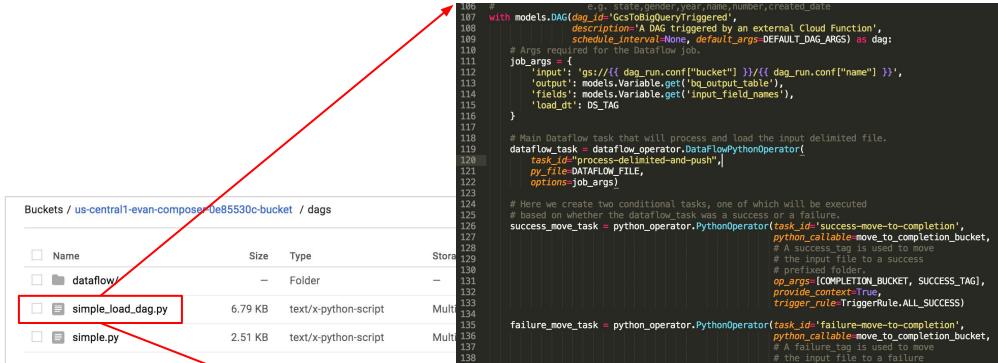
Orchestrating work between GCP services with Cloud Composer

- Apache Airflow Environment
- **DAGs and Operators**
- Workflow Scheduling
- Monitoring and Logging



Now that you're familiar with the basic environment setup, it's time to discuss your primary artifact which is your DAG and the operators you are using to call whichever services you want to send tasks to.

Airflow workflows are written in Python

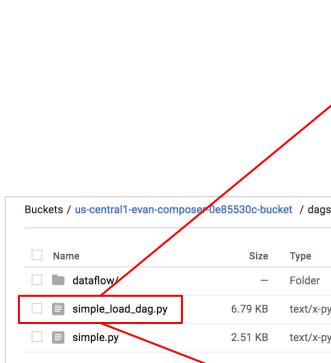


```
105
106
107     with models.DAG(dag_id='gcsToBigQueryTriggered',
108                      description='A DAG triggered by an external Cloud Function',
109                      schedule_interval=None, default_args=DEFAULT_DAG_ARGS) as dag:
110
111     # A new record for the Dataflow job.
112     job_args = [
113         'input': 'gs://{{ dag_run.conf["bucket"] }}/{{ dag_run.conf["name"] }}',
114         'output': models.Variable.get('bg_output_table'),
115         'fields': models.Variable.get('input_field_names'),
116         'load_dt': DS_TAG
117     ]
118
119     # Main Dataflow task that will process and load the input delimited file.
120     dataflow_task = dataflow_operator.DataFlowPythonOperator(
121         task_id='dataflow-delimited-and-push',
122         py_file=DATAFLOW_FILE,
123         options=job_args
124     )
125
126     # Here we create two conditional tasks, one of which will be executed
127     # based on whether the Dataflow task was a success or a failure.
128     success_move_task = python_operator.PythonOperator(task_id='success-move-to-completion',
129                                                       python_callable=move_to_completion_bucket,
130                                                       op_args=[COMPLETION_BUCKET, SUCCESS_TAG],
131                                                       provide_context=True,
132                                                       trigger_rule=TriggerRule.ALL_SUCCESS)
133
134     failure_move_task = python_operator.PythonOperator(task_id='failure-move-to-completion',
135                                                       python_callable=move_to_completion_bucket,
136                                                       op_args=[COMPLETION_BUCKET, FAILURE_TAG],
137                                                       provide_context=True,
138                                                       trigger_rule=TriggerRule.ALL_FAILED)
139
140
141     # The success_move_task and failure_move_task are both downstream from the
142     # dataflow_task.
143     dataflow_task >> success_move_task
144     dataflow_task >> failure_move_task
145
146
147
```



First, Airflow workflows are written in python. You'll have one python file for each DAG. For example, here we have `simple_load_dag.py` in our DAG Folder GCS bucket and you can see a preview of what the DAG file looks like. Don't worry about reading the code, we'll go into that later. It's sufficient enough for now to just know that there are a series of user-created tasks in each DAG file that invoke predefined operators.

Airflow workflows are written in Python



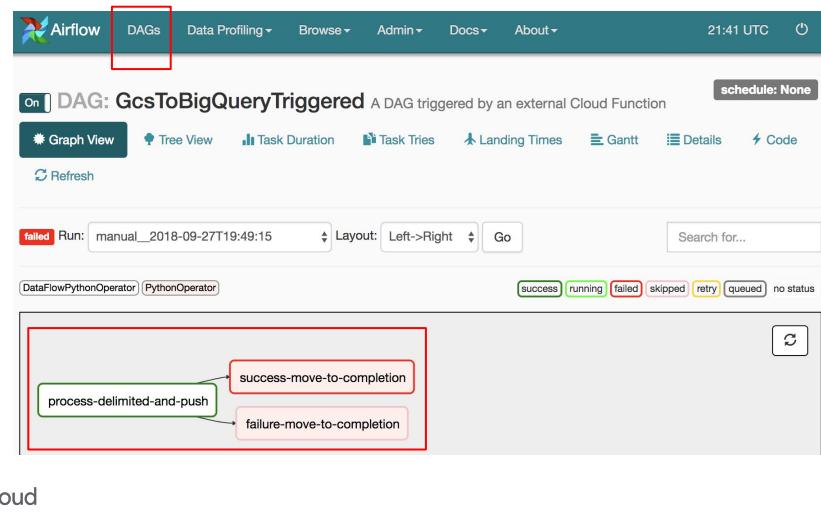
```
105
106
107     with models.DAG(dag_id='gcsToBigQueryTriggered',
108                      description='A DAG triggered by an external Cloud Function',
109                      schedule_interval=None, default_args=DEFAULT_DAG_ARGS) as dag:
110
111     # A step required for the Dataflow job.
112     job_args = [
113         'input': 'gs://{{ dag_run.conf["bucket"] }}/{{ dag_run.conf["name"] }}',
114         'output': models.Variable.get('bg_output_table'),
115         'fields': models.Variable.get('input_field_names'),
116         'load_dt': DS_TAG
117     ]
118
119     # Main Dataflow task that will process and load the input delimited file.
120     dataflow_task = dataflow_operator.DataFlowPythonOperator(
121         task_id='process-delimited-and-push',
122         py_file=DATAFLOW_FILE,
123         options=job_args)
124
125     # Here we create two conditional tasks, one of which will be executed
126     # based on whether the Dataflow task was a success or a failure.
127     success_move_task = python_operator.PythonOperator(task_id='success-move-to-completion',
128                                                       python_callable=move_to_completion_bucket,
129                                                       # A success_tag is used to move
130                                                       # the input file to a success
131                                                       # directory folder.
132                                                       op_args=[COMPLETION_BUCKET, SUCCESS_TAG],
133                                                       provide_context=True,
134                                                       trigger_rule=TriggerRule.ALL_SUCCESS)
135
136     failure_move_task = python_operator.PythonOperator(task_id='failure-move-to-completion',
137                                                       python_callable=move_to_completion_bucket,
138                                                       # A failure_tag is used to move
139                                                       # the input file to a failure
140                                                       # directory folder.
141                                                       op_args=[COMPLETION_BUCKET, FAILURE_TAG],
142                                                       provide_context=True,
143                                                       trigger_rule=TriggerRule.ALL_FAILED)
144
145     dataflow_task >> success_move_task
146     dataflow_task >> failure_move_task
147
```



Like this task which uses the DataFlowPythonOperator and is given the task_id of “process-delimited-and-push”.

We'll go over creating a DAG file and its components a little later

The python file creates a DAG



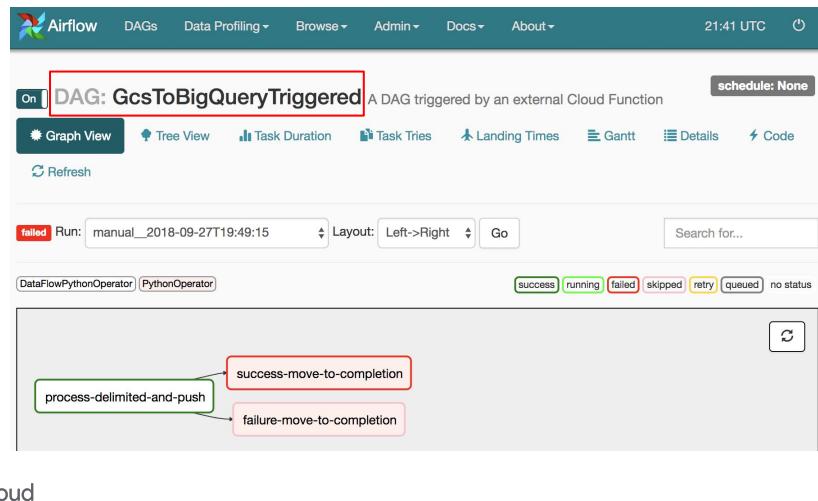
Once you've uploaded the python file to the DAGs folder, you can navigate back to the Airflow webserver and under DAGs you'll see the DAG you created with code represented visually as a directed graph (nodes and edges).

You'll remember that the python code that defined a task we called "process-delimited-and-push" is now a node in our graph here.

Let's explore a bit more of the Airflow web UI.

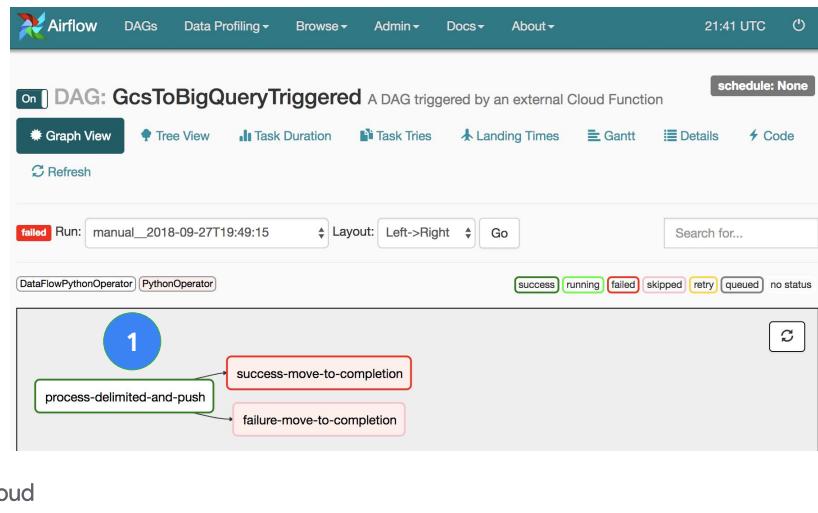
// Great DAG tutorial: <https://airflow.apache.org/tutorial.html>

Airflow webserver UI overview



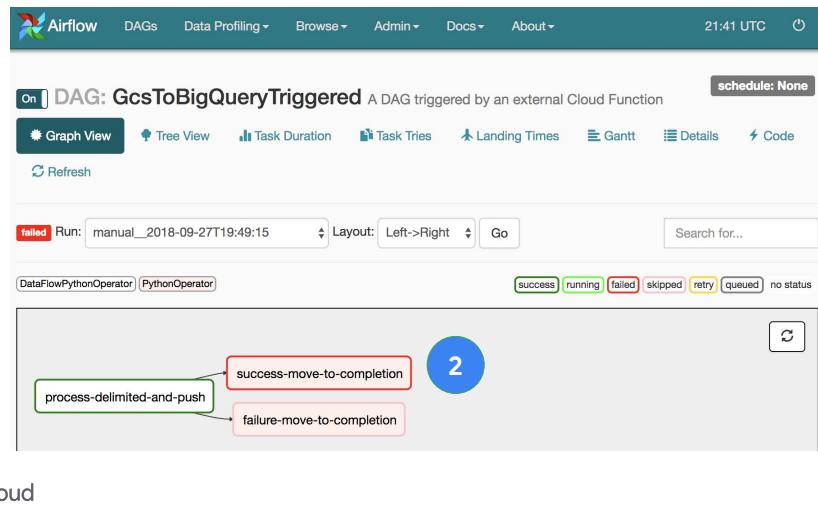
You can see that this particular workflow is called GcsToBigQueryTriggered and it has three TASKS when it runs.

The python file creates a DAG (Directed Acyclic Graph)



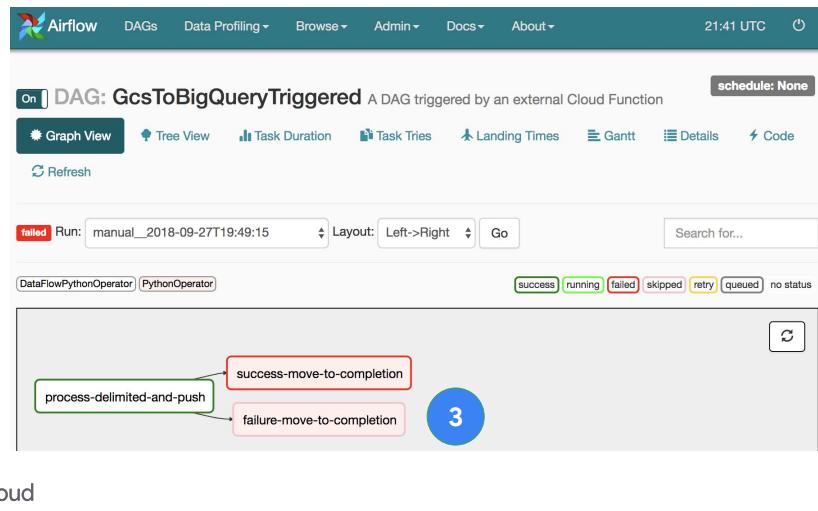
- 1) process-delimited-and-push which I just happen to know from that python file you saw earlier that this task invokes a Dataflow job to read in a new CSV file from a GCS bucket, processes it, and writes the output to BigQuery

The python file creates a DAG (Directed Acyclic Graph)



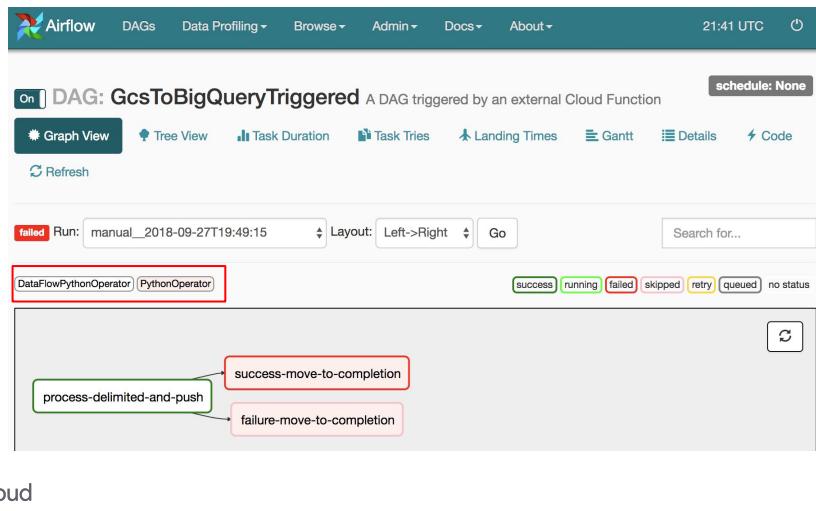
- 2) success-move-to-completion which moves the CSV file from an input GCS bucket to a processed or completed GCS bucket for archiving

The python file creates a DAG (Directed Acyclic Graph)



or 3) if the pipeline fails part way the file is moved to the completion bucket but tagged as failure. This is an example of a DAG which isn't strictly sequential. There a decision made to run one node or a different one based on the outcome of the parent node.

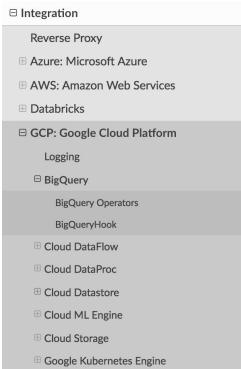
The python file creates a DAG (Directed Acyclic Graph)



But regardless of the size and shape of your workflow DAG, one common thread for all workflows is the common Operators used. If the DAG itself is HOW to run the workflow (first do step 1 then either move to step 2 or 3) the operators specify WHAT actually gets done as part of the task.

In this simple example we're calling on the DataFlowPythonOperator and the general PythonOperator. Those are by no means the only operators so let's pause here and look at all the operators at our disposal to achieve our goal of automatic retraining and deployment of our ML model.

Airflow uses operators in your DAG to orchestrate other GCP services



View all GCP services that Airflow can orchestrate:

<http://airflow.apache.org/integration.html#gcp>

+ sample code



Airflow has many operators which you can invoke the TASKS you want to complete. Operators are usually atomic in a task (which means generally you only see one operator per task). I've taken this list directly from the Apache Airflow docs of all services that Airflow can orchestrate to. Let's take a look at the ones that are most relevant to us as ML Engineers.

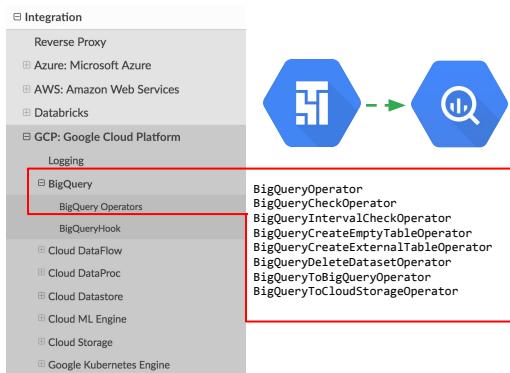
List of all operators:

https://cloud.google.com/composer/docs/how-to/using/writing-dags#gcp_name_operators

<https://airflow.apache.org/integration.html#gcp>

<https://airflow.apache.org/concepts.html#operators>

BigQuery Operators are popular for updating your ML training dataset



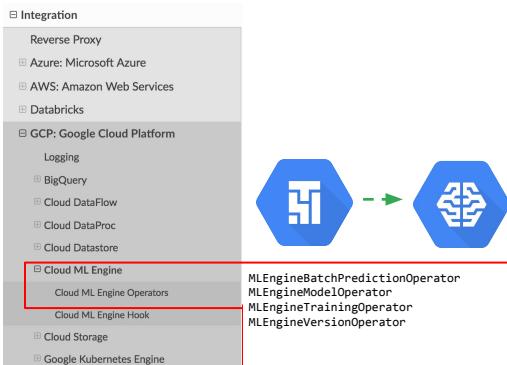
As you might have guessed, we'll certainly be making use of the BigQuery operators since our workflows live and die by the data that is fed into them through GCS and BigQuery. Here's a list of the specific operators that we can invoke in a task to call on the BigQuery service for querying and other data-related tasks. You'll be mainly working with the first three in this course but I encourage you to skim the resource link on all the operators so you can get a feel for what is possible.

List of all operators:

https://cloud.google.com/composer/docs/how-to/using/writing-dags#gcp_name_operators

<https://airflow.apache.org/integration.html#gcp>

Cloud AI Platform (formerly Cloud ML Engine) operators can launch training and deployment jobs



<https://airflow.apache.org/integration.html#cloud-ml-engine-operators>



Once we have our training data in a good place, the next logical step in our workflow is to retrain and redeploy our model. In the same DAG file after the BigQuery operators complete we can make a service call through a Cloud ML Engine operator to kick off a new training job and manage our model like incrementing the version.

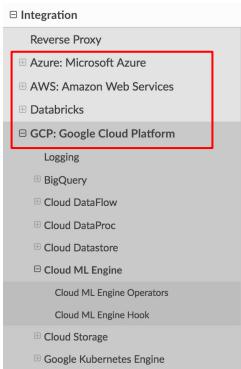
<https://airflow.apache.org/integration.html#cloud-ml-engine-operators>

List of all operators:

https://cloud.google.com/composer/docs/how-to/using/writing-dags#gcp_name_operators

<https://airflow.apache.org/integration.html#gcp>

Apache Airflow is open source and cross-platform for hybrid pipelines



<https://airflow.apache.org/integration.html#cloud-ml-engine-operators>



You might have noticed that your Airflow DAG can have operators that send tasks out to other cloud providers. This is great for hybrid workflows where you have components across multiple cloud platforms or even on-premise. Apache Airflow is open source and continually adds more operators to other services so be sure to check out the list in the documentation periodically if you're waiting for a new service to be added.

<https://airflow.apache.org/integration.html#cloud-ml-engine-operators>

Example: Common Machine Learning workflow DAG Operators

```
bq_rec_training_data → bq_export_op → ml_engine_training_op → app_engine_deploy_version
```

```
# update training data  
t1 = BigQueryOperator(  
)  
  
# BigQuery training data export to GCS  
t2 = BigQueryToCloudStorageOperator(  
)  
  
# AI Platform training job  
t3 = MLEngineTrainingOperator(  
)  
  
# App Engine deploy new version  
t4 = AppEngineVersionOperator(  
)  
  
# DAG dependencies  
t2.set_upstream(t1)  
t3.set_upstream(t2)  
t4.set_upstream(t3)
```



Here you see four tasks t1, t2, t3, and t4 and four operators corresponding to four Google Cloud Platform services. The names should look familiar and you can probably start to guess what this pipeline does at a high level just by reading them in order.

BigQuery and Cloud Storage operators get us fresh training data



```
# update training data  
t1 = BigQueryOperator(  
)  
  
# BigQuery training data export to GCS  
t2 = BigQueryToCloudStorageOperator(  
)  
  
# AI Platform training job  
t3 = MLEngineTrainingOperator(  
)  
  
# App Engine deploy new version  
t4 = AppEngineVersionOperator(  
)  
  
# DAG dependencies  
t2.set_upstream(t1)  
t3.set_upstream(t2)  
t4.set_upstream(t3)
```



The first two are concerned with *getting fresh model data* from a BigQuery dataset and into Cloud Storage for consumption by our ML model later in the pipeline.

In the lab you're going to work on later the dataset will be the one you're already familiar with, the Google Analytics news articles sample dataset.

Let's see the parameters the BigQueryOperator takes.

Use the BigQueryOperator to run SQL

```
from airflow.contrib.operators import bigquery_operator

# constants or can be dynamic based on Airflow macros
max_query_date = '2018-02-01'
min_query_date = '2018-01-01'

    # Query recent StackOverflow questions.
bq_recent_questions_query = bigquery_operator.BigQueryOperator(

    task_id='bq_recent_questions_query',

    bql="""
SELECT owner.display_name, title, view_count
FROM `bigquery-public-data.stackoverflow.posts_questions`
WHERE creation_date < CAST('{max_date}' AS TIMESTAMP)
    AND creation_date >= CAST('{min_date}' AS TIMESTAMP)
ORDER BY view_count DESC
LIMIT 100
""".format(max_date=max_query_date, min_date=min_query_date),

    use_legacy_sql=False,

    destination_dataset_table=bq_recent_questions_table_id)
```



The BigQueryOperator allows you to specify a SQL query to run against a BigQuery dataset.

BigQuery Operator

https://cloud.google.com/composer/docs/how-to/using/writing-dags#gcp_name_operators

<https://airflow.apache.org/integration.html#bigqueryoperator>

StackOverflow email notifier:

https://github.com/GoogleCloudPlatform/python-docs-samples/blob/master/composer/workflows/bq_notify.py

SQL commands can hold parameters (from Python)

```
from airflow.contrib.operators import bigquery_operator  
  
# constants or can be dynamic based on Airflow macros  
max_query_date = '2018-02-01'  
min_query_date = '2018-01-01'  
  
# Query recent StackOverflow questions.  
bq_recent_questions_query = bigquery_operator.BigQueryOperator(  
  
    task_id='bq_recent_questions_query',  
  
    bqj=""  
    SELECT owner.display_name, title, view_count  
    FROM `bigquery-public-data.stackoverflow.posts_questions`  
    WHERE creation_date < CAST('{max_date}' AS TIMESTAMP)  
        AND creation_date >= CAST('{min_date}' AS TIMESTAMP)  
    ORDER BY view_count DESC  
    LIMIT 100  
    """.format(max_date=max_query_date, min_date=min_query_date),  
  
    use_legacy_sql=False,  
  
    destination_dataset_table=bq_recent_questions_table_id)
```



In this quick example we're passing in a query which returns the top 100 most popular Stackoverflow posts from the BigQuery public dataset for a specified date range you see there in the WHERE clause. Notice anything different about the filters in the WHERE clause?

Yes! they are parameters in this case for max_date and min_date

Note the Python constants for a date range here

```
from airflow.contrib.operators import bigquery_operator  
  
# constants or can be dynamic based on Airflow macros  
max_query_date = '2018-02-01'  
min_query_date = '2018-01-01'  
  
# Query recent StackOverflow questions.  
bq_recent_questions_query = bigquery_operator.BigQueryOperator(  
  
    task_id='bq_recent_questions_query',  
  
    bqj=""  
    SELECT owner.display_name, title, view_count  
    FROM `bigquery-public-data.stackoverflow.posts_questions`  
    WHERE creation_date < CAST('{max_date}' AS TIMESTAMP)  
        AND creation_date >= CAST('{min_date}' AS TIMESTAMP)  
    ORDER BY view_count DESC  
    LIMIT 100  
    """.format(max_date=max_query_date, min_date=min_query_date),  
  
    use_legacy_sql=False,  
  
    destination_dataset_table=bq_recent_questions_table_id)
```



You can parameterize pieces of the SQL statement like what we did here to only return posts for January 2018 with `min_query_date` and `max_query_date`.

You could even set a rolling window with macros

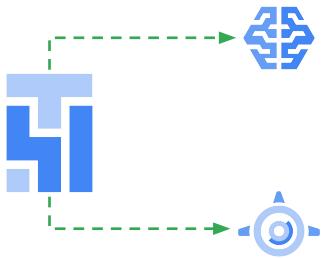
```
from airflow.contrib.operators import bigquery_operator
# constants or can be dynamic based on Airflow macros
max_query_date = {{ macros.ds_add(ds, -1) }} # one day prior
min_query_date = {{ macros.ds_add(ds, -7) }} # seven days prior

# Query recent StackOverflow questions.
bq_recent_questions_query = bigquery_operator.BigQueryOperator(
    task_id='bq_recent_questions_query',
    bqsql="""
        SELECT owner.display_name, title, view_count
        FROM `bigquery-public-data.stackoverflow.posts_questions`
        WHERE creation_date < CAST('{{max_date}}' AS TIMESTAMP)
            AND creation_date >= CAST('{{min_date}}' AS TIMESTAMP)
        ORDER BY view_count DESC
        LIMIT 100
    """.format(max_date=max_query_date, min_date=min_query_date),
    use_legacy_sql=False,
    destination_dataset_table=bq_recent_questions_table_id)
```



What's really neat is you can even make the parameters dynamic (instead of the static ones shown here) and have them be based on the DAG schedule date like `{{ macros.ds_add(ds, -7) }}` which is a week before the DAG scheduled run date.

AI Platform and App Engine operators retrain and redeploy our model

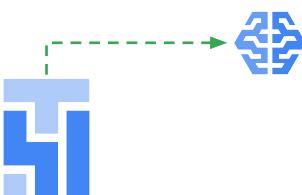


```
# update training data  
t1 = BigQueryOperator(  
)  
  
# BigQuery training data export to GCS  
t2 = BigQueryToCloudStorageOperator(  
)  
  
# AI Platform training job  
t3 = MLEngineTrainingOperator(  
)  
  
# App Engine deploy new version  
t4 = AppEngineVersionOperator(  
)  
  
# DAG dependencies  
t2.set_upstream(t1)  
t3.set_upstream(t2)  
t4.set_upstream(t3)
```



The next two operators handle re-training the model by submitting a job to Cloud Machine Learning Engine and then deploying the updated model to App Engine.

Use Cloud ML Engine operators to periodically submit new training jobs



```
t3 = MLEngineTrainingOperator(  
    task_id='ml_engine_training_op',  
    project_id=PROJECT_ID,  
    job_id=job_id,  
    package_uris=[PACKAGE_URI],  
    training_python_module='trainer.task',  
    training_args=training_args,  
    region=REGION,  
    scale_tier='CUSTOM',  
    master_type='complex_model_m_gpu',  
    dag=dag  
)
```



The `MLEngineTrainingOperator` is how Cloud Composer interacts with Cloud Machine Learning Engine and thus gives you the ability to periodically schedule new jobs to be sent to CMLE as part of your automated workflow. Take a look through the parameters you provide to the operator.

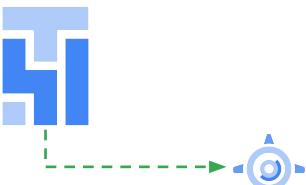
First we specify an ID for the task that we're running and then the project we want to run it on. For our lab, this will be the same project that is also managing the Airflow Instance and everything else. Then you have the option of creating your own job id, which usually involves a concatenated timestamp or similar unique identifier. After that is path to the actual model code which we'll have in a zipped file which is created by a shell script after training is run in our lab.

Then the actual Python module name within your code to run within the `MLEngine` training job after installing the package. In this example we've called it `trainer.task`

Next is a series of arguments for training which we store as an array called `"training_args"` which includes the location of the job directory, where the training files are, where the output directory is, and what data we're using. The next three: `region`, `scale_tier`, and `master_type` are GCP infrastructure parameters for where you want the job to be ran and with what speciality hardware like GPUs or TPUs if any.

https://airflow.apache.org/_modules/airflow/contrib/operators/mlengine_operator.html

Use App Engine operators to periodically deploy and redeploy models



```
t4 = AppEngineVersionOperator(  
    task_id='app_engine_deploy_version',  
    project_id=PROJECT_ID,  
    service_id='default',  
    region=REGION,  
    service_spec=None,  
    dag=dag  
)
```



Lastly, once your model is retrained and ready to be delivered as an API endpoint for serving, we need to redeploy our app engine project with the latest model. Here the parameters are simple since we're just redeploying our existing app engine instance associated with our project. And that's it! Well almost. That's the last TASK but at the end of most all DAG files you'll find

Manage pipelines and dependencies as code

```
# update training data
t1 = BigQueryOperator(
)

# BigQuery training data export to GCS
t2 = BigQueryToCloudStorageOperator(
)

# AI Platform training job
t3 = MLEngineTrainingOperator(
)

# App Engine deploy new version
t4 = AppEngineVersionOperator(
)

# DAG dependencies
t2.set_upstream(t1)
t3.set_upstream(t2)
t4.set_upstream(t3)
```



The actual order in which we want these operators to run. This is what the D in DAG is for: Directed. For our example, t2 or task 2 won't run until t1 has completed. This is what gives our graph the dependencies of a workflow. I can probably guess what you're thinking at this point, you could build some cool branching of multiple child nodes per upstream parent and that is totally possible. Just don't forget to comment your code so you know where one branch begins, where it's going, and what tasks are involved. As a tip, after you load your DAG file into the DAG folder you can see the visualization of your DAG in the airflow UI as a directed graph, a gantt chart, or a list if you wanted. Reviewing the visual representation of the ordered tasks will help you confirm your tasks are ordered properly.

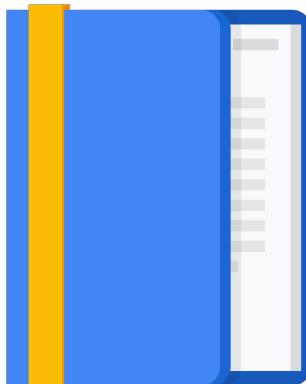
Agenda

Building Batch Data Pipelines visually with Cloud Data Fusion

- Components
- UI Overview
- Building a Pipeline
- Exploring Data using Wrangler

Orchestrating work between GCP services with Cloud Composer

- Apache Airflow Environment
- DAGs and Operators
- Workflow Scheduling
- Monitoring and Logging



Now that you're familiar with the Cloud Composer and Apache Airflow environments AND the basics of building a list of tasks for GCP services in your DAG, it is time to discuss a really important topic: workflow scheduling.

Two scheduling options for Cloud Composer workflows



Periodic

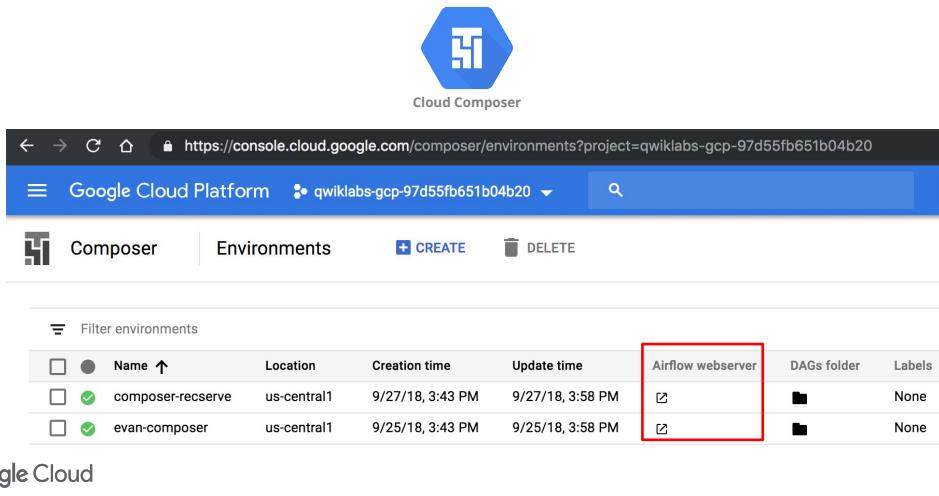


Event-driven



As we hinted at earlier, there are two different ways your workflow can be ran without you sitting there manually clicking “run DAG”. The first and most common is a set schedule or periodic run of a workflow like once a day at 6am or weekly on Saturdays. The second way is trigger-based like if you wanted to run your workflow whenever new CSV data files were loaded into a GCS bucket or if new data came in from a PUB/SUB topic you’ve subscribed to.

Launch the Airflow webserver to interact with your DAGs



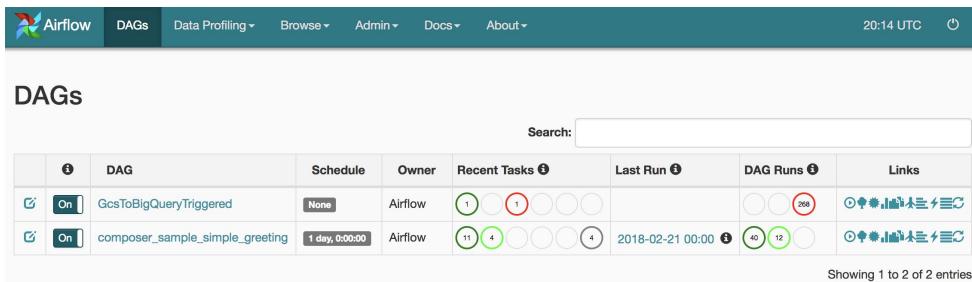
The screenshot shows the Google Cloud Platform interface for Cloud Composer environments. At the top, there's a navigation bar with the URL <https://console.cloud.google.com/composer/environments?project=qwiklabs-gcp-97d55fb651b04b20>. Below it is a blue header bar with the text "Google Cloud Platform" and "qwiklabs-gcp-97d55fb651b04b20". The main content area has a title "Cloud Composer" with a logo. Below the title, there are tabs for "Composer" and "Environments", and buttons for "CREATE" and "DELETE". A search bar is also present. The main table lists environments:

Name	Location	Creation time	Update time	Airflow webserver	DAGs folder	Labels
composer-recserve	us-central1	9/27/18, 3:43 PM	9/27/18, 3:58 PM	View	Edit	None
evan-composer	us-central1	9/25/18, 3:43 PM	9/25/18, 3:58 PM	View	Edit	None

The last two columns, "Airflow webserver" and "Labels", are highlighted with a red box. At the bottom left, there's a "Google Cloud" logo.

To view the schedules for your DAGs you first launch the Airflow webserver from within Cloud Composer. This is a great link to bookmark.

Airflow scheduling basics



The screenshot shows the Airflow web interface with the 'DAGs' tab selected. The top navigation bar includes links for Airflow, DAGs, Data Profiling, Browse, Admin, Docs, About, and a power icon. The time is listed as 20:14 UTC. The main content area is titled 'DAGs' and contains a search bar. A table lists two DAG entries:

	DAG	Schedule	Owner	Recent Tasks	Last Run	DAG Runs	Links
<input checked="" type="checkbox"/>	On GcsToBigQueryTriggered	None	Airflow	1 1 1 1 1 1		268	View Details
<input checked="" type="checkbox"/>	On composer_sample_simple_greeting	1 day, 0:00:00	Airflow	11 4 1 1 1 1 4	2018-02-21 00:00	40 12	View Details

Below the table, a message indicates "Showing 1 to 2 of 2 entries".



Then navigate to the DAGs tab to view the existing workflows that you have python DAG files for. Here we have two DAGs. The bottom one, `composer_sample_simple_greeting` has a daily schedule but..

Why is this DAG missing a schedule?

DAGs							Search:
	DAG	Schedule	Owner	Recent Tasks	Last Run	DAG Runs	Links
	GcsToBigQueryTriggered		Airflow				
	composer_sample_simple_greeting		Airflow		2018-02-21 00:00		



..why is this top DAG missing a schedule? How would it ever get run?

Option 1: Event-driven scheduling with Cloud Functions

DAGs

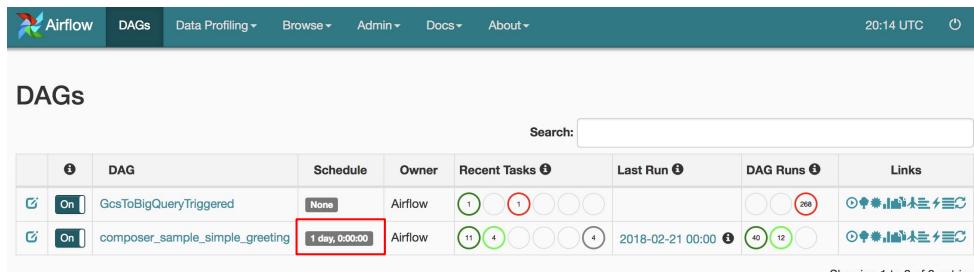
	DAG	Schedule	Owner	Recent Tasks	Last Run	DAG Runs	Links
<input checked="" type="checkbox"/>	On GcsToBigQueryTriggered	None	Airflow	1 1		268	
<input checked="" type="checkbox"/>	On composer_sample_simple_greeting	1 day, 0:00:00	Airflow	11 4	2018-02-21 00:00	40 12	

Showing 1 to 2 of 2 entries



The answer is the fact that it's not on a set schedule at all. It's event driven. The driver of when this workflow runs is a Cloud Function that we create. In the next lesson, we'll actually create our own Cloud Function that watches a GCS bucket for new CSV files.

Option 2: Specify pipeline schedule_interval in your DAG



	DAG	Schedule	Owner	Recent Tasks	Last Run	DAG Runs	Links
<input checked="" type="checkbox"/>	GcsToBigQueryTriggered	None	Airflow	1 1 1 1 1		268	
<input checked="" type="checkbox"/>	composer_sample_simple_greeting	1 day, 0:00:00	Airflow	11 4 1 1 1 4	2018-02-21 00:00	40 12	

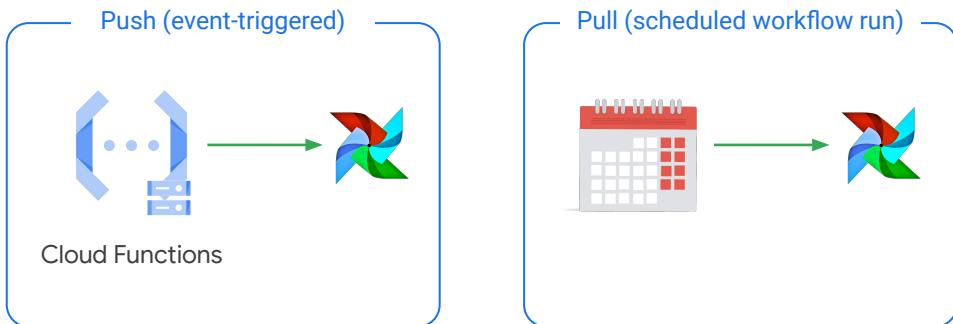
Showing 1 to 2 of 2 entries

```
with models.DAG(  
    'composer_sample_simple_greeting',  
    schedule_interval=datetime.timedelta(days=1),  
    default_args=default_dag_args) as dag:
```



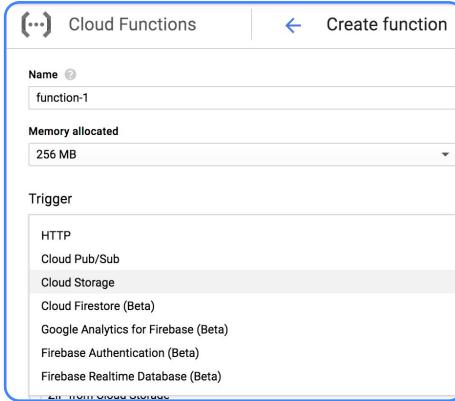
If you wanted to go the regular schedule route you can specify the `schedule_interval` in your DAG code like what you see here. By the way, clicking on the schedule of 1 day here in the UI won't allow you to edit it there but instead will take you to the history of all the runs for that workflow.

Two types of workflow ETL patterns



As you saw earlier in this course, there are two general patterns for ETL workflows. Event triggered or push (as-in I push a new file to Cloud Storage and my workflow kicks off) or PULL which is where Airflow at a set time could look in my Cloud Storage folder and take all the contents that are found for it's workflow run.

Use Cloud Functions to create an event-based push architecture

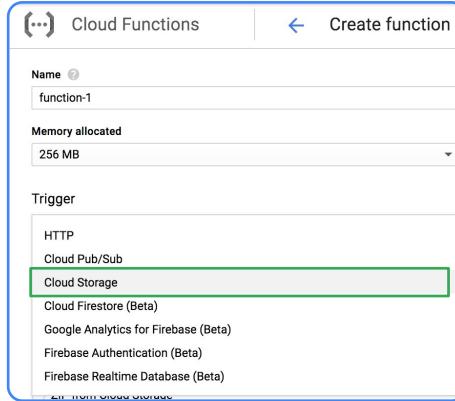


We can use Cloud Functions to create our event-driven or push architecture workflow. I mentioned triggering on events within a GCS bucket but you can also trigger based on HTTP requests, Pub/sub, Firestore, Firebase and more as you see here.

Generally, push technology is GREAT when wanting to distribute transactions as-they-happen. Stock tickers, and other types of financial institution transactions are very important when it comes to push technology. How about disasters and notification? Again, important.

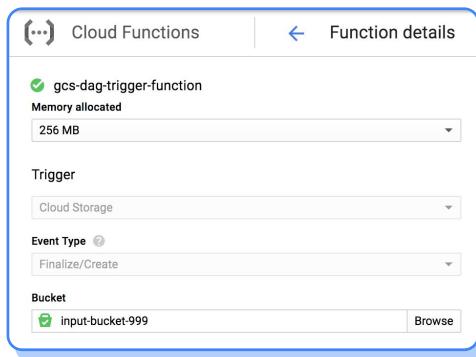
For ML workflows where your upstream data doesn't arrive at a regular pace (like get all the transactions at the end of each day) consider experimenting with a push architecture. Your final lab, since it's based on regular Google Analytics news article data will be a pull architecture but I've added in an optional lab for you to get practice with Cloud Functions and event-driven workflows for those interested so let's talk through it more now.

Example: Trigger an event when new data is loaded to Cloud Storage



For our example let's assume we have a CSV file or set of files loaded to Cloud Storage so we'll choose a Cloud Storage trigger for our function.

Creating a Cloud Function to trigger an Airflow DAG



Then we specify an Event Type (Finalize/Create new files) and a bucket to watch.

<https://cloud.google.com/composer/docs/how-to/using/triggering-with-gcf>

Creating a Cloud Function GCS trigger an Airflow DAG



Cloud Functions

```
index.js  package.json

14  * @param {Object} event The Cloud Functions event.
15  * @param {Function} callback The callback function.
16  */
17 exports.triggerDag = function triggerDag (event, callback) {
18   // Fill in your Composer environment information here.
19
20   // The project that holds your function
21   const PROJECT_ID = 'wikilabs-gcp';
22   // Navigate to your webserver's login page and get this from the URL
23   const CLIENT_ID = '';
24   // This should be part of your webserver's URL:
25   const WEBSERVER_ID = 'b9...-tp';
26   // The name of the DAG you wish to trigger
27   const DAG_NAME = 'GcsToBigQueryTriggered';
28
29   // Other constants
30   const WEBSERVER_URL = `https://${WEBSERVER_ID}.appspot.com/api/experimental/dags/${DAG_NAME}/dag`;
31   const USER_AGENT = 'gcf-event-trigger';
32   const BODY = {'conf': JSON.stringify(event.data)};
33
34   // Make the request
35   authorizeIap(CLIENT_ID, PROJECT_ID, USER_AGENT)
36     .then(function iapAuthorizationCallback (iap) {
37       makeIapPostRequest(WEBSERVER_URL, BODY, iap.idToken, USER_AGENT, iap.jwt);
38     })
39     .then(_ => callback(null))
40     .catch(callback);
41   };
42 };
43
```



As part of the Cloud Function, we need to well... create an actual function is javascript that we want called. The good news is most of this code for triggering Airflow DAGs in a function is all boilerplate for you to copy from as a starting point.

Creating a Cloud Function GCS trigger an Airflow DAG



Cloud Functions

```
index.js  package.json

14 * @param {Object} event The Cloud Functions event.
15 * @param {Function} callback The callback function.
16 */
17 exports.triggerDag = function triggerDag(event, callback) {
18 // Fill in your Composer environment information here.
19
20 // The project that holds your function
21 const PROJECT_ID = 'wikilabs-gcp';
22 // Navigate to your webserver's login page and get this from the URL
23 const CLIENT_ID = '';
24 // This should be part of your webserver's URL:
25 const WEBSERVER_ID = 'b9...-tp';
26 // The name of the DAG you wish to trigger
27 const DAG_NAME = 'GcsToBigQueryTriggered';
28
29 // Other constants
30 const WEBSERVER_URL = `https://${WEBSERVER_ID}.appspot.com/api/experimental/dags/${DAG_NAME}/dag`;
31 const USER_AGENT = 'gcf-event-trigger';
32 const BODY = {'conf': JSON.stringify(event.data)};
33
34 // Make the request
35 authorizeIap(CLIENT_ID, PROJECT_ID, USER_AGENT)
36 .then(function iapAuthorizationCallback (iap) {
37   makeIapPostRequest(WEBSERVER_URL, BODY, iap.idToken, USER_AGENT, iap.jwt);
38 })
39 .then(_ => callback(null))
40 .catch(callback);
41 };
42
43;
```



Here we specify a name for our function called triggerDag

Creating a Cloud Function GCS trigger an Airflow DAG



Cloud Functions

```
index.js  package.json

14  * @param {Object} event The Cloud Functions event.
15  * @param {Function} callback The callback function.
16  */
17 exports.triggerDag = function triggerDag (event, callback) {
18   // Fill in your Composer environment information here.
19
20   // The project that holds your function
21   const PROJECT_ID = 'wikilabs-gcp';
22   // Navigate to your webserver's login page and get this from the URL
23   const CLIENT_ID = '1054011111111-000000000000.apps.googleusercontent.com';
24   // This should be part of your webserver's URL:
25   // (tenant-project-id).appspot.com
26   const WEBSERVER_ID = 'b9...-tp';
27   // The name of the DAG you wish to trigger
28   const DAG_NAME = 'GcsToBigQueryTriggered';
29
30   // Other constants
31   const WEBSERVER_URL = `https://${WEBSERVER_ID}.appspot.com/api/experimental/dags/${DAG_NAME}/dag`;
32   const USER_AGENT = 'gcf-event-trigger';
33   const BODY = {'conf': JSON.stringify(event.data)};
34
35   // Make the request
36   authorizeIap(CLIENT_ID, PROJECT_ID, USER_AGENT)
37     .then(function iapAuthorizationCallback (iap) {
38       makeIapPostRequest(WEBSERVER_URL, BODY, iap.idToken, USER_AGENT, iap.jwt);
39     })
40     .then(_ => callback(null))
41     .catch(callback);
42 };
43
```



Then we tell it where your Airflow environment is to be triggered and which DAG in that Airflow environment. In this case it's looking for one called GcsToBigQueryTriggered

Creating a Cloud Function GCS trigger an Airflow DAG



Cloud Functions

```
index.js  package.json

14 * @param {Object} event The Cloud Functions event.
15 * @param {Function} callback The callback function.
16 */
17 exports.triggerDag = function triggerDag (event, callback) {
18 // Fill in your Composer environment information here.
19
20 // The project that holds your function
21 const PROJECT_ID = 'wikilabs-gcp';
22 // Navigate to your webserver's login page and get this from the URL
23 const CLIENT_ID = '';
24 // This should be part of your webserver's URL:
25 // (tenant-project-id).appspot.com
26 const WEBSERVER_ID = 'b9';
27 // The name of the DAG you wish to trigger
28 const DAG_NAME = 'GcsToBigQueryTriggered';
29
30 // Other constants
31 const WEBSERVER_URL = `https://${WEBSERVER_ID}.appspot.com/api/experimental/dags/${DAG_NAME}/dag`;
32 const USER_AGENT = 'gcf-event-trigger';
33 const BODY = {'conf': JSON.stringify(event.data)};
34
35 // Make the request
36 authorizeIap(CLIENT_ID, PROJECT_ID, USER_AGENT)
37 .then(function iapAuthorizationCallback (iap) {
38   makeIapPostRequest(WEBSERVER_URL, BODY, iap.idToken, USER_AGENT, iap.jwt);
39 })
40 .then(_ => callback(null))
41 .catch(callback);
42 };
43 }
```



Keep in mind you can have multiple workflows or DAGs in a single Airflow environment so be sure you specify the correct DAG_NAME to trigger!

Creating a Cloud Function GCS trigger an Airflow DAG



Google Cloud

Then we have a few constants that are provided which construct the Airflow URL that we're going to trigger a POST request to as well as who's making the request and what is the body of the request is.

Creating a Cloud Function GCS trigger an Airflow DAG



Cloud Functions

```
index.js  package.json

14  * @param {Object} event The Cloud Functions event.
15  * @param {Function} callback The callback function.
16  */
17 exports.triggerDag = function triggerDag (event, callback) {
18   // Fill in your Composer environment information here.
19
20   // The project that holds your function
21   const PROJECT_ID = 'wikilabs-gcp';
22   // Navigate to your webserver's login page and get this from the URL
23   const CLIENT_ID = '';
24   // This should be part of your webserver's URL:
25   const WEBSERVER_ID = 'b9...-tp';
26   // The name of the DAG you wish to trigger
27   const DAG_NAME = 'GcsIosigQueryTriggered';
28
29   // Other constants
30   const WEBSERVER_URL = `https://${WEBSERVER_ID}.appspot.com/api/experimental/dags/${DAG_NAME}/dag`;
31   const USER_AGENT = 'gcf-event-trigger';
32   const BODY = {'conf': JSON.stringify(event.data)};
33
34
35   // Make the request
36   authorizeIap(CLIENT_ID, PROJECT_ID, USER_AGENT)
37     .then(function iapAuthorizationCallback (iap) {
38       makeIapPostRequest(WEBSERVER_URL, BODY, iap.idToken, USER_AGENT, iap.jwt);
39     })
40     .then(_ => callback(null))
41     .catch(callback);
42 }
43
```



Lastly, the triggerDag function makes the actual request against the Airflow server to kick-off a workflow DAG.

Cloud Functions Function details

ZIP from Cloud Storage Cloud Source repository

Runtime: Node.js 6

index.js package.json

```
1 'use strict';
2
3 const fetch = require('node-fetch');
4 const FormData = require('form-data');
5
6 /**
7 * Triggered from a message on a Cloud Storage bucket.
8 * IAP authorization based on:
9 * https://stackoverflow.com/questions/45787676/how-to-authenticate-a-cloud-storage-trigger-without-a-key
10 * https://cloud.google.com/iap/docs/authentication-howto
11 * https://cloud.google.com/functions/authn
12 */
13 * @param {Object} event The Cloud Functions event.
14 * @param {Function} callback The callback function.
15 */
16
17 exports.triggerDag = function triggerDag(event, callback) {
18 // Fill in your Configuration information here.
19
20 // The project that holds your function.
21 const PROJECT_ID = 'qwiksite-gcp-37d85fb651b04b20';
22 // Navigate to https://console.cloud.google.com/project/<PROJECT_ID>/logins and get this:
23 const CLIENT_ID = '941500678485-gde61d87qtdm91t17809uj';
24 // This should be part of your webserver's URL.
25 }
```

Function to execute triggerDag

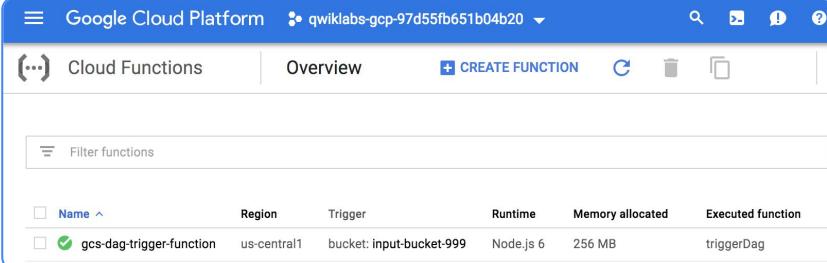
Google Cloud

Be sure to specify the function you want Cloud Functions to call in your script

Once you've got the Cloud Function code ready in your index.js file and the metadata about the function in package.json (which contains code dependency and versioning information) you still need to specify which function you actually want executed. In this case we created one called triggerDag so we just copy that down.

I'll also save you about 20 minutes of frustration and tell you that the 'function to execute' box is case sensitive so all capital letters DAG is different than capital D lowercase a and g.

Cloud Function is triggered whenever a new file is loaded to a specific Cloud Storage bucket



The screenshot shows the Google Cloud Platform Cloud Functions Overview page. At the top, there is a navigation bar with the Google Cloud logo, the text "Google Cloud Platform", the project name "qwiklabs-gcp-97d55fb651b04b20", and standard navigation icons (search, refresh, help). Below the navigation bar, there are tabs for "Cloud Functions" (selected), "Overview", and a "CREATE FUNCTION" button. To the right of these tabs are icons for creating, deleting, and cloning functions. A search bar labeled "Filter functions" is present. The main area displays a table of functions. The table has columns: "Name" (with a sorting arrow), "Region", "Trigger", "Runtime", "Memory allocated", and "Executed function". One row is visible, showing the function "gcs-dag-trigger-function" in the "Name" column, "us-central1" in "Region", "bucket: input-bucket-999" in "Trigger", "Node.js 6" in "Runtime", "256 MB" in "Memory allocated", and "triggerDag" in "Executed function". The "Name" column header has an upward arrow indicating it is sorted in descending order.

<input type="checkbox"/> Name ^	Region	Trigger	Runtime	Memory allocated	Executed function
<input checked="" type="checkbox"/> gcs-dag-trigger-function	us-central1	bucket: input-bucket-999	Node.js 6	256 MB	triggerDag



And there you have it! Your Cloud Function has been created and is actively watching your Cloud Storage bucket for file uploads. But how can you be sure everything is working as intended? For that, check out the next topic on monitoring and logging.

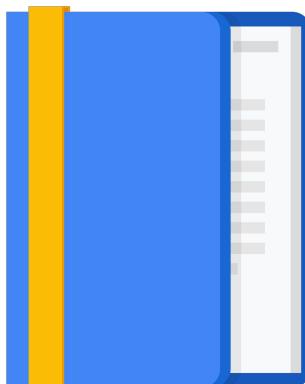
Agenda

Building Batch Data Pipelines visually with Cloud Data Fusion

- Components
- UI Overview
- Building a Pipeline
- Exploring Data using Wrangler

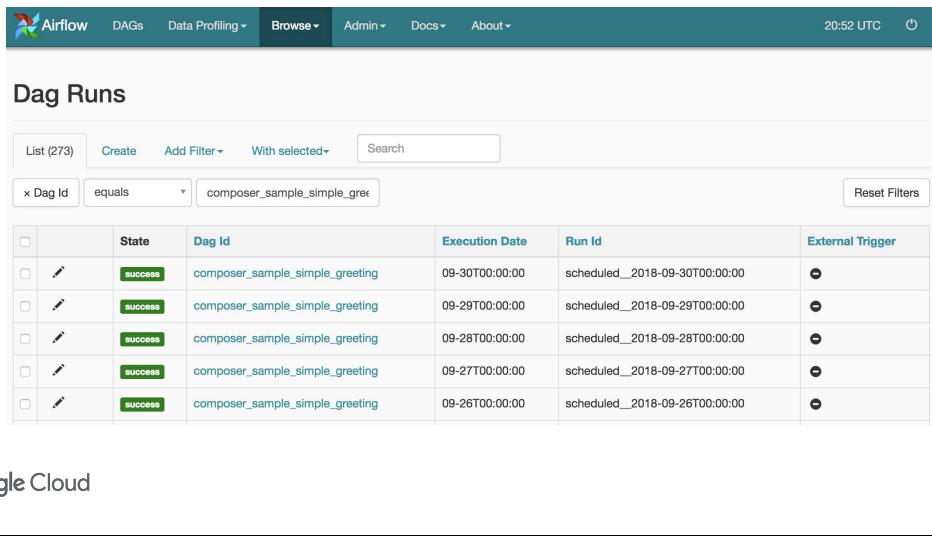
Orchestrating work between GCP services with Cloud Composer

- Apache Airflow Environment
- DAGs and Operators
- Workflow Scheduling
- Monitoring and Logging



By this point, we've got our environment setup with our DAGs running at predefined schedule or with triggered events. The last topic we'll cover before you practice what you've learned in your labs is how to monitor and troubleshoot your cloud functions and Airflow workflows.

Monitor current and historical workflow progress



The screenshot shows the Airflow interface for monitoring DAG runs. At the top, there's a navigation bar with links for Airflow, DAGs, Data Profiling, Browse, Admin, Docs, and About, along with a timestamp of 20:52 UTC and a user icon. Below the navigation is a search bar and filter options. A specific filter is applied for 'Dag Id' equals 'composer_sample_simple_greeting'. The main area displays a table of historical runs for this DAG, with columns for State, Dag Id, Execution Date, Run Id, and External Trigger. All five runs listed are in the 'SUCCESS' state.

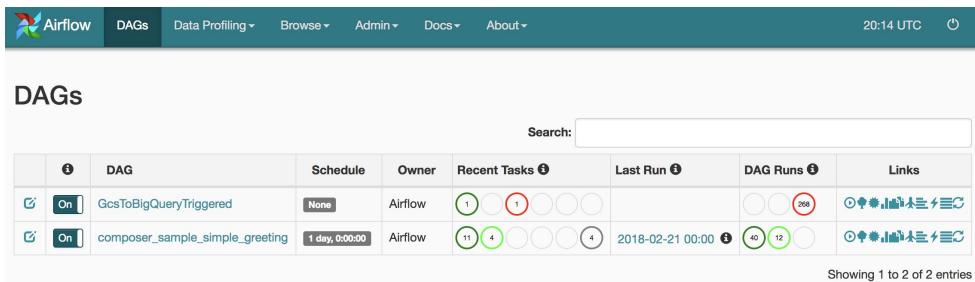
	State	Dag Id	Execution Date	Run Id	External Trigger
✓	SUCCESS	composer_sample_simple_greeting	09-30T00:00:00	scheduled_2018-09-30T00:00:00	↻
✓	SUCCESS	composer_sample_simple_greeting	09-29T00:00:00	scheduled_2018-09-29T00:00:00	↻
✓	SUCCESS	composer_sample_simple_greeting	09-28T00:00:00	scheduled_2018-09-28T00:00:00	↻
✓	SUCCESS	composer_sample_simple_greeting	09-27T00:00:00	scheduled_2018-09-27T00:00:00	↻
✓	SUCCESS	composer_sample_simple_greeting	09-26T00:00:00	scheduled_2018-09-26T00:00:00	↻

Google Cloud

One of the most common reasons you'll want to investigate the historical runs of your DAGs is in the event that your workflow simply stops working. Note that you can have it auto-retry for a set number of attempts in case it's a transient bug but sometimes you just can't get your workflow to run at all in the beginning.

In the Dag Runs you can monitor when your pipelines ran and in what State like success, running, or failure. The quickest way to get to this page is clicking on the schedule for any of your DAGs from the main DAGs page. Here we have 5 successful runs over 5 days for this DAG so this one seems to be running just fine.

Quickly gauge pipeline health of DAG Runs



Airflow DAGs Data Profiling ▾ Browse ▾ Admin ▾ Docs ▾ About ▾ 20:14 UTC ⚡

DAGs

Search:

	DAG	Schedule	Owner	Recent Tasks 1	Last Run 2018-02-21 00:00	DAG Runs 0	Links
GcsToBigQueryTriggered	On	None	Airflow	1 1		268	View Details
composer_sample_simple_greeting	On	1 day, 0:00:00	Airflow	11 4	2018-02-21 00:00	40 12	View Details

Showing 1 to 2 of 2 entries

Which pipeline is healthier?



Back on the main page for DAGs, we see some red which indicates trouble with some of our recent DAG runs.

Quickly gauge pipeline health of DAG Runs

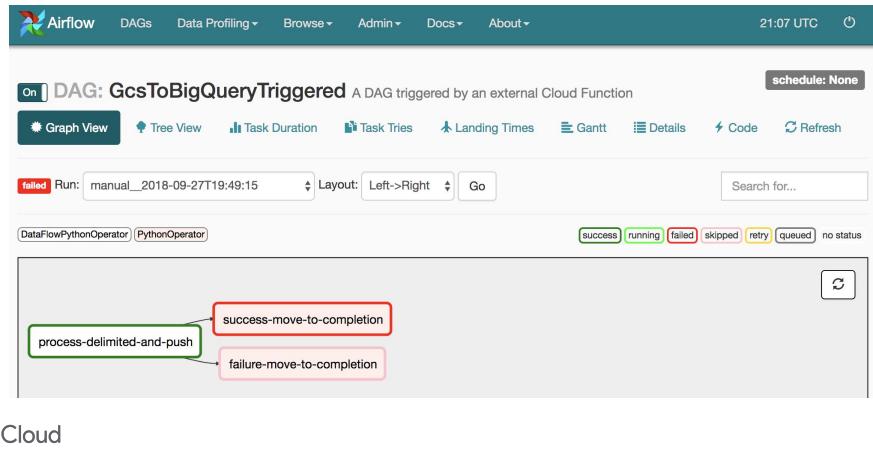
	DAG	Schedule	Owner	Recent Tasks	Last Run	DAG Runs	Links
<input checked="" type="checkbox"/>	On GcsToBigQueryTriggered	None	Airflow	1 (green), 1 (red), 0 (white)		268 (red box)	
<input checked="" type="checkbox"/>	On composer_sample_simple_greeting	1 day, 0:00:00	Airflow	11 (green), 4 (white), 0 (red)	2018-02-21 00:00	40 (green), 12 (white)	

Showing 1 to 2 of 2 entries



Speaking of DAG Runs, you'll note the three circles below which indicates how many runs passed, are currently active, or have failed. It certainly doesn't look good for 268 runs failed and 0 passed for this first DAG. Let's see what happened. We click on the name of the DAG to get to the visual representation.

Quickly gauge pipeline health of DAG Runs



It looks like the first task is succeeding, judging by the green border, but the next task ‘success-move-to-completion’ is failing. Note that the lighter pink color for the ‘failure-move-to-completion’ node means that node was skipped. So reading into this a bit, the CSV file was correctly processed by Dataflow in the first task but there is some issue moving the CSV file to a different Cloud Storage bucket as part of task two.

To troubleshoot, click on the node of a particular task and then click LOGS.

Monitor and troubleshoot Airflow step errors in logs

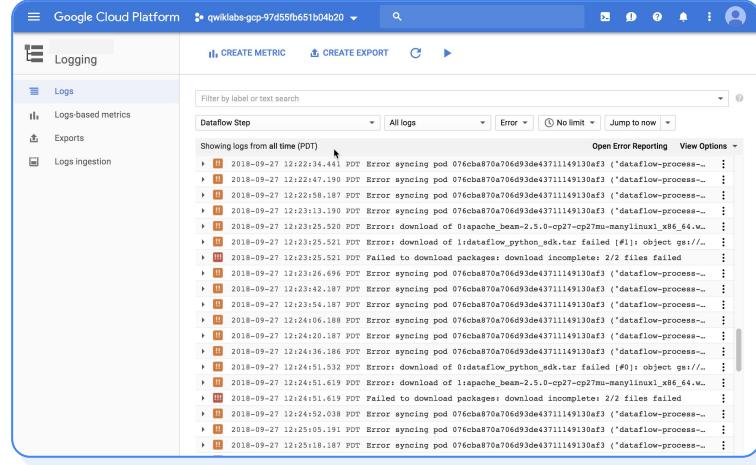
The screenshot shows the Airflow web interface with the following details:

- DAG:** GcsToBigQueryTriggered
- Schedule:** None
- Task Instance:** success-move-to-completion (2018-09-27 19:49:15)
- Log tab selected:** Task Instance Details, Rendered Template, Log (highlighted), XCom
- Log by attempts:** Attempt 1
- Log content:** Shows several INFO-level log entries related to reading a remote log from Google Cloud Storage and running on an airflow worker.
- Starting attempt 1 of:** [redacted]
- Final error message:** errors.HttpError: <HttpError 400 when requesting https://www.googleapis.com/storage/v1/b/input-bucket-999/o/usa_names.csv/copyTo/b/gs%3A%2Foutput-bucket-999%2F/o/success%2F2018-09-27%2Fusa_names.csv?alt=json returned "Invalid bucket name: 'gs://output-bucket-999/'>



Here you will find the logs for that specific Airflow run. I search for the word “error” and then start my diagnosis there. Here this was a pretty simple error where it was trying to copy a file from an input bucket to an output bucket and the output bucket didn’t exist or was named poorly.

Monitor Dataflow health in Cloud Logging



The screenshot shows the Google Cloud Platform Logs viewer interface. The left sidebar has sections for Logging, Logs-based metrics, Exports, and Logs ingestion. The main area is titled 'Dataflow Step' and shows a list of log entries. A filter bar at the top allows filtering by label or text search, log type (All logs, Error), and time range (No limit, Jump to now). The log entries are timestamped and show various error messages related to Dataflow steps.

Timestamp	Message
2018-09-27 12:12:22,34,441 PDT	Error syncing pod 076cba870a706d93de43711149130af3 ('dataflow-process-...
2018-09-27 12:12:22,47,190 PDT	Error syncing pod 076cba870a706d93de43711149130af3 ('dataflow-process-...
2018-09-27 12:12:22,58,187 PDT	Error syncing pod 076cba870a706d93de43711149130af3 ('dataflow-process-...
2018-09-27 12:12:23,13,190 PDT	Error syncing pod 076cba870a706d93de43711149130af3 ('dataflow-process-...
2018-09-27 12:12:31,25,520 PDT	Error: download of /apache_beam-2.5.0-0-p27mu-manylinux_x86_64.w...
2018-09-27 12:12:31,25,521 PDT	Error: download of /dataflow_python_sdk.tgz failed (#1): object g...
2018-09-27 12:12:31,25,521 PDT	Failed to download packages: download incomplete: 2/2 files failed
2018-09-27 12:12:31,26,436 PDT	Error syncing pod 076cba870a706d93de43711149130af3 ('dataflow-process-...
2018-09-27 12:12:31,42,187 PDT	Error syncing pod 076cba870a706d93de43711149130af3 ('dataflow-process-...
2018-09-27 12:12:31,44,187 PDT	Error syncing pod 076cba870a706d93de43711149130af3 ('dataflow-process-...
2018-09-27 12:12:40,06,188 PDT	Error syncing pod 076cba870a706d93de43711149130af3 ('dataflow-process-...
2018-09-27 12:12:42,20,187 PDT	Error syncing pod 076cba870a706d93de43711149130af3 ('dataflow-process-...
2018-09-27 12:12:43,16,186 PDT	Error syncing pod 076cba870a706d93de43711149130af3 ('dataflow-process-...
2018-09-27 12:12:43,16,186 PDT	Error: download of /dataflow_python_sdk.tgz failed (#1): object g...
2018-09-27 12:12:43,16,186 PDT	Failed to download packages: download incomplete: 2/2 files failed
2018-09-27 12:12:43,51,619 PDT	Error: download of /apache_beam-2.5.0-0-p27mu-manylinux_x86_64.w...
2018-09-27 12:12:43,51,619 PDT	Failed to download packages: download incomplete: 2/2 files failed
2018-09-27 12:12:45,52,038 PDT	Error syncing pod 076cba870a706d93de43711149130af3 ('dataflow-process-...
2018-09-27 12:12:45,52,038 PDT	Error syncing pod 076cba870a706d93de43711149130af3 ('dataflow-process-...
2018-09-27 12:12:45,52,038 PDT	Error syncing pod 076cba870a706d93de43711149130af3 ('dataflow-process-...
2018-09-27 12:12:45,51,187 PDT	Error syncing pod 076cba870a706d93de43711149130af3 ('dataflow-process-...

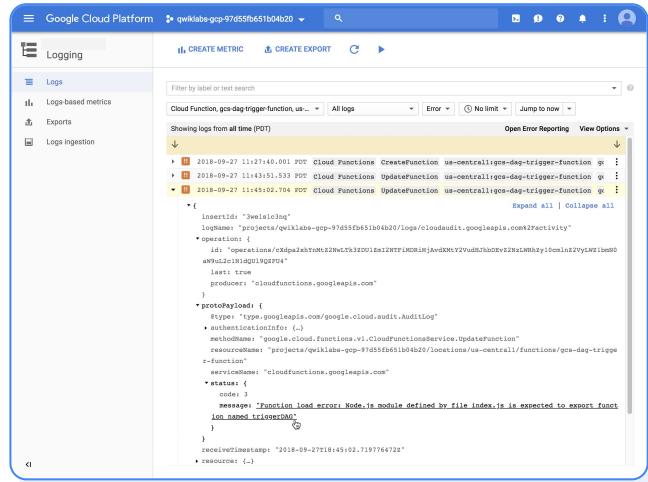


Another tool in your toolkit for diagnosing Airflow failures is the general Google Cloud logs. Since Airflow launches other Google Cloud services through tasks, you can see and filter for errors for those services in Cloud Logging as you would debugging any other normal application. Here I've filtered for Dataflow step errors to troubleshoot why my workflow is failing.

It turns out that I had not changed the name of the output bucket for the CSV file so after the file was processed by Dataflow as part of step 1, it dumped the completed file back into the input bucket which triggered another Dataflow job for processing and so on.

https://console.cloud.google.com/logs/viewer?project=qwiklabs-gcp-97d55fb651b04b20&minLogLevel=400&expandAll=false×tamp=2018-10-01T21:03:39.539000000Z&customFacets=&limitCustomFacetWidth=true&interval=NO_LIMIT&resource=cloud_function%2Ffunction_name%2Fqcs-dag-trigger-function%2Fregion%2Fus-central1&scrollTimestamp=2018-09-27T18:45:02.704000000Z

Monitor Cloud Function health in Cloud Logging



The screenshot shows the Google Cloud Platform Logging interface. The left sidebar has options for Logs, Log-based metrics, Experts, and Log ingestion. The main area shows a log entry for a Cloud Function named 'gcs-dag-trigger-function'. The log entry details the creation of the function, its update, and an error message indicating a function load error due to a missing 'triggerDAG' function. The log entry also includes metadata such as timestamp, log name, and resource information.

```
2018-09-27 11:27:40.001 PDT Cloud Functions CreateFunction us-central1:gcs-dag-trigger-function ge ...
2018-09-27 11:40:51.539 PDT Cloud Functions UpdateFunction us-central1:gcs-dag-trigger-function ge ...
2018-09-27 11:45:02.707 PDT Cloud Functions UpdateFunction us-central1:gcs-dag-trigger-function ge ...
{
  insertId: "2ewel1sl0g"
  logName: "projects/qwiklabs-gcp-97d55fb651b04b20/logs/cloudaudit.googleapis.com%2factivity"
  operation: {
    id: "operations/cx0pa2ahytrN1ZNGC5J2UO1mz1ZBT1LNU1MIAwvXKmT2vUdnJhAbvV2x1Mhby1oelnz2VyyLW1l0m
    aWuLz2c13iq1q1gPgr7v4"
    last: true
    producer: "cloudfunctions.googleapis.com"
  }
  protoPayload: {
    type: "type.googleapis.com/google.cloud.audit.AuditLog"
    authenticationInfo: {}
    methodName: "google.cloud.functions.v1.CloudFunctionsService.UpdateFunction"
    resourceNames: "projects/qwiklabs-gcp-97d55fb651b04b20/locations/us-central1/functions/gcs-dag-trigge
    rFunction"
    serviceAccount: "cloudfunctions.googleapis.com"
  }
  attributes: {
    code: 3
    message: "Function load error: Node.js module defined by file index.js is expected to export func
    tion named 'triggerDAG'"
  }
}
receiveTimestamp: "2018-09-27T18:45:02.718776472Z"
resource: {}
```



You might be wondering, if there's an error with my Cloud Function my Airflow instance would never have been triggered or issued any logs at all (since it was unaware we were trying to trigger it) and you're exactly right. If you're using Cloud Functions, be sure to check the normal Google Cloud logs for errors and warnings in addition to your Airflow logs.

In this example, each time I uploaded a CSV file to my Cloud Storage bucket hoping to trigger a Cloud Function and then my DAG I got an error expected to find function triggerDAG. Remember way back when I said Cloud Functions were case-sensitive? Looking for a function with capital DAG doesn't exist if it's capital D lowercase a and g so be sure to be mindful when setting up your Cloud Function for the first time.

https://console.cloud.google.com/logs/viewer?project=qwiklabs-gcp-97d55fb651b04b20&minLogLevel=400&expandAll=false×tamp=2018-10-01T21:03:39.539000000Z&customFacets=&limitCustomFacetWidth=true&interval=NO_LIMIT&resource=cloud_function%2Ffunction_name%2Fgcs-dag-trigger-function%2Fregion%2Fus-central1&scrollTimestamp=2018-09-27T18:45:02.704000000Z



An Introduction to Cloud Composer

Objectives

- Use GCP Console to create a Cloud Composer environment
- View and run a DAG in the Airflow web interface
- View the results of the wordcount job in storage

Module Summary

Cloud Data Fusion for building data pipelines.

Cloud Composer, Cloud Functions and Cloud Scheduler are the glue for your data pipelines.

