# Streaming data from Twitter to GCP ●

11/3/2019    Practicum                                                            Mia Tadić

**Google Cloud Platform (GCP)** is a suite of public cloud computing services offered by Google. The platform includes a range of hosted services for compute, storage and application development that run on Google hardware. Services used in this blog are Google Cloud Pub/Sub and Google BigQuery.

- *Google Cloud Pub/Sub* is a managed, scalable and real-time messaging ingestion service that allows messages to be exchanged between applications.
- *Google BigQuery* is a scalable, managed enterprise data warehouse for analytics. It is a big data service for data processing and analytics for SQL-like queries made against multi-terabyte data sets.
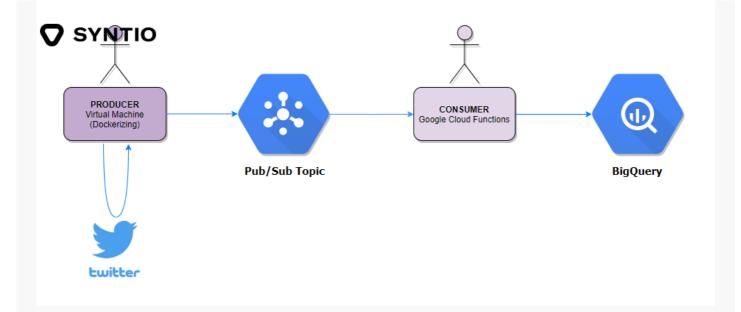
## Task idea

This blog demonstrates the task of ingesting data from remote API (i.e. Twitter's API) to cloud (i.e. Google Cloud).

It was done using Producer and Consumer.

- *Producer* sent tweets containing hashtag „dataengineering" from Twitter to Pub/Sub topic. Producer was dockerized and then ran on Google Cloud's Virtual Machine.

- *Consumer* streamed messages from Pub/Sub topic to BigQuery table. Consumer was created using Google Cloud Functions, and it was **automatically** triggered by any new messages.

Described process is shown in following diagram:

# Procedure

Here are the prerequisites and steps of our procedure.

Prerequisits for the task are:

1. Twitter Developer account
2. Google Cloud Platform account
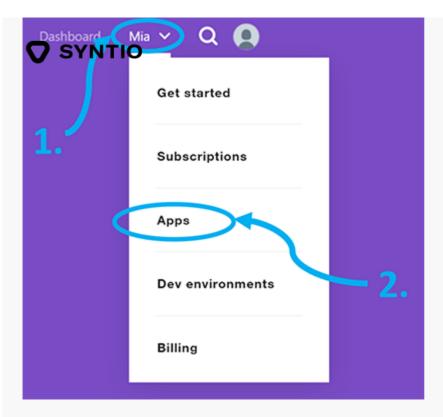3. Installed Docker

Main steps of procedure are:

1. Obtaining Twitter credentials
2. Creating a Producer
3. Creating a Consumer
4. Dockerizing a Producer

Let's dig into steps of the procedure.

## Obtain Twitter credentials

App's **client key** and **secret key**, **access token** and **token secret** can be found in Twitter Dashboard.

Those were needed in later step.

## Create a Producer

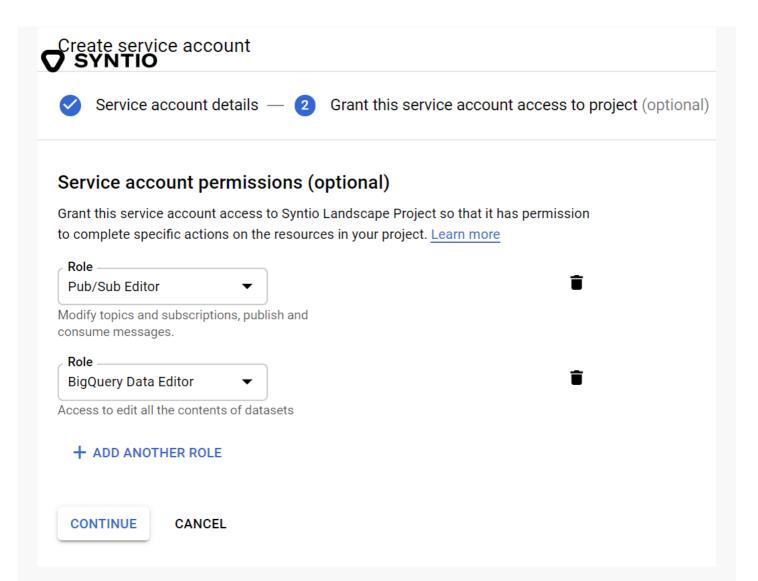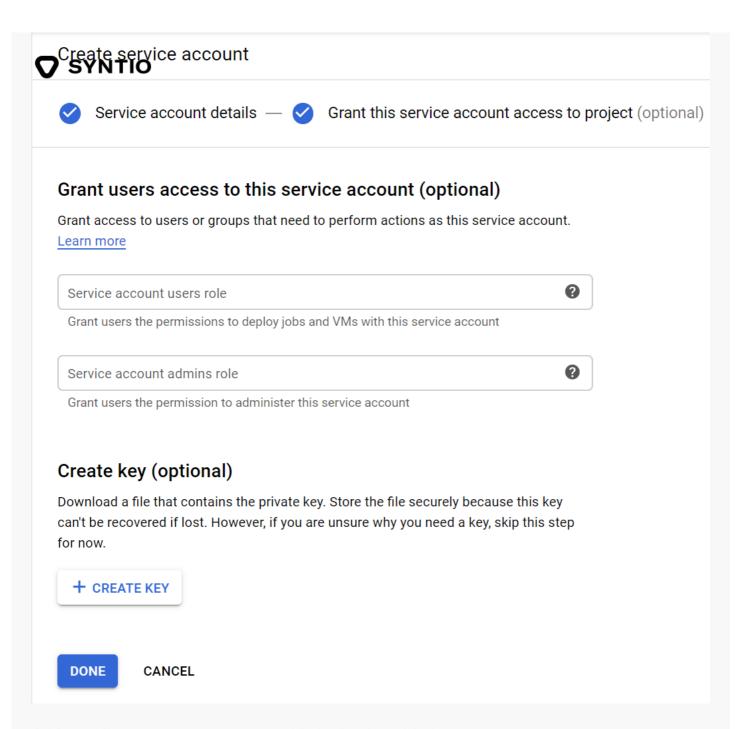First, following resources were created:

### 1. Service account

Service account authorized access to needed GCP resources (Pub/Sub, BigQuery).

In GCP console, *Service accounts* can be found under *IAM & Admin* section. By clicking on **CREATE SERVICE ACCOUNT** button, several steps were required:

- Step 1: Filling in desired details
- Step 2: Adding „Pub/Sub Editor" and „BigQuery Data Editor" roles

# Create service account

✓ Service account details — ② Grant this service account access to project (optional)

## Service account permissions (optional)

Grant this service account access to Syntio Landscape Project so that it has permission to complete specific actions on the resources in your project. Learn more

**Role**

Pub/Sub Editor ▼                                    🗑

Modify topics and subscriptions, publish and consume messages.

**Role**

BigQuery Data Editor ▼                               🗑

Access to edit all the contents of datasets

**+ ADD ANOTHER ROLE**

**CONTINUE**    CANCEL

---

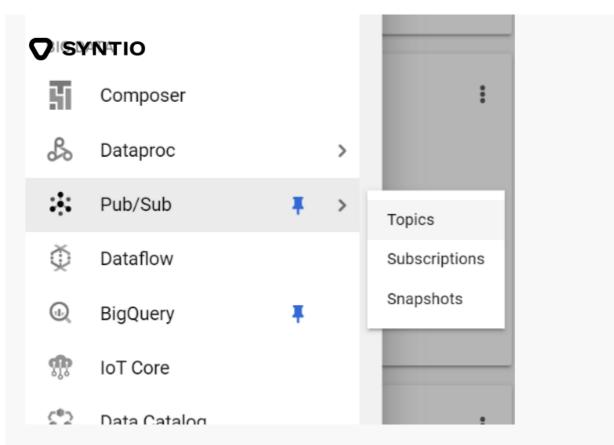- Step 3: Downloading JSON key file by choosing **CREATE KEY** option

Additionally, user environment variable had to be added:

name=GOOGLE_AUTH_CREDENTIALS

value=path_to_downloaded_json_file

## 2. Pub/Sub Topic

Pub/Sub topic was the place where Producer sent tweets to.

## 3. Pub/Sub Subscription

Subscription pulled data from desired Pub/Sub topic. This was a way to see how data looked like. _Subscriptions_ can be found in the picture above.

Finally, it was all set for Producer to be created.

As mentioned, the Producer's task was to send tweets containing hashtag "dataengineering" from Twitter to Pub/Sub topic.

**Note:** Keys and tokens obtained from Twitter were needed here.

Firstly, Pub/Sub package needed to be installed:

```
pip install –upgrade google-cloud-pubsub
```

To begin with Python script, following packages were imported:

```python
import json

from tweepy import Stream
from tweepy import OAuthHandler
from tweepy.streaming import StreamListener
from google.cloud import pubsub_v1
```

Since Producer was said to send tweets to Pub/Sub topic, connection to Pub/Sub topic had to be established somehow. It was done this way:

```python
client = pubsub_v1.PublisherClient()
```

```
topic_path = client.topic_path(<project_id>, <topic_name>)
```

On the other side, since Producer accessed Twitter to grab its data, it was necessary to make connection to Twitter too. That connection was made this way:

```
auth = OAuthHandler(<twitter-client-key>, <twitter-secret-key>)
auth.set_access_token(<twitter-access-token>, <twitter-token-secret>)
twitterStream = Stream(auth, listener())
twitterStream.filter(track=["dataengineering"])
```

This code called *tweepy*'s implemented methods. Additionally, object *listener* was introduced in the code for the purposes of the task. It was a class which inherited *tweepy*'s *StreamListener* class.

For the demonstration purposes, only following parts of the grabbed tweets were taken to save them to Pub/Sub topic:

- created_at
- id
- text
- source
- retweet_count
- user_name
- user_location
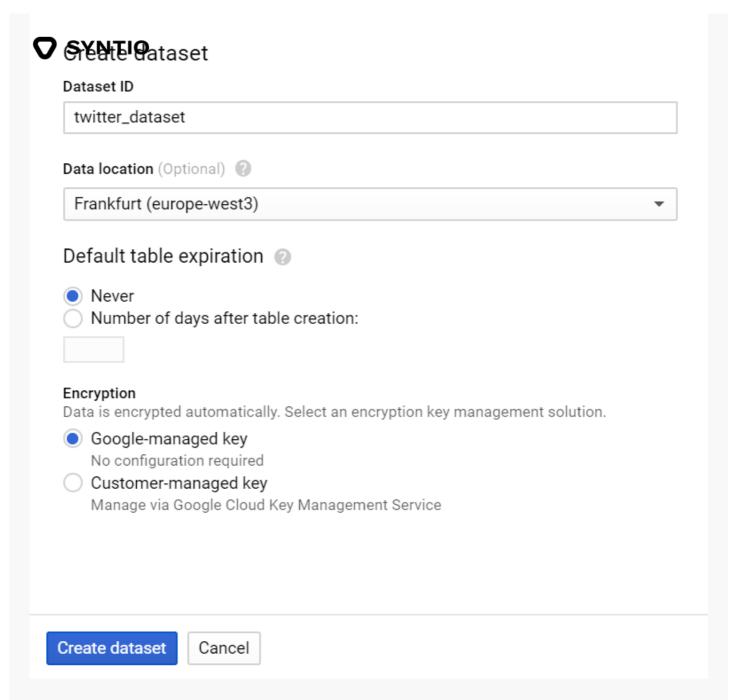- user_followers_count
- user_statuses_count

Extracted parts were concatenated in *tweet_message* after which *tweet_message* was published to Pub/Sub topic using following command:

```
self.client.publish(self.topic_path, data=tweet_message.encode('utf-8'))
```

## CREATE A CONSUMER

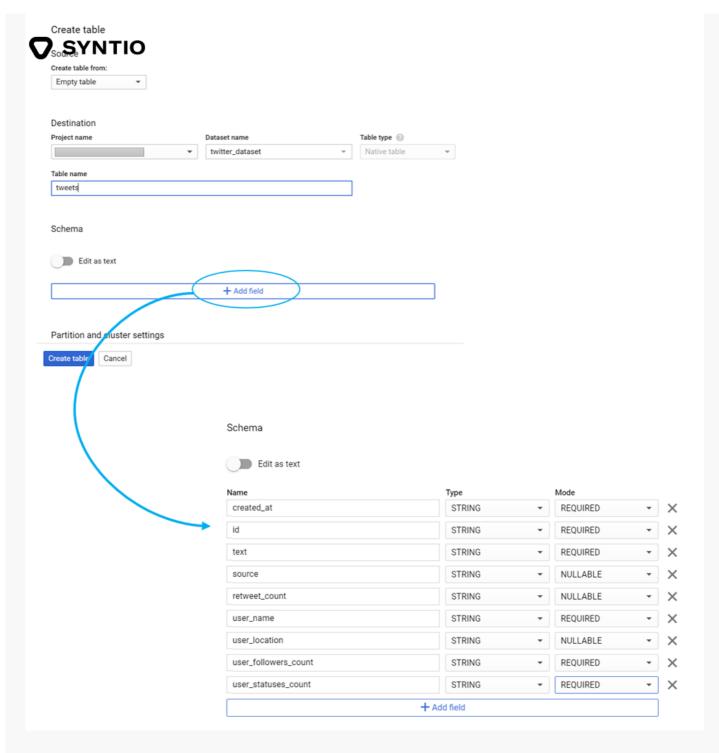As stated before, Consumer's task was to stream messages from Pub/Sub topic to BigQuery table – so a BigQuery table had to be created.

BigQuery tables are generally contained in datasets. Therefore, dataset had to be created first. It was created in GCP console, by finding *BigQuery* under *Big Data* section, choosing project and clicking on **CREATE DATASET**. It popped a form which can be seen in a picture below:

## Create dataset

**Dataset ID**

twitter_dataset

**Data location** (Optional) ?

Frankfurt (europe-west3) ▼

## Default table expiration ?

◉ Never

◯ Number of days after table creation:

[ ]

**Encryption**

Data is encrypted automatically. Select an encryption key management solution.

◉ Google-managed key

No configuration required

◯ Customer-managed key

Manage via Google Cloud Key Management Service

[ Create dataset ]  [ Cancel ]

Then, by entering the created dataset and simply clicking on **CREATE TABLE** button, table was created. Based on extracted tweet's parts, table's schema must have looked like this:
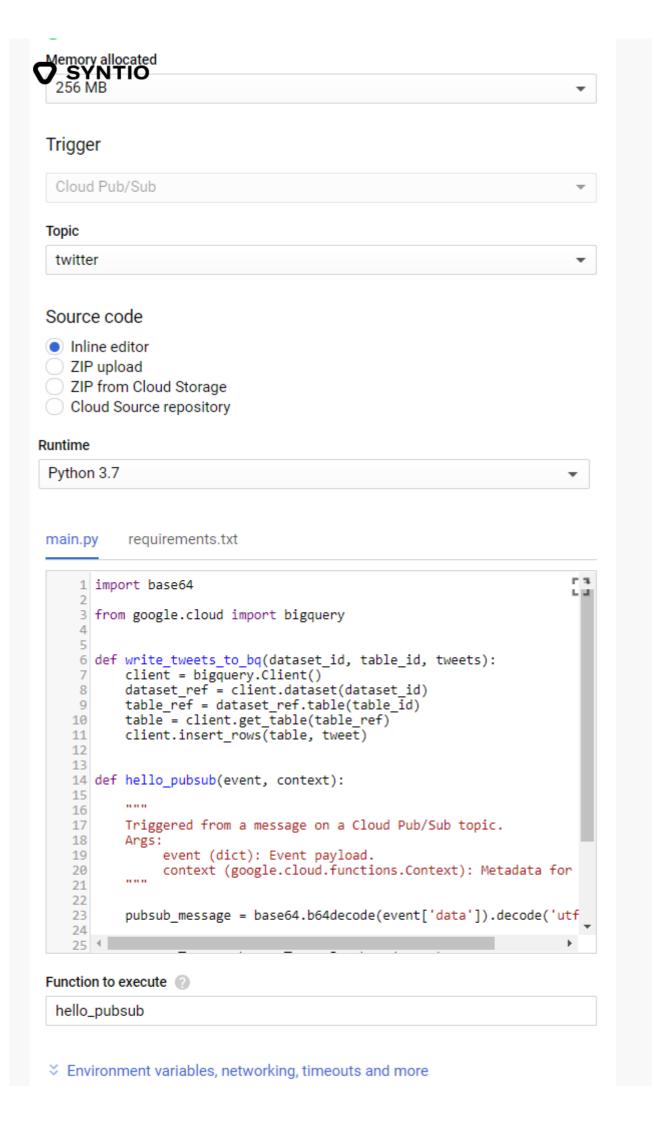
At last, Consumer was ready to be created. As mentioned, Consumer was created using Google Cloud Functions.



Missing values needed to be filled in:

**Memory allocated**

```
256 MB                                          ▼
```

## Trigger

```
Cloud Pub/Sub                                   ▼
```

**Topic**

```
twitter                                         ▼
```

## Source code

- ● Inline editor
- ○ ZIP upload
- ○ ZIP from Cloud Storage
- ○ Cloud Source repository

**Runtime**

```
Python 3.7                                       ▼
```

**main.py**   requirements.txt

```python
1  import base64
2
3  from google.cloud import bigquery
4
5
6  def write_tweets_to_bq(dataset_id, table_id, tweets):
7      client = bigquery.Client()
8      dataset_ref = client.dataset(dataset_id)
9      table_ref = dataset_ref.table(table_id)
10     table = client.get_table(table_ref)
11     client.insert_rows(table, tweet)
12
13
14 def hello_pubsub(event, context):
15
16     """
17     Triggered from a message on a Cloud Pub/Sub topic.
18     Args:
19          event (dict): Event payload.
20          context (google.cloud.functions.Context): Metadata for
21     """
22
23     pubsub_message = base64.b64decode(event['data']).decode('utf
24
25
```

**Function to execute** ❓

```
hello_pubsub
```

⌄ Environment variables, networking, timeouts and more

File *main.py* contained Consumer code, while *requirements.txt* contained function dependencies (packages with their versions).

Google Cloud Function was triggered from a message on a Cloud Pub/Sub topic. This event was passed to *hello_pubsub* function as an argument and then function collected data (tweets) from the event:

```
pubsub_message = base64.b64decode(event['data']).decode('utf-8')
```

Since Consumer's task was to write messages to BigQuery table, Consumer had to connect to BigQuery table. To establish the connection, the following code was used:

```
client = bigquery.Client()
dataset_ref = client.dataset(<dataset_id>)
table_ref = dataset_ref.table(<table_id>)
table = client.get_table(table_ref)
client.insert_rows(table, <tweet>)
```

# Dockerize

Following step in our *steps of procedure* was to dockerize Producer. The idea was to pack the Producer into Docker image, send the image to virtual machine and then run it from there.

For the purposes of this task, Docker Desktop and Docker Hub were used and the instructions and installation file for various operational systems can be found here.

For Docker, two files (*Dockerfile* and *requirements.txt*) must be created and saved into the **same Python project** as producer code.

Dockerfile should look like this:

```
# Use an official Python runtime as a parent image
FROM python:3.7-slim

# Set the working directory to /app
WORKDIR /app

# Copy the current directory contents into the container at /app
ADD . /app

# Install any needed packages specified in requirements.txt
RUN pip install --trusted-host pypi.python.org -r requirements.txt

# Define environment variable
```
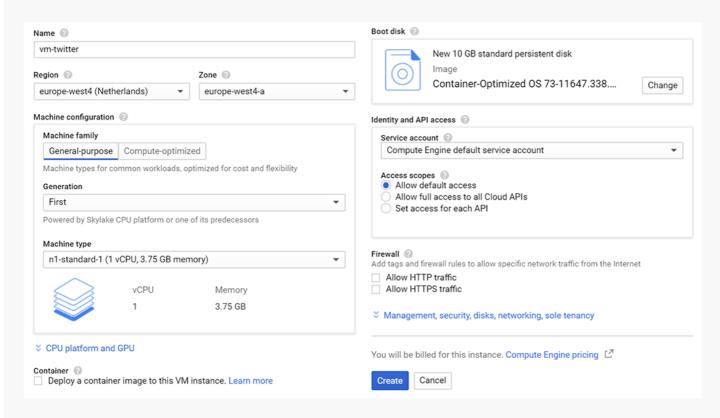
```
ENV GOOGLE_APPLICATION_CREDENTIALS="<json-key-file>"
```



```
# Run app.py when the container launches
CMD ["python", "producer.py"]
```

## Note:

- After checking which Python version is being used, it is fine to just add suffix "slim" with no additional installation.
- JSON key file must be copied into the **same Python project** directory (it is Dockerfile's working directory now, so that is the place where Docker searches for environment variable)

As described in the beginning of the blog, Docker image was deployed to the virtual machine. Therefore, in this preparation for dockerizing, next step was to create GCP's virtual machine.

To create VM, *VM instances* can be found in GCP console, under *Compute* section and by choosing *Compute Engine* and then clicking **CREATE INSTANCE**. For this demonstration, following details were sufficient:



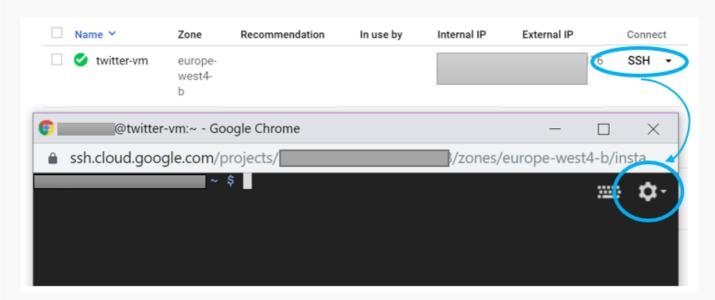At last, steps for deploying Docker image to VM were:

1. Build Docker image

```
docker build -t producer .
```

2. Zip Docker image

```
docker save -o producer.tar producer
```

**Note:** Make sure to run both steps from producer's Python project directory

SYNTIO

3. Upload zip file to virtual machine

By clicking first on **SSH** button in VM instances in GCP console, and afterwards on **Settings** button, *Upload file* option was chosen and *producer.tar* file was uploaded.
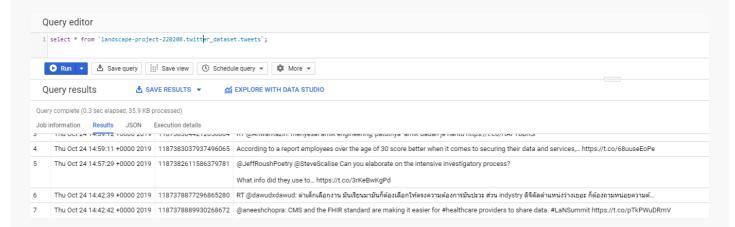


4. Unzip Docker image

```
docker load -i producer.tar
```

5. Run Docker image

```
docker run -d producer
```

With this, the process of ingesting messages to Pub/Sub topic was completed.

To check if created Producer and Consumer work properly, BigQuery table can be checked. Table from this demonstration looked like this:



# Conclusion and ideas

In this step-by-step tutorial, it was presented how to stream data to cloud from remote API using a simple Producer-Consumer flow. Instead of GCP and Twitter, there are plenty more cloud platforms and API's to use. There are also other Google Cloud resources to play with, like Google Cloud Storage or Google Cloud SQL. Another tip is to send data from Cloud Pub/Sub topic to BigQuery using Google Cloud Dataflow.

**Note**: Don't forget to delete the resources when they are not needed anymore.

# SYNTIO
# DATA ENGINEERING CO.

© Syntio, 2020.

LinkedIn   Privacy

European Union
Together with EU funds

EUROPEAN STRUCTURAL
AND INVESTMENT FUNDS

Operational programme
COMPETITIVENESS
AND COHESION