# Use Python and Google Cloud to Schedule API data Import into BigQuery

In this tutorial, I'm going to show you how to set up a serverless data pipeline in GCP that will do the following.

1. Schedule the download of a csv file from the internet
2. Import the data into BigQuery

Note - This tutorial generalizes to any similar workflow where you need to import data into BigQuery.

Here's the workflow. There are a few moving parts, but it's not too complicated.



**The Workflow**

We will need to create a few things to get started. You'll need a Google Cloud account and a project, as well as API access to Cloud Storage, Cloud Functions, Cloud Scheduler, Pub/Sub, and BigQuery.
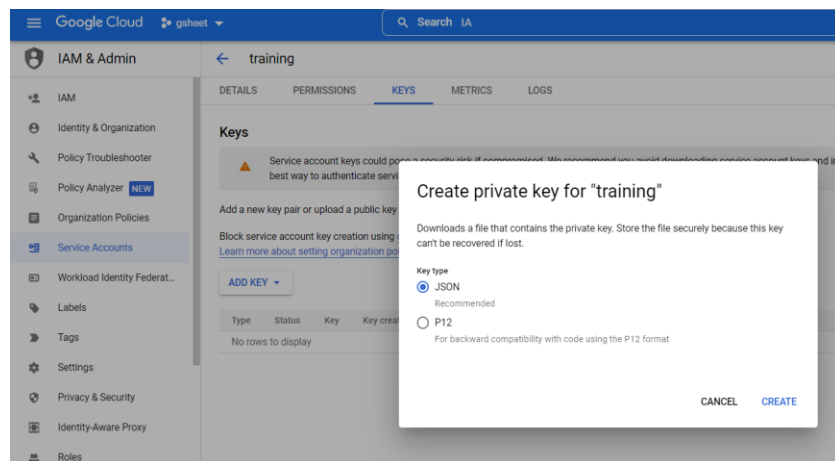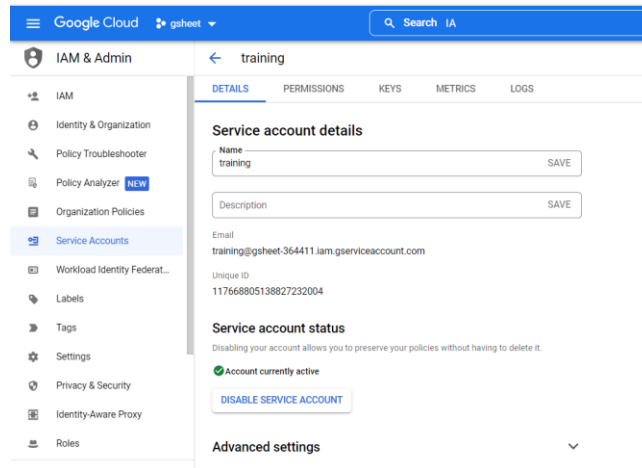
Here's a summary of what we're going to build.

1. One BigQuery dataset and table
2. One Cloud Functions
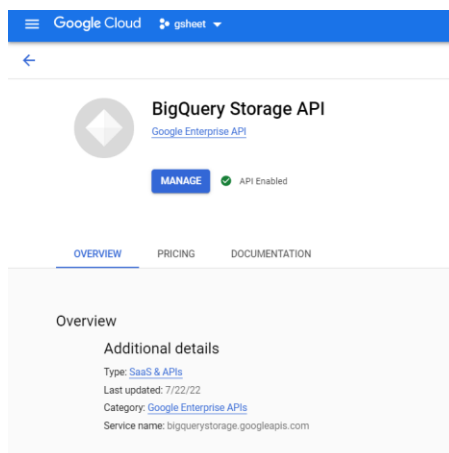3. One Cloud Scheduler job

All the code examples are in this [GitHub repo](#).

Even though this example is for the Thailand air quality dataset, it can be used for any situation where you import data into BigQuery.

1. Create new GCP project or used an existing project.

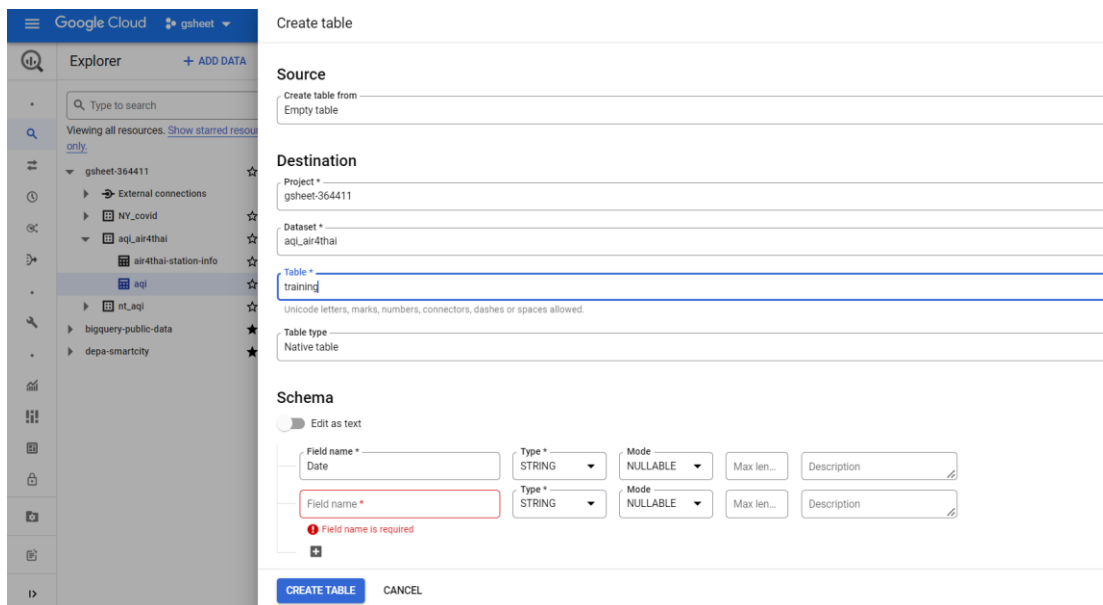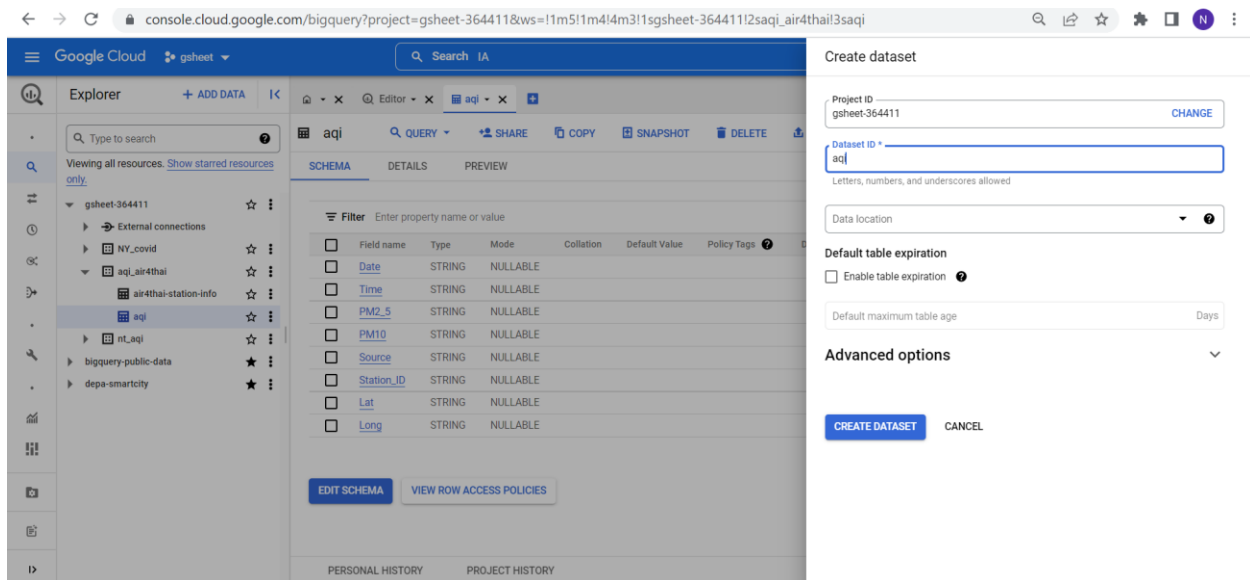2. create a service account and create json key file that uses for remote access from the local host

## 3. Enable BigQuery API



## 4. Create A BigQuery Dataset and Table

Creating a BigQuery dataset and table is very simple. Just remember that you first create a dataset, then create a table.

When you create your BigQuery table, you'll need to create a schema with the following fields. These BigQuery fields match the fields in the Thailand air quality json API's header.

Create the BigQuery table, which should have a schema that looks like this.

5. Create google cloud function, this example used trigger type as cloud pub/sub then fill all variables and click save and next.



6. select runtime as Python 3.8, write your python code and necessary python package in requirement.txt

Note *entry point should be the same as your main function

**Runtime**
Python 3.8

**Entry point ***
main

⚠ The specified entry point might not be present in your source code. Please ensure the entry point in your code matches the input field.

**Source code**
Inline Editor

📄 main.py

📄 requirements.txt

```python
import base64
import requests
import pandas as pd
from pandas.io import gbq

def api_to_df(url):
    data = requests.get(url)
    json = data.json()
    df = pd.json_normalize(json['stations'])
    cols = df.columns
    cols=cols.str.replace('LastUpdate.', '',regex = True)
    cols=cols.str.replace('.value', '',regex = True)
    cols=cols.str.replace('.aqi', '',regex = True)
    cols=cols.str.replace('AQI.Level', 'AQI_Level',regex = True)
    df.columns=cols
    df.rename(columns = {'stationID':'Station_ID','date':'Date', 'time':'Time','lat':'Lat','long':'Long','PM25':'PM2_5'}, inplace = True)
    df['Source'] = 'Air4Thai'
    data_out = df[['Date','Time','PM2_5','PM10','Source','Station_ID','Lat','Long']]
    return data_out

def main(data, context):
    df = api_to_df('http://air4thai.pcd.go.th/services/getNewAQI_JSON.php')

    df.to_gbq(destination_table='gsheet-364411.aqi_air4thai.aqi',project_id='gsheet-364411',
              if_exists='append')
```

**Runtime**
Python 3.8

**Entry point ***
main

⚠ The specified entry point might not be present in your source code. Please ensure the entry point in your code matches the input field.

**Source code**
Inline Editor

📄 main.py

📄 requirements.txt

```
# Function dependencies, for example:
# package>=version
google-cloud-bigquery
pandas_gbq
```

7. deploy your google cloud function, it will taking 2-3 minute.

8. if no error or malfunction, your function will appear green color

≡ **Google Cloud**  gsheet ▾ 🔍 Search clo ✕

(···) Cloud Functions | Functions | ➕ CREATE FUNCTION | ⟳ REFRESH

☰ Filter  Filter functions

| | | Environment | Name ↑ | Last deployed | Region | Trigger | Runtime | Memory allocated | Executed function | Authentication ❓ | Actions |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ☐ | ✅ | 1st gen | load_air4thai | Nov 7, 2022, 4:23:42 PM | us-central1 | Topic: aqi | Python 3.8 | 256 MB | main | | ⋮ |

9. check your data in your BigQuery table, your function get the data from API and writes it to BigQuery directly.

10. Create A Cloud Scheduler job - Cloud Scheduler is a GCP utility that allows you to schedule tasks. If you've used cron, then you'll have a good idea of how this works. If not, no worries. Cron is very simple to figure out. You basically just need to know that cron goes something like this. You've got 5 options that you can populate however you like. In our example, we'll set a Cloud Scheduler cron frequency for every ten minutes (this is obviously overkilled for a data set that updates once a day, but since we're in example-land, the point stands).

```
# ┌─────────────── min (0 - 59)
# │ ┌───────────── hour (0 - 23)
# │ │ ┌─────────── day of month (1 - 31)
# │ │ │ ┌───────── month (1 - 12)
# │ │ │ │ ┌─────── day of week (0 - 6) (0 to 6 are Sunday to
# │ │ │ │ │          Saturday, or use names; 7 is also Sunday)
# │ │ │ │ │
# │ │ │ │ │
# * * * * *  command to execute
```

Cloud Scheduler needs to trigger something to execute. For this tutorial, we'll target the Pub/Sub topic that we created in the Cloud Function we deployed above. Once you get the frequency populated, hit create (or update). Here's what your Cloud Scheduler should look like.

Then in the Cloud Scheduler console, hit "Run" to test the workflow. If everything works, you should see your BigQuery table populated with data.

Let's See the Data, Now that we've got the workflow pieced together and working, let's have a look at the air quality data in BigQuery. Here's what you should see in BigQuery that hourly update.