How to create the YouTube API?

 ➢ Go to https://console.cloud.google.com/
 ➢ Sign into your Google Account
 ➢ Create Project
 ➢ Go to Library -> select the required YouTube API from the list (we are using YouTube Data API v3)
 ➢ Enable API
 ➢ Create API key -> select Credentials > Create Credentials  > API Key

YouTube Documentation:
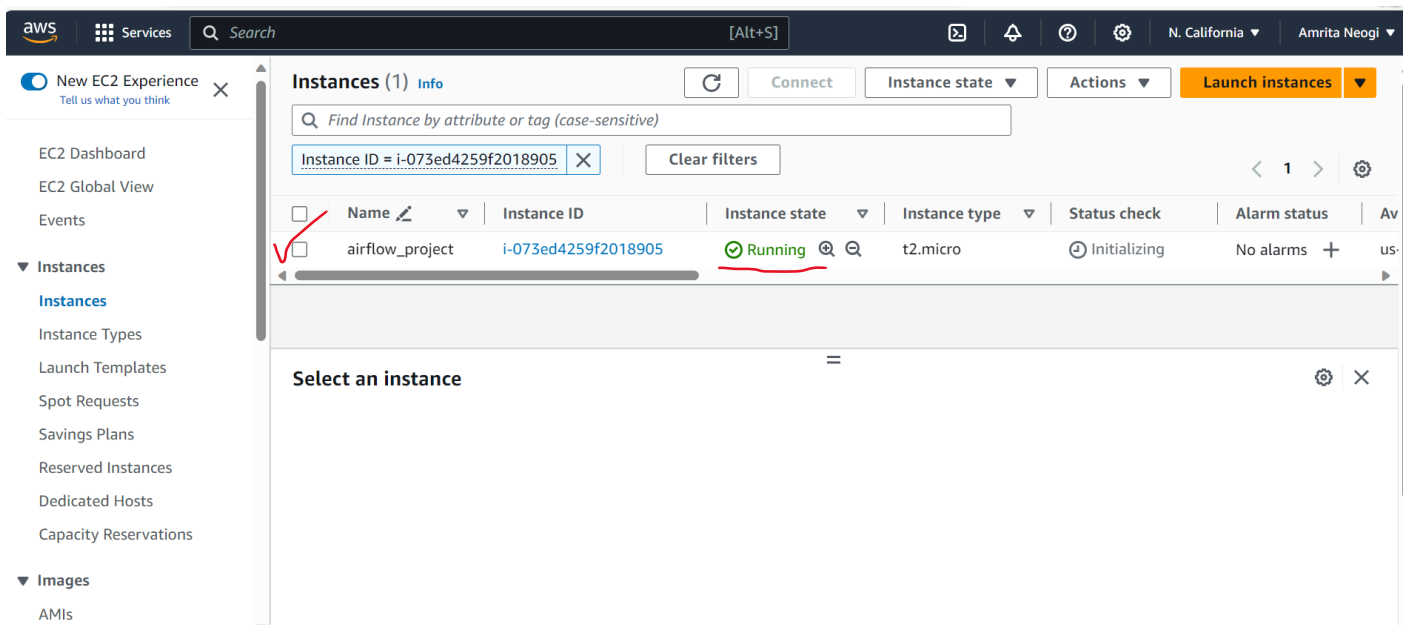
https://developers.google.com/youtube/v3

Prerequisites:

Python 2.7 or Python 3.5+

The pip package management tool

The Google APIs Client Library for Python:

pip install --upgrade google-api-python-client

 ➢ Execute the Python script youtube_etl.py to extract the YouTube data.

 ➢ Create EC2 instant that is the online machine and deploy Airflow on that.

   - Login to AWS account
   - Select EC2 > Instances > Launch Instances > provide Name and tags (e.g. airflow_project)
   - Choose Application and OS Image → Ubuntu
   - Choose Instance Type → t3_medium (additional charge required, else go with t2_micro it is free)
   - Create key-pair to access the EC2 instance:
       airflow_ec2_key →key will be downloaded → keep the key in the same folder
   - Allow SSH, HTTP, HTTPS traffic from internet
   - Launch Instance



 ➢ Connect to Airflow:

- Select the Instance > Click on 'Connect' > SSH client > copy the autogenerated example



- Open cmd and past the key, make sure to be in the same folder



Ubuntu console will open: (run command twice if the ubuntu console is not up the first time)

- Install the following SSH commands in the ubuntu console:
  sudo apt-get update
  sudo apt install python3-pip
  sudo pip install apache-airflow
  sudo pip install pandas
  sudo pip install s3fs
  sudo pip install tweepy→ for twitter data
  sudo pip install --upgrade google-api-python-client → for YouTube
- Check if everything is installed properly write 'airflow' in the ubuntu console:



- Connect to the airflow server:
  airflow standalone

➢ Create a 'youTube_dag.py' file (refer file for the code)→ import the ETL function from 'youTube_etl.py' file
➢ Create a S3 bucket (S3 buckets | S3 | Global (amazon.com))
➢ Save the file we created in the S3 bucket in the 'youTube_etl.py' file
  e.g., channel_stat.to_csv('s3://amrita-neogi-yt-bucket/youtube_stat.csv')
➢ Connect to the airflow again →make changes to 'airflow.cfg'

```
[core]
# The folder where your airflow pipelines live, most likely a
# subfolder in a code repository. This path must be absolute.
dags_folder = /home/ubuntu/airflow/dags   ←
# Hostname by providing a path to a callable, which will resolve the hostname.
# The format is "package.function".
```

➢ Create a new folder inside S3 copy the etl.py code from local folder to the EC2 machine
- mkdir youTube_dag
- cd youTube_dag
- sudo nano youTube_dag.py
- copy code from 'youTube_dag.py' and paste here
- ctrl + X to save the file
- repeat same process for etl file
  sudo nano youTube_etl.py
  copy code from 'youTube_etl.py' and paste here
  ctrl + X to save the file

** to stop airflow server → ctrl + C

➢ Make sure to have permission from the EC2 to write on the S3 bucket.



- Create IAM role



Give S3 and EC2 full access:

- Create role – 's3_ec2_airflow_role'
- Update IAM role