Google Cloud

# NOAA Public Datasets on Google Cloud

01/28/2024

Google Cloud

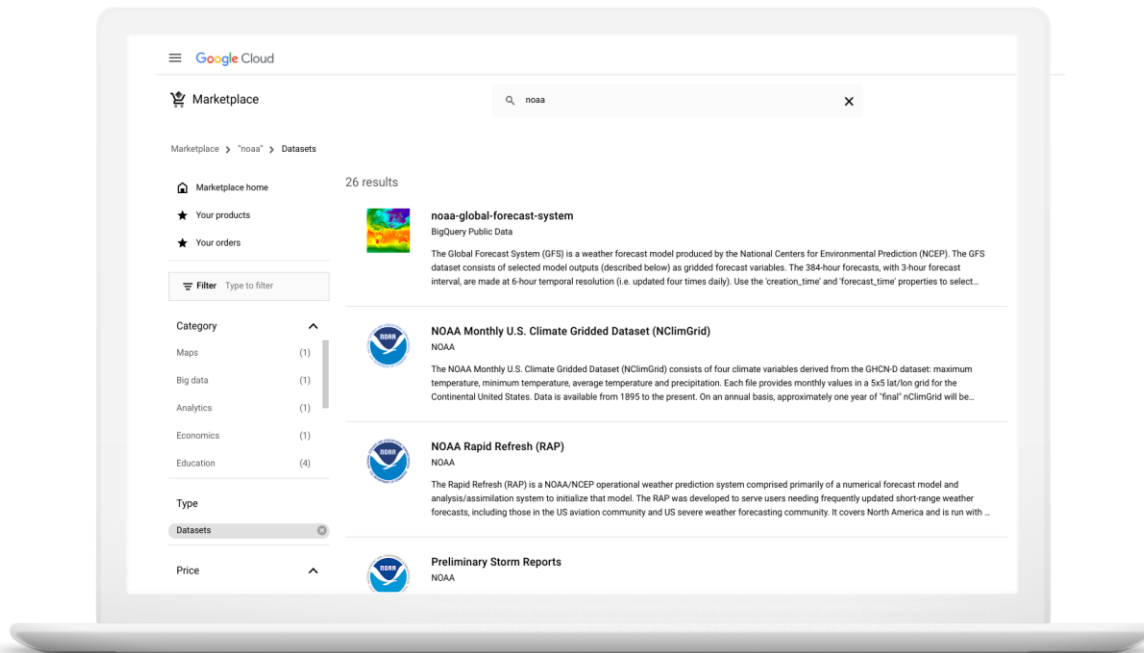# 01

## Public Datasets on Google Cloud

Google Cloud Public Datasets is a carefully curated and (mostly) Google managed dataset catalog from various sources all around the world, including weather data, shopping data, crypto, and even Google's own Search data.

# Datasets in Marketplace

- Google Cloud Marketplace is the source of truth for datasets in GCP
- Can search and filter through what's available
- No login required to browse dataset entries
- All consumption of raw (non-tabular, bucket) data is free
- BigQuery tabular data is charged per query



Google Cloud

An explosion of satellite data
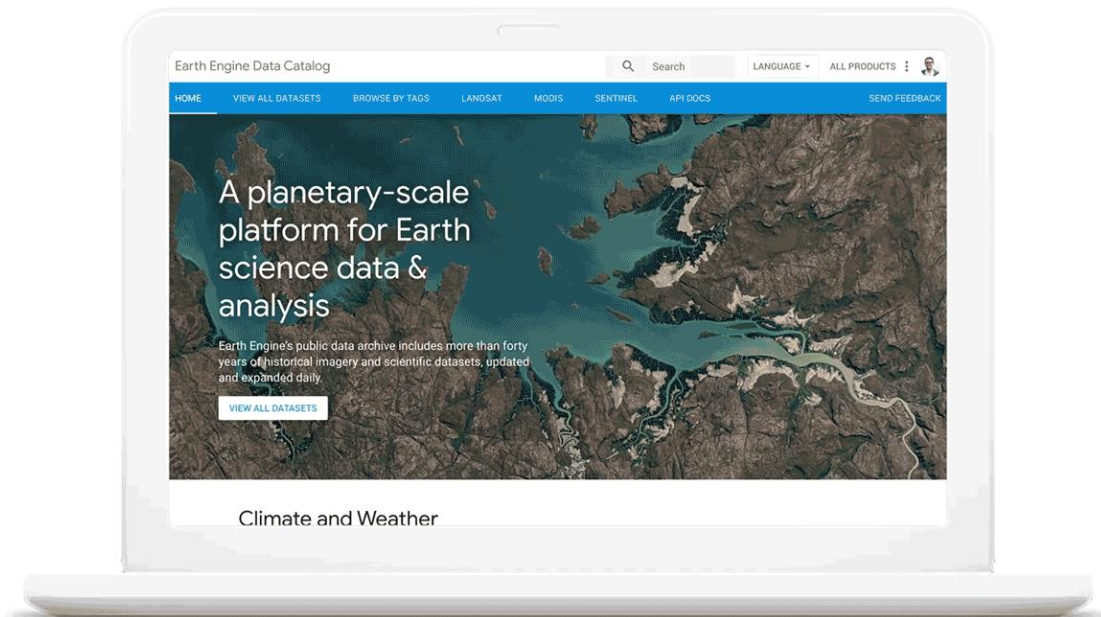
Source: NASA

# 70+ Petabytes
Growing daily

# 1 Petabyte
Monthly growth rate

# 700+
Curated datasets

# Continuously updated in near real-time



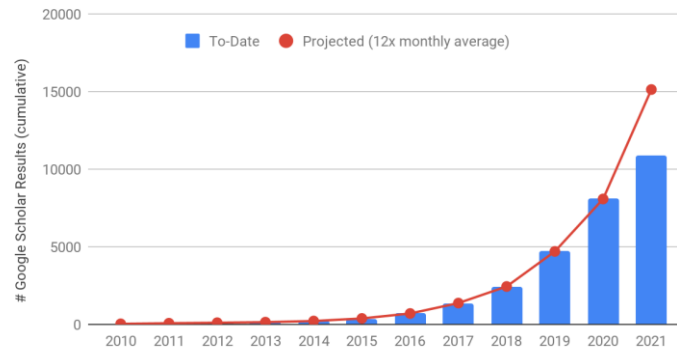developers.google.com/earth-engine/datasets/

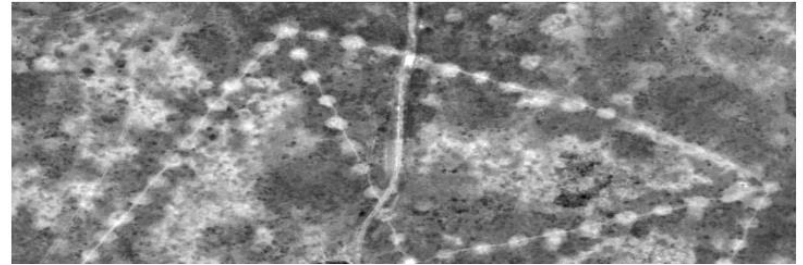# Google Earth has helped users find amazing things



*The New York Times*

*Hundreds of Mysterious Stone 'Gates' Found in Saudi Arabia's Desert*

How Google Earth helped find Mozambique's lost forest of Mount Mabu | video

*The New York Times*

*NASA Adds to Evidence of Mysterious Ancient Earthworks*

News

**Enthusiast uses Google to reveal Roman ruins**

Google Earth programme leads to remains of ancient villa.

# What is Earth Engine? | Code Editor



Your Assets

Search

Your Code

Data Inspector

API Docs

Output Console

Your scripts,
Example scripts

Batch Tasks

Drawing
Tools

Output Map

## code.earthengine.google.com

# Featured NOAA datasets



| | **Global Summary of the Day (GSOD)**<br><br>A dozen daily averages computed from global hourly station data, covering 1929 to present. | **Next Generation Radar (NEXRAD)**<br><br>High-resolution S-band Doppler weather radars operated by the National Weather Service (NWS). |
|---|---|---|
| **Global Forecast System (GEFS)**<br><br>A weather model created by the National Centers for Environmental Prediction (NCEP) that generates 21 separate forecasts to address underlying uncertainties in the input data. | **High Resolution Rapid Refresh (HRRR)**<br><br>3-km resolution hourly updated, cloud-resolving, convection-allowing atmospheric model. | **Global Historical Climatology Network (GHCN)**<br><br>Integrated database of climate summaries from land surface stations across the globe. |

Google Cloud

# 02

## Accessing and Using Public Datasets

# Google Cloud Datasets Marketplace



Google Cloud

# Where are public datasets stored?

**Google Cloud Storage** is a managed service for storing unstructured data.

Buckets contain objects (a.k.a. files and folders) that contain the data and how they're organized.

There are more than 80 buckets containing public datasets in various file formats.

**BigQuery** is Google's fully managed, serverless data warehouse for structured data.

It supports querying using a dialect of SQL.

There are more than 300 public BigQuery datasets spanning thousands of tables.

# GEFS data in the GCP Marketplace



https://console.cloud.google.com/marketplace/product/noaa-public/gfs-ensemble-forecast-system

# Access GEFS data using a browser

1. The Google Cloud Storage bucket that stores the data is gfs-ensemble-forecast-system

2. Using a web browser, access the root of the bucket with the following URI*
   https://console.cloud.google.com/storage/browser/gfs-ensemble-forecast-system

3. You can examine multiple levels of the bucket by appending the path to the URI above.

   For example, to access the path gefs.20230815/06/atmos/bufr, use the URI
   https://console.cloud.google.com/storage/browser/gfs-ensemble-forecast-system/gefs.20230815/06/atmos/bufr

*You will be asked to sign in if you are not currently signed in

https://console.cloud.google.com/storage/browser/gfs-ensemble-forecast-system

# Access GEFS data using the command line

**gsutil** is a Python application that lets you access Cloud Storage buckets and contents from the command line.

To list objects from the root of the bucket:

```
$ gsutil ls gs://gfs-ensemble-forecast-system

gs://gfs-ensemble-forecast-system/gefs.20210101/
gs://gfs-ensemble-forecast-system/gefs.20210102/
gs://gfs-ensemble-forecast-system/gefs.20210103/
gs://gfs-ensemble-forecast-system/gefs.20210104/
gs://gfs-ensemble-forecast-system/gefs.20210105/
gs://gfs-ensemble-forecast-system/gefs.20210106/
gs://gfs-ensemble-forecast-system/gefs.20210107/
gs://gfs-ensemble-forecast-system/gefs.20210108/
 ...
```

Using gcloud storage has a similar effect:

```
$ gcloud storage ls gs://gfs-ensemble-forecast-system
```

# Access GEFS data using the command line

To copy an entire prefix (directory tree) and its contents to the current directory*

```
$ gsutil -m cp gs://gfs-ensemble-forecast-system/gefs.20230812 .
```

*The -m flag enables multiprocessing to parallelize object downloads. Note that data for a single date (i.e. a `gefs.YYYYMMDD` folder) is more than 100 GB in size.

Again, using gcloud storage has a similar effect (without the -m flag):

```
$ gcloud storage cp gs://gfs-ensemble-forecast-system/gefs.20230812 .
```

(For more info, see https://cloud.google.com/sdk/gcloud/reference/storage)

# 03

## Use Cases and Journeys

# Weather Satellites: An Invaluable Resource

- Nearly unlimited use cases. Weather observations for agriculture, transportation, finance, and energy sectors
- Climate change monitoring
- Disaster/risk management
  - Wildfires
  - Extreme Floods
  - Hurricanes, Tropical Storms, and Extreme Weather

OPERATIONAL DATA
SEA LEVEL PRESSSUE (mb)  01-DAY MEAN FOR:
Tue JAN 23 2024

NOAA/Physical Sciences Laboratory

# Global Weather Models: Another Invaluable Resource

- Global weather model ensembles available
  - Use 30 perturbed + 1 control forecast to increase your certainty in how much uncertainty a model has!



- Global weather models available in high resolution 4x per day for:
  - Business Analytics
  - Operational Forecast Needs
  - ML training and validation

# Quick code example for JPSS ATMS data on GCP

## Try this on your own Colab!

## Setup

Setup your environment for all the tools you will need to accomplish your task.

```
!pip install -q zarr xarray[complete] fsspec aiohttp requests gcsfs cartopy

from google.colab import auth
from google.cloud import storage

from datetime import datetime
import xarray as xr
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import os
import fsspec
import gcsfs

auth.authenticate_user()
fs = gcsfs.GCSFileSystem(project='something')
```

## Instantiate client and fetch blobs

Start up your data set and find the data you want to work with. In this case we are working with JPSS MIRS data.

```python
# Instantiates a client
storage_client = storage.Client()

# The bucket name for the JPSS VIIRS data
bucket_name = "noaa-nesdis-n20"

def list_blobs(bucket_name, prefix, delimiter=None):
    """Lists all the blobs in the bucket."""

    storage_client = storage.Client()
    blobs = storage_client.list_blobs(bucket_name,
prefix=prefix, delimiter=delimiter)
    # Note: The call returns a response only when the
iterator is consumed.
    results = []
    for blob in blobs:
        # print(blob.name)
        results.append(blob.name)
    return results

results = list_blobs(bucket_name,
"NPR_MIRS_IMG/2023/09/15", None)

# Alternative fetching method
# results = !gsutil ls 'gs://noaa-nesdis-
n20/NPR_MIRS_IMG/2023/09/15/'

# Get rid of 33 min data for this example
results = [results[i] for i in range(len(results)) if
results[i].find('33min')<0]
print('Number of files: ', len(results))
```

Google Cloud

THE FILES ARE INSIDE THE COMPUTER

Identify the files we want to work with…

Isolate the files you want to work with and add some additional parsed data.

```python
def parse_dates(string_names):
  string_names = string_names[48:62]
  year = int(string_names[0:4])
  month = int(string_names[4:6])
  day = int(string_names[6:8])
  hour = int(string_names[8:10])
  minute = int(string_names[10:12])
  seconds= int(string_names[12:14])
  string_dt = datetime(year, month, day, hour, minute,
seconds)
  return string_dt

# Create a dataframe and add datetime field from filenames
dfr = pd.DataFrame(results, columns=['Files'])
dfr['Date'] = dfr.Files.apply(parse_dates)

#Fetch data from an important time
lets_get = dfr[(dfr.Date >= "2023-09-15 15:30:00") &
               (dfr.Date < "2023-09-15 23:30:00")]

# Get Filenames
Lets_get = lets_get.Files.to_list()
print('Filtered to:', len(lets_get))

# Show first 5 files
lets_get[0:5]
```

Google Cloud

Open and process each file

Process netcdf files and plot using xarray and matplotlib.

```python
datasets = []
for file in lets_get:
    data_path = 'gs://' + bucket_name + '/' + file
    ds3 = xr.open_dataset(fs.open(data_path), engine='h5netcdf')

    # Reduce data set to Skin temperature data
    datasets.append(ds3['TSkin'])

# Concatenate all the data to make one image from all data sources
combined = xr.concat(datasets, dim='Field_of_view')

from cartopy import config
import matplotlib.pyplot as plt
import cartopy.crs as ccrs

ax = plt.axes(projection=ccrs.PlateCarree())
combined.plot.pcolormesh(x='Longitude', y='Latitude',
cmap='rainbow', transform=ccrs.PlateCarree())

ax.coastlines()
plt.show()
```

# Quick code example for JPSS VIIRS Longwave IR data on GCP

Try this on your own Colab! https://colab.research.google.com

# Imports and Setup

Setup your kernel for all the tools you will need to accomplish these processing tasks.

```
!pip install -q zarr xarray[complete] fsspec aiohttp requests gcsfs cartopy

import h5py
import gcsfs
import matplotlib.pyplot as plt
from google.colab import auth
from google.cloud import storage
from datetime import datetime
import pandas as pd
import numpy as np
from cartopy import config
import matplotlib.pyplot as plt
import cartopy.crs as ccrs

auth.authenticate_user()
```

# Setup variables for data processing steps

Making some useful variables to limit data processing to reasonable size for processing on basic Colab.

```python
bucket_name = "noaa-nesdis-n21"
target_data = "VIIRS-I4-IMG-EDR"
target_data_geo = "VIIRS-IMG-GTM-EDR-GEO"

year = '2023'
month = '09'
day = '15'
start_hour = '17'
start_minute = '13'
end_hour = '17'
end_minute = '17'

start_limiter = datetime(int(year), int(month), int(day),
int(start_hour), int(start_minute), 0)
end_limiter = datetime(int(year), int(month), int(day),
int(end_hour), int(end_minute), 0)
fs = gcsfs.GCSFileSystem(anon=True)
```

# Instantiate client to fetch IR data blobs

Find all the relevant data from the IR sensor. In this case we are working with JPSS VIIRS IR data so we are focused on the I4 band and during our time range.

```python
# Instantiates a client
storage_client = storage.Client()

def list_blobs(bucket_name, prefix, delimiter=None):
    """Lists all the blobs in the bucket."""

    storage_client = storage.Client()
    blobs = storage_client.list_blobs(bucket_name, prefix=prefix, delimiter=delimiter)
    # Note: The call returns a response only when the iterator is consumed.
    results = []
    for blob in blobs:
        # print(blob.name)
        results.append(blob.name)
    return results

results = list_blobs(bucket_name, f"{target_data}/{year}/{month}/{day}/", None)

print('Number of files: ', len(results))

def parse_dates(s):
  s = s[38:56]
  year = int(s[1:5])
  month = int(s[5:7])
  day = int(s[7:9])
  hour = int(s[11:13])
  minute = int(s[13:15])
  seconds= int(s[15:17])
  string_dt = datetime(year, month, day, hour, minute, seconds)
  return string_dt

dfr = pd.DataFrame(results, columns=['Files'])
dfr['Date'] = dfr.Files.apply(parse_dates)

lets_get = dfr[(dfr.Date >= start_limiter) & (dfr.Date < end_limiter)].Files.to_list()
print('Filtered to:', len(lets_get))
lets_get[0:5]
```

## Instantiate client to fetch geo blobs

Start up your data set and find the data you want to work with. In this case we are working with grabbing GEO data.

```python
# Instantiates a client
storage_client = storage.Client()

def list_blobs(bucket_name, prefix, delimiter=None):
    """Lists all the blobs in the bucket."""

    storage_client = storage.Client()
    blobs = storage_client.list_blobs(bucket_name, prefix=prefix, delimiter=delimiter)
    # Note: The call returns a response only when the iterator is consumed.
    results = []
    for blob in blobs:
        # print(blob.name)
        results.append(blob.name)
    return results

results_geo = list_blobs(bucket_name, f"{target_data_geo}/{year}/{month}/{day}/", None)

print('Number of files: ', len(results))

def parse_dates(s):
  s = s[43:61]
  year = int(s[1:5])
  month = int(s[5:7])
  day = int(s[7:9])
  hour = int(s[11:13])
  minute = int(s[13:15])
  seconds= int(s[15:17])
  string_dt = datetime(year, month, day, hour, minute, seconds)
  return string_dt

dfr_geo = pd.DataFrame(results_geo, columns=['Files'])
dfr_geo['Date'] = dfr_geo.Files.apply(parse_dates)

lets_get_geo = dfr_geo[(dfr_geo.Date >= start_limiter) & (dfr_geo.Date <
end_limiter)].Files.to_list()
print('Filtered to:', len(lets_get_geo))
lets_get_geo[0:5]
```

Google Cloud

Process all the data you need for your project.



Process data to numpy arrays (online) and prep for visualization.

```python
agg_data_lat = []
agg_data_lon = []

# One file is ALOT of data
for file in lets_get_geo[1:2]:
  r = 'gs://' + bucket_name + '/' + file
  print(r)
  f = h5py.File(fs.open(r), 'r')
  latitude = f['All_Data']['VIIRS-IMG-GTM-EDR-GEO_All']['Latitude'][:]
  print(np.array(latitude).shape)
  longitude = f['All_Data']['VIIRS-IMG-GTM-EDR-GEO_All']['Longitude'][:]
  print(np.array(longitude).shape)

  # If processing multiple files at once
  for row in range(latitude.shape[0]):
    agg_data_lat.append(latitude[row])
  for row in range(longitude.shape[0]):
    agg_data_lon.append(longitude[row])

agg_data = []
counter = 0

for file in lets_get[1:2]:
  r = 'gs://' + bucket_name + '/' + file
  print(r)
  f = h5py.File(fs.open(r), 'r')
  if counter < 1:
    products = f['All_Data']['VIIRS-I4-IMG-EDR_All']

  arr = f['All_Data']['VIIRS-I4-IMG-EDR_All']['BrightnessTemperature'][:]
  print(np.array(arr).shape)
  for row in range(0, arr.shape[0]):
    agg_data.append(arr[row])

 longitude[longitude<-900] = np.nan
 latitude[latitude<-900] = np.nan
```
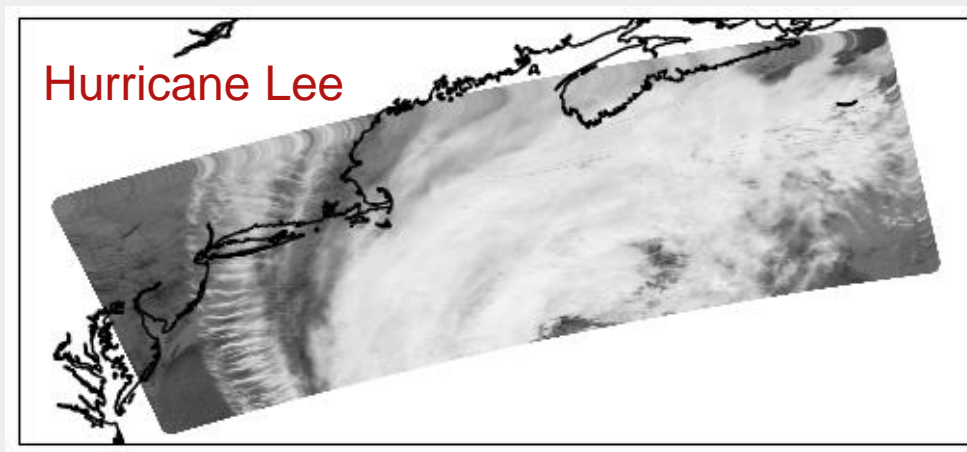
# Visualize Numpy Arrays

Plot numpy arrays using matplotlib and cartopy.

```python
from cartopy import config
import matplotlib.pyplot as plt
import cartopy.crs as ccrs

ax = plt.axes(projection=ccrs.PlateCarree())
plt.scatter(longitude[:, 2500:6500], latitude[:, 2500:6500],
c=arr[:, 2500:6500], cmap='Greys',
transform=ccrs.PlateCarree())

ax.coastlines()
plt.show()
```



Hurricane Lee

# Thank you.

Questions? Email us at:
cloud-public-dataset-conferences@google.com

Google Cloud