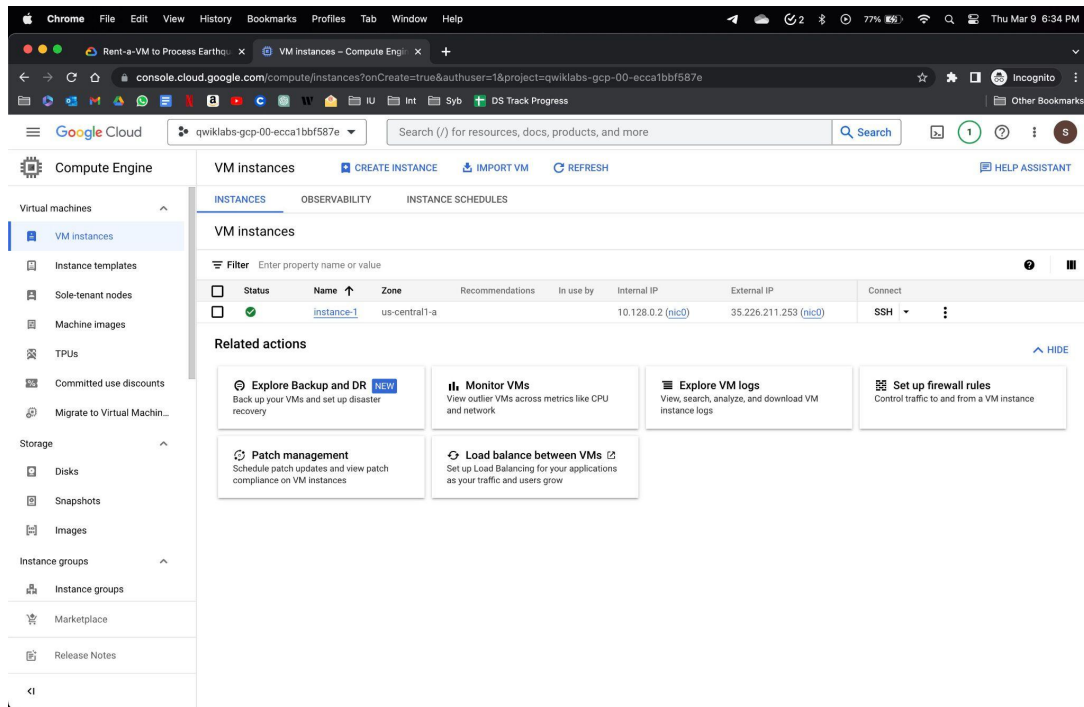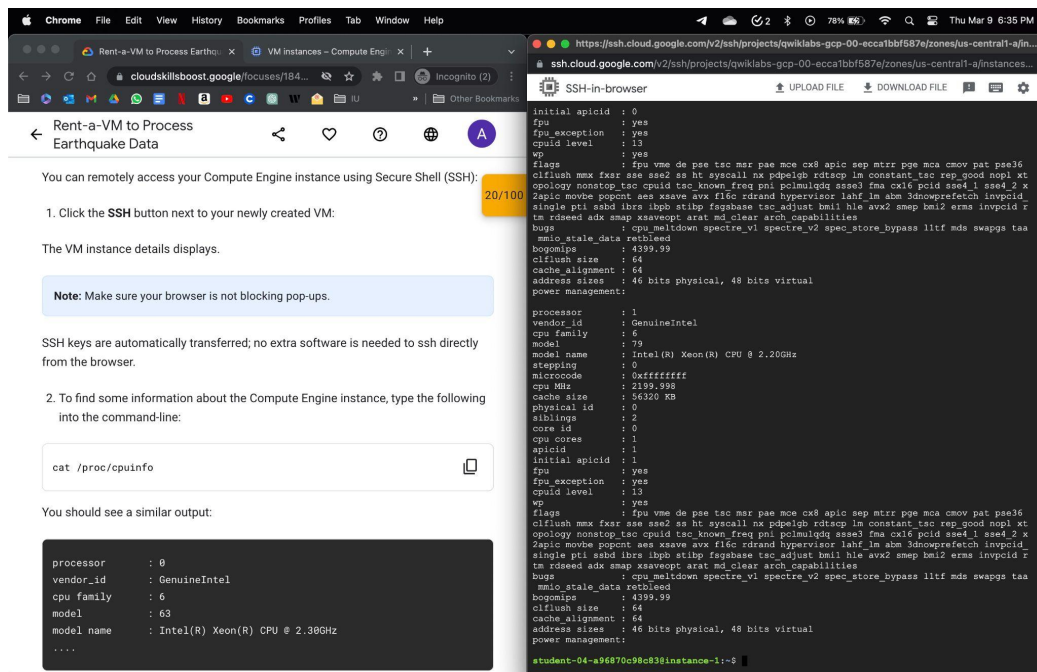# Lifecycle and pipelines

- **Aditya Sanjay Mhaske**

## Rent-a-VM to Process Earthquake Data

### Task 1. Create Compute Engine instance with the necessary API access



### Task 2. SSH into the instance

Task 3. Install software
Task 4. Ingest USGS data
Task 5. Transform the data
Task 6. Create a Cloud Storage bucket



Task 7. Store data



Task 8. Publish Cloud Storage files to web
Congratulations!

**Describe the high-level architecture of this pipeline by covering the following:**

- Data sources and types
- Data storage technologies
- Data transformations tools and techniques (make sure to look into the Python code to understand that)
- Other data components that are included (e.g., visualization, reporting, etc.)

1. Data sources and types:

The pipeline processes earthquake data from the US Geological Survey (USGS) website, which provides real-time data on earthquakes worldwide. The data is in the form of CSV files containing information such as each earthquake's date, time, location, magnitude, and depth.

2. Data storage technologies:

The pipeline uses several Google Cloud storage technologies, including Google Cloud Storage (GCS), Google Bigtable, and Google Pub/Sub. The CSV files are stored in GCS buckets, and the pipeline uses Bigtable to store and manage earthquake metadata. Pub/Sub is used to decouple the different stages of the pipeline and facilitate communication between them.

3. Data transformations tools and techniques:

The pipeline uses a combination of Python libraries and Google Cloud services to transform and process earthquake data. The main Python libraries used are pandas, NumPy, and Matplotlib, which are used for data manipulation, processing, and visualization. The pipeline uses Google Dataflow to parallelize the data processing and apply transformations such as filtering, grouping, and aggregating. Additionally, the pipeline uses Google Cloud Functions to trigger the data processing when new data is available in the GCS buckets.

4. Other data components included:

The pipeline includes several data components, including data visualization and reporting. Matplotlib is used to create various plots and visualizations of the earthquake data, such as scatter plots and histograms. The pipeline also generates earthquake reports, including summaries of the earthquake data by magnitude and location. The reports are stored in GCS buckets and can be accessed through a web application built using Google App Engine.

**This pipeline covers several aspects of the data lifecycle, including:**

In this week's module, the USGS model was explained. The following steps were taken:

- Defining the data plan, which involved identifying the required process and resources.
- Collecting the data from the USGS website into the VM.
- Processing the data using the transform.py file and implementing the mentioned transformations.
- Analyzing the data by creating visualizations on a map.
- Storing the data on the bucket.
- Ensuring accessibility of the data by making it available on the bucket and generating an earthquakes.htm file to facilitate sharing/publishing.



Citation:
1. https://www.usgs.gov/data-management/data-lifecycle
2. https://docs.gcp.databricks.com/getting-started/data-pipeline-get-started.html

**Discuss what could be added to the pipeline using this module's materials.**

There are several components that could be added to this pipeline using the Google Cloud Module's materials:

1.  Real-time data streaming: Currently, the pipeline processes earthquake data that is collected periodically from the USGS website. By adding real-time data streaming capabilities using Google Cloud Pub/Sub, the pipeline could process and analyze earthquake data in real time, which would be useful for disaster response and emergency management.

2.  Machine learning models: The pipeline could be enhanced by adding machine learning models that could predict the likelihood and severity of earthquakes based on historical data. Google Cloud's AutoML and AI Platform could be used to train and deploy these models.

3.  Anomaly detection: The pipeline could include an anomaly detection system to detect unusual patterns or trends in the earthquake data, which could be an indication of seismic activity that requires further investigation.

4.  Integration with other data sources: The pipeline could be expanded to include data from other sources, such as weather data or satellite imagery, to provide more context and insights into the earthquake data.