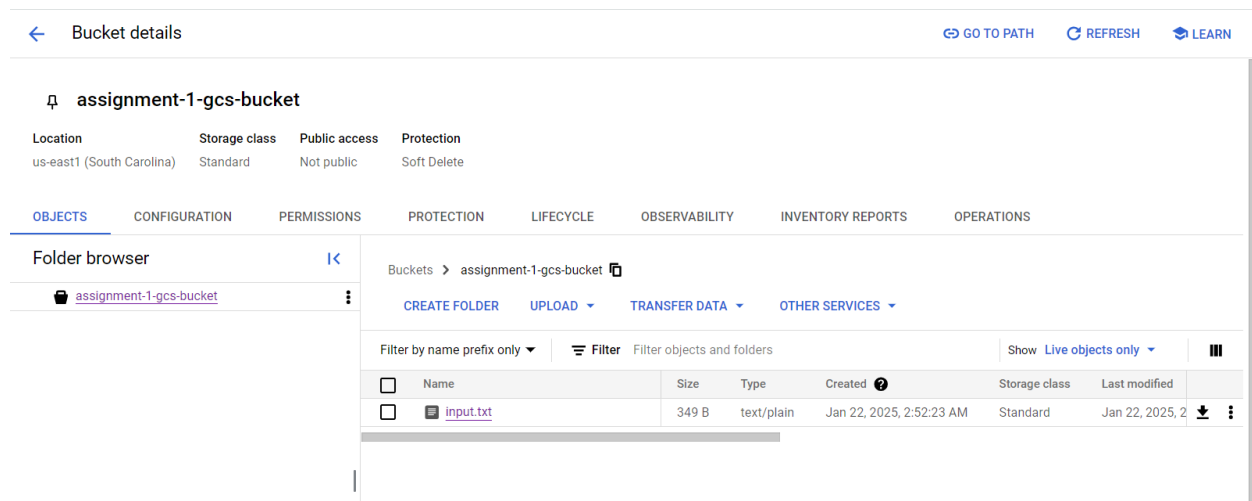# Question : Spin Up a VM and write a python program to count lines of a file placed in GCS.

**Submitted By: Shramana Sinha, 23F1002703**

# 1. Setup and Prerequisites

## 1.1 Creating a GCS Bucket

1.  A GCS bucket named "assignment-1-gcs-bucket" was created with default settings, except for data storage location set to "Lowest latency within a single region".
2.  A file named "input.txt" was uploaded to the created bucket.



*Screenshot 1: GCS bucket and file*

## 1.2 Spinning up a Virtual Machine

A Virtual Machine (VM) was created with default settings, except for using Ubuntu 22.04 LTS as the Operating System.

*Screenshot 2: VM creation settings*

## 1.3 Installing Required Software

The necessary softwares were installed on the VM:
- Package lists were updated: `sudo apt update`
- Python 3 and pip were installed: `sudo apt install python3 python3-pip`
- Google Cloud Storage library was installed: `pip3 install google-cloud-storage`
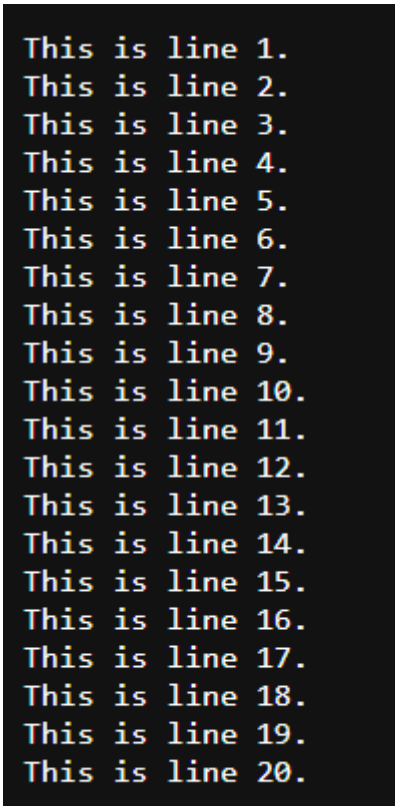
# 2. Code Execution and Results

The script was executed using the command:
`python3 line_counter.py`



*Screenshot 3: Terminal showing the execution of the script and its output*

The output matches the number of lines in the input file.

*Screenshot 4: File showing the number of lines*

# 3. Code Explanation

The Python script consists of two main functions:

## 3.1 count_lines_in_gcs(bucket_name: str, blob_name: str) -> int

This function is responsible for connecting to GCS, downloading the specified file, and counting its lines. Here's a breakdown of its operations:

1. Initialize the GCS client using `storage.Client()`. It uses the default authentication method, which assumes that the VM has the necessary permissions to access the GCS bucket.
2. Get the specified bucket and blob (file) using the provided names.
3. Download the file content as text using `blob.download_as_text()`.
4. Split the text content into a list of lines using `splitlines()`. This method handles different line ending characters (`\\n`, `\\r`, or `\\r\\n`) automatically.
5. Count the number of lines using `len()`.
6. Return the line count.

## 3.2 main()

This function serves as the entry point of the script:

1. Define the bucket name and blob name (file path).
2. Call `count_lines_in_gcs()` with the specified bucket and blob names.
3. Print the result or any error messages.