

# Virtualization and Cloud computing

Amresh Kumar M23CSA004

## Assignment- 03

---

### Github

link: [https://github.com/m23csa004/vcc\\_assignment3](https://github.com/m23csa004/vcc_assignment3)

Video link: 📺 m23csa004\_vcc\_assign3.mkv

### 1. Introduction

This report outlines the step-by-step implementation of a local VM-based resource monitoring system that triggers autoscaling in Google Cloud Platform (GCP). The system monitors CPU utilization and automatically provisions additional instances when the CPU usage crosses a defined threshold.

I have used **GCP** as Public cloud and **virtual box** as hypervisor and **psutil** as a resource monitoring tool.

### 2. Step-by-Step Implementation

#### Step 1: Setting Up a Local VM

We begin by creating a virtual machine (VM) locally using VirtualBox:

- Download and install **VirtualBox**.
- Use an **Ubuntu Desktop 24.04** ISO image to install Ubuntu.
- Complete the installation and upgrade **Guest Additions** for better performance.

#### Step 2: Install Google Cloud SDK

- Install Google Cloud SDK: **sudo snap install google-cloud-sdk --classic**
- Verify installation: **gcloud version**
- Authenticate with Google Cloud: **gcloud auth login**
- Get project lists and project id: **gcloud projects list**

I have used the same project as i have used in the assignment 2 My First Project

- Set project configuration: **gcloud config set project <your-project-id>**
- Enable Compute Engine API: **gcloud services enable compute.googleapis.com**

```

Building dependency tree... Done
Reading state information... Done

No apt package "google-cloud-sdk", but there is a snap with that name.
Try 'snap install google-cloud-sdk'

E: Unable to locate package google-cloud-sdk
Amresh@Amresh:~/assign3$ sudo snap install google-cloud-sdk --classic
google-cloud-sdk 515.0.0 from Cloud SDK (google-cloud-sdk) installed
Amresh@Amresh:~/assign3$ gcloud version
Google Cloud SDK 515.0.0
alpha 2025.03.14
beta 2025.03.14
bq 2.1.14
bundled-python3-unix 3.12.8
core 2025.03.14
gcloud-crc32c 1.0.0
gsutil 5.33
minikube 1.35.0
skaffold 2.14.1
Amresh@Amresh:~/assign3$ gcloud config list
[core]
disable_usage_reporting = True

Your active configuration is: [default]
Amresh@Amresh:~/assign3$ gcloud auth login
Your browser has been opened to visit:

https://accounts.google.com/o/oauth2/auth?response_type=code&client_id=32555940559.apps.googleusercontent.com&redirect_url=http%3A%2F%2Flocalhost%3A8085%2F&scope=openid+https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fuserinfo.email+https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fcloud-platform+https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fappengine.admin+https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fsqlservice.login+https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fcompute+https%3A%2F%2Fwww.googleapis.com%2Fauth%2Faccounts.reauth&state=61TD7kVkp2EUn27nZQptFuZjXUZM2&access_type=offline&code_challenge=QoWe2ZP7pZRYNNJCZycl6_cWDPlLpz5yloJYKZ1Mv4&code_challenge_method=S256

You are now logged in as [m23csa004@itj.ac.in].
Your current project is [None]. You can change this setting by running:
$ gcloud config set project PROJECT_ID
Amresh@Amresh:~/assign3$
  
```

## Step 3: Setting Up Python Environment

- Install Python and required packages:
  - **sudo apt update && sudo apt install -y python3-pip python3-venv**
- Create a virtual environment:
  - **python3 -m venv venv**
  - **source venv/bin/activate**
- Install dependencies:
  - **pip install flask psutil google-auth-oauthlib google-auth-http2 google-api-python-client**
  - **Flask:** A lightweight web framework for building web applications in Python.
  - **psutil:** A library for retrieving system and process-related information such as CPU, memory, and disk usage.
  - **google-auth-oauthlib:** A library for handling OAuth 2.0 authentication in Python applications.
  - **google-auth-http2:** Provides authentication support for http2, enabling secure HTTP requests using Google authentication.

- **google-api-python-client**: A client library for accessing various Google APIs, such as Drive, Gmail, and Sheets.

## Step 4: Implementing Resource Monitoring

We develop a **Flask-based monitoring application** (`app.py`) that:

- Monitors CPU usage.
- Runs a background thread that checks CPU utilization.
- Triggers the GCP VM creation script when CPU usage exceeds **75%**.

## Step 5: Implementing Auto-Scaling Script

We create a Python script (`create_gcp_vm.py`) that

- Defines a **startup script** for VM instances.
- Creates an **Instance Group** with autoscaling rules.
- Uses **gcloud commands** to launch new instances when CPU utilization exceeds 75%.

## Steps to Generate the Autoscaling Command

### 1. Create an Instance Template in GCP:

- Go to **Compute Engine > Instance Templates** in the GCP Console.
- Click **Create Instance Template** and configure the following:
  - **Name**: Choose a unique identifier.
  - **Region**: Select your preferred region.
  - **Machine Configuration**: Choose **e2-micro** for cost efficiency.
  - **Firewall Settings**: Enable **Allow HTTP and HTTPS traffic**.
- Under **Advanced > Management > Automation**, paste the **startup script** that initializes the instance.
- Click **Create**.

## 2. Create a Managed Instance Group:

- Navigate to **Compute Engine > Instance Groups**.
- Click **Create Instance Group** and select the previously created **Instance Template**.
- Enable **Autoscaling** with the following parameters:
  - **Minimum instances:** 1
  - **Maximum instances:** 5
  - **CPU utilization threshold:** 75%
- Save the configuration.

## 3. Generate the Equivalent gcloud Command:

- Once the instance group is configured, click **Equivalent Code** (top-right corner of the GCP Console).
- Copy the generated command and include it in the Python script.

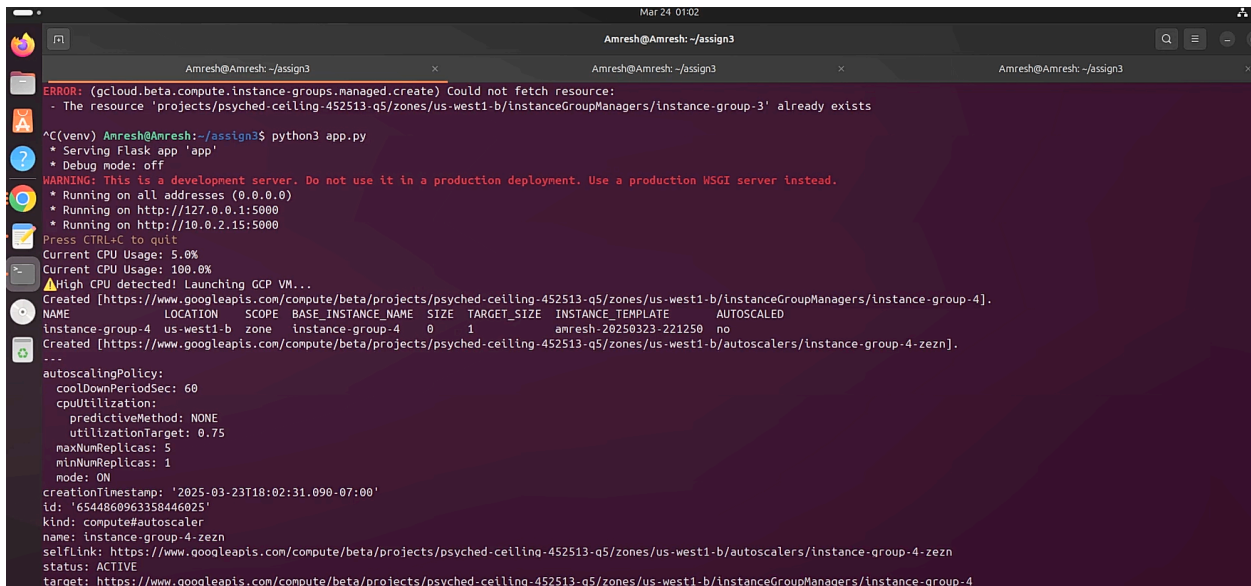
This ensures that the autoscaling behavior is dynamically managed based on CPU usage.

The screenshot shows the Google Cloud Platform console interface for creating a managed instance group. The left sidebar displays the navigation menu with 'Compute Engine / Instance Groups / Create instance group' selected. The main content area is titled 'Create Instance Group' and includes a 'Health check' dropdown menu, a 'Default action' dropdown menu, and a section for 'Updates during VM instance repair' with two radio button options: 'Keep the same instance configuration' (selected) and 'Update the instance configuration'. Below this is a 'Port mapping' section with a '+ Add port' button and a 'Show advanced configuration' link. On the right side, the 'Equivalent code' tab is active, showing a gcloud command for creating a managed instance group with autoscaling. The command is as follows:

```
1 gcloud beta compute instance-groups managed create instance-group-4 \
2 --project=psyched-ceiling-452513-q5 \
3 --base-instance-name=instance-group-4 \
4 --template=projects/psyched-ceiling-452513-q5/regions/us-west1/instanceTemplates/
5 amresh-20250323-221258 \
6 --size=1 \
7 --zone=us-west1-b \
8 --default-action-on-vn-failure=repair \
9 --action-on-vn-failed-health-check=default-action \
10 --no-force-update-on-repair \
11 --standby-policy-mode=manual \
12 --list-managed-instances-result=pageless \
13 && \
14 gcloud beta compute instance-groups managed set-autoscaling instance-group-4 \
15 --project=psyched-ceiling-452513-q5 \
16 --zone=us-west1-b \
17 --mode=on \
18 --min-num-replicas=1 \
19 --max-num-replicas=5 \
20 --target-cpu-utilization=0.75 \
21 --cpu-utilization-predictive-method=none \
22 --cool-down-period=300
```

## Step 6: Running the Application Locally

- Start the Flask app: **python3 app.py**
- Open another terminal and simulate high CPU usage:
  - **sudo apt install -y stress**
  - **stress --cpu 4 --timeout 60s**
  -
- Observe that CPU usage spikes and triggers the GCP VM creation script.



```
Mar 24 01:02
Amresh@Amresh: ~/assign3
ERROR: (gcloud.beta.compute.instance-groups.managed.create) Could not fetch resource:
- The resource 'projects/psyched-ceiling-452513-q5/zones/us-west1-b/instanceGroupManagers/instance-group-3' already exists
^C(venv) Amresh@Amresh:~/assign3$ python3 app.py
* Serving Flask app 'app'
* Debug mode: off
WARNING: This is a development server. Do not use it in a production deployment. Use a production WSGI server instead.
* Running on all addresses (0.0.0.0)
* Running on http://127.0.0.1:5000
* Running on http://10.0.2.15:5000
Press CTRL+C to quit
Current CPU Usage: 5.0%
Current CPU Usage: 100.0%
High CPU detected! Launching GCP VM...
Created [https://www.googleapis.com/compute/beta/projects/psyched-ceiling-452513-q5/zones/us-west1-b/instanceGroupManagers/instance-group-4].
NAME LOCATION SCOPE BASE_INSTANCE_NAME SIZE TARGET_SIZE INSTANCE_TEMPLATE AUTOSCALED
instance-group-4 us-west1-b zone instance-group-4 0 1 amresh-20250323-221250 no
Created [https://www.googleapis.com/compute/beta/projects/psyched-ceiling-452513-q5/zones/us-west1-b/autoscalers/instance-group-4-zezn].
...
autoscalingPolicy:
  coolDownPeriodSec: 60
  cpuUtilization:
    predictiveMethod: NONE
    utilizationTarget: 0.75
  maxNumReplicas: 5
  minNumReplicas: 1
  mode: ON
creationTimestamp: '2025-03-23T18:02:31.090-07:00'
id: '6544860963358446025'
kind: compute#autoscaler
name: instance-group-4-zezn
selfLink: https://www.googleapis.com/compute/beta/projects/psyched-ceiling-452513-q5/zones/us-west1-b/autoscalers/instance-group-4-zezn
status: ACTIVE
target: https://www.googleapis.com/compute/beta/projects/psyched-ceiling-452513-q5/zones/us-west1-b/instanceGroupManagers/instance-group-4
```

## Why Simulate Load?

Load testing helps verify the effectiveness of auto-scaling policies and ensures that the infrastructure can handle high traffic scenarios without performance degradation.

I have performed only a cpu load test, other load tests like memory, network and disk can be performed.

## Step 7: Verifying Autoscaling in GCP

- Check if the instance group is created: **gcloud compute instance-groups managed list**
- List instances within the group:
  - **gcloud compute instance-groups managed list-instances instance-group-4 --zone=us-west1-b**
- Connect to an instance using SSH:

- **gcloud compute ssh <instance-name> --zone=us-west1-b**
- Monitor autoscaling activity:
  - **gcloud compute instance-groups managed describe instance-group-4 --zone=us-west1-b**

<input type="checkbox"/> Status	Name ↑	Zone	Recommendations	In use by	Internal IP	External IP	Connect
<input type="checkbox"/>	<a href="#">instance-2</a>	asia-south2-c			10.190.0.2 ( <a href="#">nic0</a> )		SSH ▾ ⋮
<input type="checkbox"/>	<a href="#">instance-20250302-140538</a>	us-central1-c			10.128.0.2 ( <a href="#">nic0</a> )		SSH ▾ ⋮
<input type="checkbox"/>	<a href="#">instance-group-1-cpgs</a>	us-central1-c		<a href="#">instance-group-1</a>	10.128.0.24 ( <a href="#">nic0</a> )		SSH ▾ ⋮
<input type="checkbox"/>	<a href="#">instance-group-1-nsgd</a>	us-central1-c		<a href="#">instance-group-1</a>	10.128.0.51 ( <a href="#">nic0</a> )		SSH ▾ ⋮
<input type="checkbox"/>	<a href="#">instance-group-2-4psb</a>	asia-south2-a		<a href="#">instance-group-2</a>	10.190.0.8 ( <a href="#">nic0</a> )		SSH ▾ ⋮
<input type="checkbox"/>	<a href="#">instance-group-2-nqtg</a>	asia-south2-a		<a href="#">instance-group-2</a>	10.190.0.12 ( <a href="#">nic0</a> )		SSH ▾ ⋮
<input type="checkbox"/>	<a href="#">instance-group-3-59qb</a>	us-west1-b		<a href="#">instance-group-3</a>	10.138.0.6 ( <a href="#">nic0</a> )	<a href="#">35.197.85.140</a> ( <a href="#">nic0</a> )	SSH ▾ ⋮
<input type="checkbox"/>	<a href="#">instance-group-4-24jt</a>	us-west1-b		<a href="#">instance-group-4</a>	10.138.0.14 ( <a href="#">nic0</a> )	<a href="#">34.82.214.156</a> ( <a href="#">nic0</a> )	SSH ▾ ⋮
<input type="checkbox"/>	<a href="#">instance-group-4-446l</a>	us-west1-b		<a href="#">instance-group-4</a>	10.138.0.16 ( <a href="#">nic0</a> )	<a href="#">34.118.199.32</a> ( <a href="#">nic0</a> )	SSH ▾ ⋮
<input type="checkbox"/>	<a href="#">instance-group-4-dr4h</a>	us-west1-b		<a href="#">instance-group-4</a>	10.138.0.9 ( <a href="#">nic0</a> )	<a href="#">34.169.192.204</a> ( <a href="#">nic0</a> )	SSH ▾ ⋮
<input type="checkbox"/>	<a href="#">instance-group-4-lrb5</a>	us-west1-b		<a href="#">instance-group-4</a>	10.138.0.13 ( <a href="#">nic0</a> )	<a href="#">34.127.52.119</a> ( <a href="#">nic0</a> )	SSH ▾ ⋮
<input type="checkbox"/>	<a href="#">instance-group-4-s8ql</a>	us-west1-b		<a href="#">instance-group-4</a>	10.138.0.15 ( <a href="#">nic0</a> )	<a href="#">34.105.114.45</a> ( <a href="#">nic0</a> )	SSH ▾ ⋮

Rows per page:

Related actions

## Step 8: Testing Deployment

1. Visit `http://<instance-external-ip>:5000` to check the application response.
2. Verify autoscaling by running the **stress** command and observing instance scaling behavior.

## 3. Conclusion

This implementation successfully demonstrates:

- Local VM-based CPU monitoring.
- Automated VM provisioning in GCP upon high CPU utilization.
- Deployment of a Flask-based application on auto-scaled instances.
- Firewall configuration for external accessibility.

By following these steps, we can efficiently scale resources dynamically based on workload demands, ensuring optimized cloud resource utilization.

You can also use tools like Prometheus, Grafana or node exporter for better graphics.

### 3. Architecture Design:

