# Assignment - 2

# GCP VM Auto-Scaling & Security Configuration

**Submitted By**

Anuj CHINCHOLIKAR
B22ES018

**March 1, 2025**

# Contents

# 1    Introduction

> **Project Overview**
>
> This report documents the implementation of a Google Cloud Platform (GCP) Virtual Machine (VM) with auto-scaling policies and robust security measures. The project aims to create a scalable and secure infrastructure that can adapt to changing workload demands while maintaining strict access controls.

## 1.1    Project Objectives

- Create and configure a VM instance with Apache web server

- Implement auto-scaling using Managed Instance Groups (MIG)

- Configure security measures including firewall rules and IAM roles

- Test the implementation using various tools and techniques

## 1.2    Technical Stack

- Google Cloud Platform (GCP)

- Bash scripts for automation

- Apache HTTP Server

- Testing tools: curl, nmap, stress

- Optional: Terraform for Infrastructure-as-Code

# 2    Architecture Overview

> The architecture consists of a Managed Instance Group (MIG) with an underlying instance template that defines VM configurations. Auto-scaling policies are applied to the MIG to dynamically adjust the number of instances based on CPU utilization. Security is enforced through firewall rules that restrict access to specific IP addresses and IAM roles that provide controlled access to GCP resources.

## 2.1 System Components

- VM Instance: Base infrastructure running Apache web server

- Instance Template: Blueprint for creating identical VM instances

- Managed Instance Group: Group of VMs managed as a single entity

- Auto-scaling Policy: Rules for scaling instances based on CPU utilization

- Firewall Rules: Network security rules to control traffic

- IAM Roles: Identity and access management controls

# 3 Implementation

## 3.1 Creating a VM Instance

To create a VM instance with Apache installed, we use the following script:

```bash
#!/bin/bash
# Create a single VM instance on GCP with Apache installed

# Set your project ID and zone
PROJECT_ID="your-project-id"
ZONE="us-central1-a"
INSTANCE_NAME="assignment-2-gcp-vm"
MACHINE_TYPE="e2-medium"
IMAGE_FAMILY="ubuntu-2004-lts"
IMAGE_PROJECT="ubuntu-os-cloud"

# Set the active project
gcloud config set project $PROJECT_ID

# Create the VM instance with HTTP and HTTPS tags
gcloud compute instances create $INSTANCE_NAME \
  --zone=$ZONE \
  --machine-type=$MACHINE_TYPE \
  --image-family=$IMAGE_FAMILY \
  --image-project=$IMAGE_PROJECT \
  --tags=http-server,https-server \
  --metadata=startup-script='#!/bin/bash
    apt-get update
    apt-get install -y apache2
    systemctl start apache2
    systemctl enable apache2'

echo "VM instance $INSTANCE_NAME created with Apache installed."
```

Listing 1: Create VM with Apache installed

> ### VM Configuration Details
>
> - Ubuntu 20.04 LTS as the operating system
>
> - e2-medium machine type (2 vCPUs, 4 GB memory)
>
> - HTTP and HTTPS network tags for firewall rule targeting
>
> - A startup script that installs and configures Apache



Figure 1: Apache Installation on VM Instance

## 3.2    Configuring Auto-Scaling

Auto-scaling is implemented in two steps:

**Step 1:** Creating an instance template

**Step 2:** Creating a Managed Instance Group with auto-scaling policies

### 3.2.1    Creating an Instance Template

```bash
#!/bin/bash
# Create an instance template for auto-scaling

PROJECT_ID="your-project-id"
TEMPLATE_NAME="assignment-2-template"
MACHINE_TYPE="e2-medium"
IMAGE_FAMILY="ubuntu-2004-lts"
IMAGE_PROJECT="ubuntu-os-cloud"

gcloud config set project $PROJECT_ID

gcloud compute instance-templates create $TEMPLATE_NAME \
  --machine-type=$MACHINE_TYPE \
  --image-family=$IMAGE_FAMILY \
  --image-project=$IMAGE_PROJECT \
  --tags=http-server,https-server \
  --metadata=startup-script='#!/bin/bash
```

```
18      apt-get update
19      apt-get install -y apache2
20      systemctl start apache2
21      systemctl enable apache2'
22
23 echo "Instance template $TEMPLATE_NAME created successfully."
```

Listing 2: Create an instance template for auto-scaling
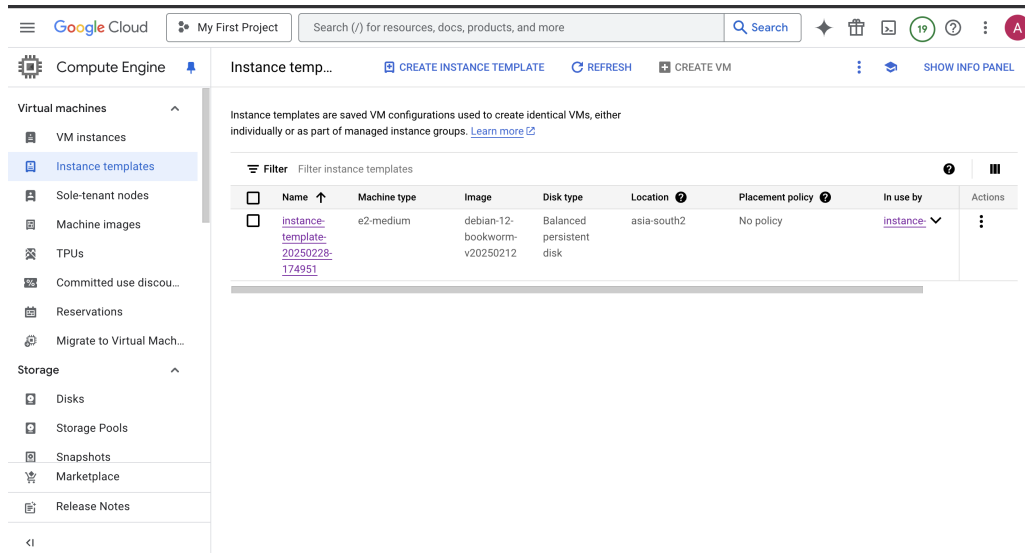


Figure 2: Created Instance Template in GCP Console

### 3.2.2    Creating a Managed Instance Group with Auto-Scaling

```
1  #!/bin/bash
2  # Create a Managed Instance Group (MIG) with auto-scaling
3
4  PROJECT_ID="your-project-id"
5  ZONE="us-central1-a"
6  MIG_NAME="assignment-2-mig"
7  TEMPLATE_NAME="assignment-2-template"
8
9  gcloud config set project $PROJECT_ID
10
11 # Create the Managed Instance Group
12 gcloud compute instance-groups managed create $MIG_NAME \
13   --base-instance-name=assignment-2-mig-instance \
14   --template=$TEMPLATE_NAME \
15   --size=1 \
16   --zone=$ZONE
17
18 # Configure auto-scaling: target CPU utilization 60%, min 1 instance, max 5
      instances, cooldown 60s
19 gcloud compute instance-groups managed set-autoscaling $MIG_NAME \
20   --zone=$ZONE \
21   --cool-down-period=60 \
22   --max-num-replicas=5 \
23   --min-num-replicas=1 \
24   --target-cpu-utilization=0.6
25
26 echo "Managed Instance Group $MIG_NAME created with auto-scaling configured.
      "
```

```
27  echo "Auto-scaling parameters:"
28  echo "  - Target CPU utilization: 60%"
29  echo "  - Min instances: 1"
30  echo "  - Max instances: 5"
31  echo "  - Cooldown period: 60 seconds"
```

Listing 3: Create a Managed Instance Group with auto-scaling

**Auto-Scaling Policy Parameters**

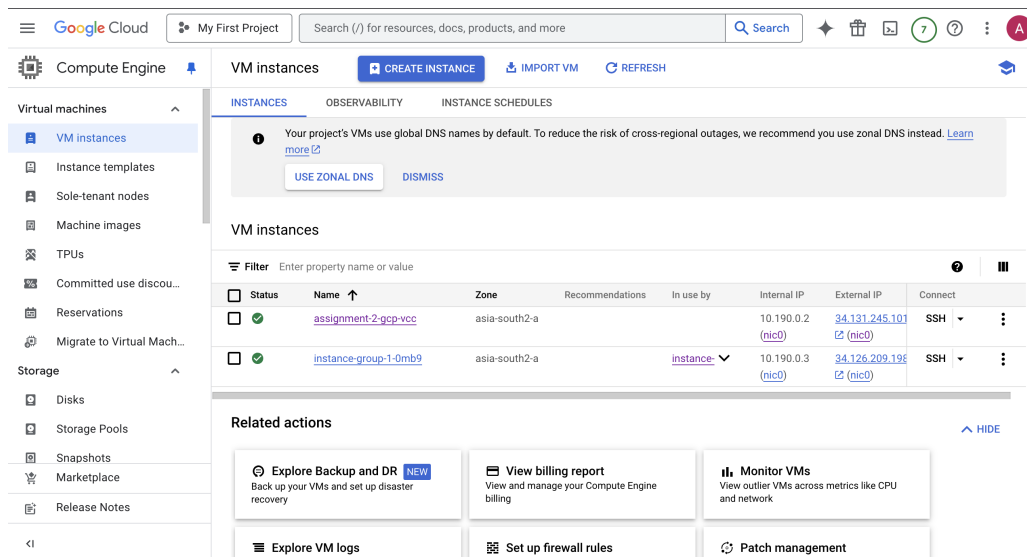| Parameter | Value |
|---|---|
| Target CPU utilization | 60% |
| Minimum instances | 1 |
| Maximum instances | 5 |
| Cooldown period | 60 seconds |



Figure 3: Managed Instance Group with Initial VM Instance

## 3.3   Implementing Security Measures

### 3.3.1   Firewall Rules

```
1  #!/bin/bash
2  # Set up firewall rules to allow HTTP/HTTPS traffic only from a specific IP
       (223.238.204.234/32)
3
4  PROJECT_ID="your-project-id"
5  ALLOWED_IP="223.238.204.234/32"
6
7  gcloud config set project $PROJECT_ID
8
9  # Create firewall rule for HTTP traffic
10 gcloud compute firewall-rules create allow-http-custom \
11   --direction=INGRESS \
12   --priority=1000 \
13   --network=default \
14   --action=ALLOW \
```

```
15    --rules=tcp:80 \
16    --source-ranges=$ALLOWED_IP \
17    --target-tags=http-server
18
19 # Create firewall rule for HTTPS traffic
20 gcloud compute firewall-rules create allow-https-custom \
21    --direction=INGRESS \
22    --priority=1000 \
23    --network=default \
24    --action=ALLOW \
25    --rules=tcp:443 \
26    --source-ranges=$ALLOWED_IP \
27    --target-tags=https-server
28
29 echo "Firewall rules created successfully."
30 echo "HTTP (port 80) and HTTPS (port 443) traffic is now allowed only from
      IP: $ALLOWED_IP"
```

Listing 4: Set up firewall rules for restricted access

### Firewall Configuration

- Allows HTTP (port 80) and HTTPS (port 443) traffic only from the specified IP address (223.238.204.234/32)

- Targets instances with the "http-server" and "https-server" tags

- Sets a priority of 1000 for the rules

### 3.3.2   IAM Roles

```
1 #!/bin/bash
2 # Set up IAM roles to grant restricted access
3
4 PROJECT_ID="your-project-id"
5 USER_EMAIL="user@example.com"  # Replace with the email of the user
6
7 # Grant the Compute Viewer role to the specified user
8 gcloud projects add-iam-policy-binding $PROJECT_ID \
9    --member="user:$USER_EMAIL" \
10   --role="roles/compute.viewer"
11
12 echo "IAM role 'Compute Viewer' assigned to $USER_EMAIL"
13 echo "Note: This role provides read-only access to all compute resources"
```

Listing 5: Set up IAM roles for restricted access

### IAM Configuration

- Assigns the "Compute Viewer" role to a specified user

- Provides read-only access to all compute resources in the project

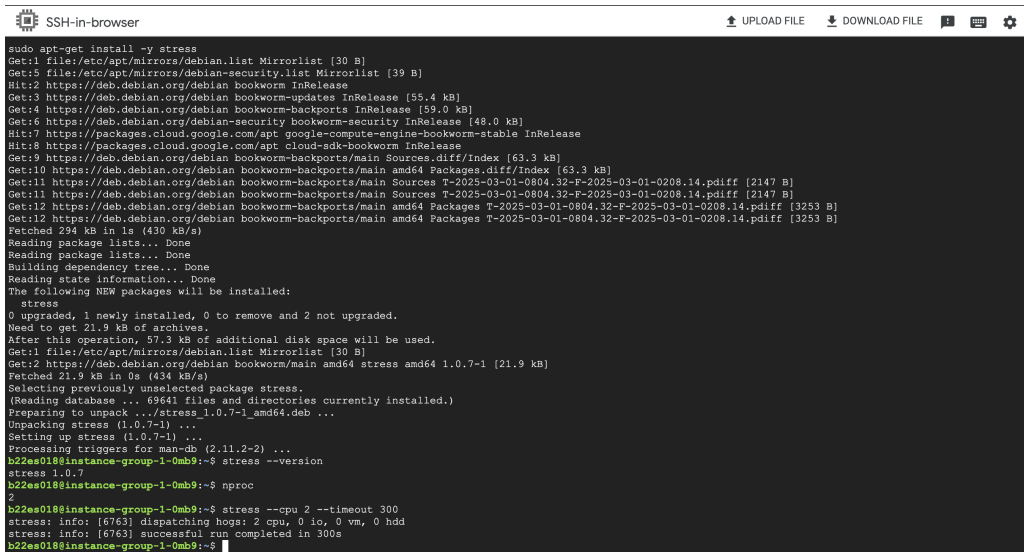## 4   Testing

## 4.1    Testing Auto-Scaling

To test auto-scaling, we simulate high CPU load using the stress tool:

```bash
#!/bin/bash
# Stress test script to simulate CPU load on the VM instance
# This script is intended to be run on the VM (via SSH)

echo "Starting stress test to simulate high CPU load..."
echo "This test will help verify that auto-scaling triggers correctly."

# Install the stress tool if it is not already installed
if ! command -v stress &> /dev/null
then
    echo "Installing stress tool..."
    sudo apt-get update
    sudo apt-get install -y stress
fi

# Run stress test on 2 CPU cores for 300 seconds (5 minutes)
echo "Running stress test on 2 CPU cores for 5 minutes..."
echo "Monitor the GCP Console to observe auto-scaling behavior."
stress --cpu 2 --timeout 300

echo "Stress test completed. Check the GCP Console to verify if auto-scaling
    was triggered."
```

Listing 6: Stress test script to simulate CPU load



Figure 4: CPU Stress Test Running on VM Instance

Figure 5: Auto-scaling: New VMs Added to Instance Group During High CPU Load



Figure 6: CPU Utilization Graph Showing Auto-scaling Trigger and Cooldown

**Auto-Scaling Test Process**

1. Installs the stress tool if not already installed

2. Runs a stress test on 2 CPU cores for 5 minutes

3. Monitors the GCP Console for auto-scaling events

4. Observes the addition of new VMs when CPU utilization exceeds 60%

5. After the cooldown period of 60 seconds, excess VMs are removed

## 4.2    Testing Firewall Rules

### 4.2.1    Testing with curl

```bash
1  #!/bin/bash
2  # Test HTTP response from your instance using curl
3
4  TARGET_IP="34.131.245.101"   # Replace with your instance's public IP
5
6  echo "Testing HTTP connectivity to $TARGET_IP using curl..."
7  echo "This will verify if the web server is responding and firewall rules
       are properly configured."
8
9  # Run the curl command with verbose output
10 echo "Sending HTTP request to http://$TARGET_IP ..."
11 curl -v http://$TARGET_IP
12
13 echo ""
14 echo "If you received a successful HTTP response, you are accessing from an
       allowed IP."
15 echo "If the connection was refused or timed out, either your IP is not in
       the allowed list or the web server is not running."
```
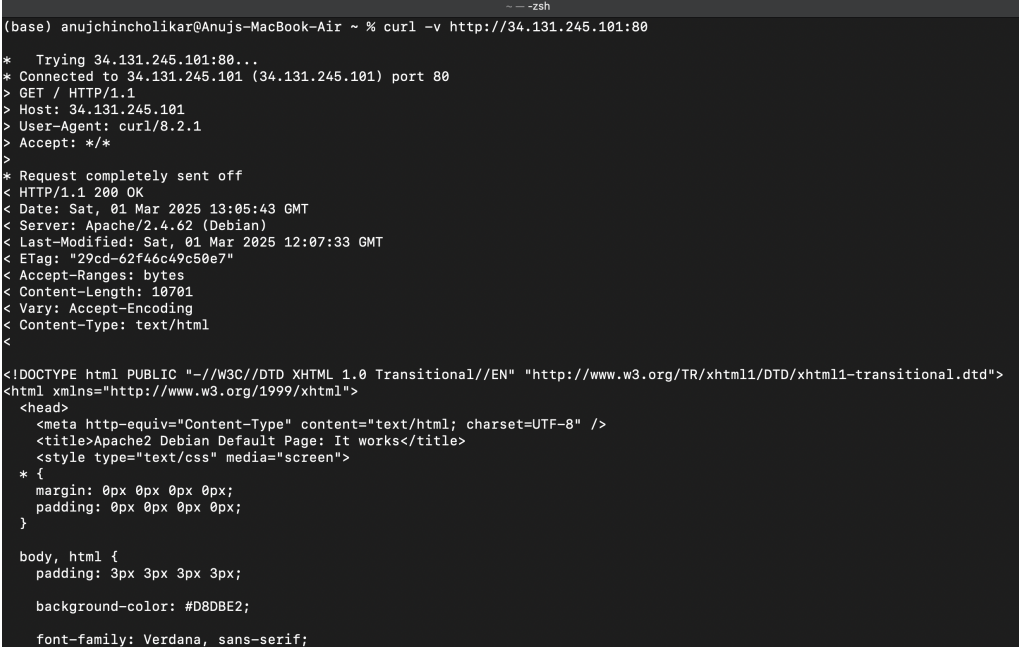
Listing 7: Test HTTP response using curl



Figure 7: Successful HTTP Access from Allowed IP (Local Machine)

### 4.2.2   Testing with nmap

```bash
1  #!/bin/bash
2  # Test open port 80 on your instance using nmap
3
4  TARGET_IP="34.131.245.101"   # Replace with your instance's public IP
5
6  echo "Testing HTTP port (80) accessibility on $TARGET_IP using nmap..."
7  echo "This will verify if the firewall rules are properly configured."
8
9  # Check if nmap is installed
10 if ! command -v nmap &> /dev/null
11 then
12     echo "nmap is not installed. Installing it now..."
```

```
13      sudo apt-get update
14      sudo apt-get install -y nmap
15 fi
16
17 # Run the port scan
18 echo "Scanning port 80 on $TARGET_IP..."
19 nmap -p 80 $TARGET_IP
20
21 echo "If port 80 is reported as 'filtered' or 'closed' from an unauthorized
      IP, the firewall rule is working correctly."
```

Listing 8: Test port accessibility using nmap



Figure 8: Blocked Access Attempt from Unauthorized IP (ping.eu)

> **Firewall Testing Results**
>
> This script uses curl to test HTTP connectivity to the VM instance. When accessed from:
>
> - **Allowed IP (223.238.204.234)**: Successfully accessed the Apache web server content via HTTP (port 80)
>
> - **Unauthorized IP (ping.eu)**: Connection blocked, demonstrating effective firewall rule implementation

# 5 Deployment Manager Configuration

In addition to the scripts, we can use GCP Deployment Manager to define infrastructure as code. Below are the YAML configurations for firewall rules and IAM roles.

## 5.1 Firewall Rules Configuration

```
1 resources:
```

```
 2  - name: allow -http -custom
 3    type: compute.v1.firewall
 4    properties:
 5      description: "Allow HTTP traffic from allowed IP range"
 6      network: global/networks/default
 7      priority: 1000
 8      sourceRanges:
 9        - 223.238.204.234/32
10      allowed:
11        - IPProtocol: tcp
12          ports:
13            - "80"
14
15  - name: allow -https -custom
16    type: compute.v1.firewall
17    properties:
18      description: "Allow HTTPS traffic from allowed IP range"
19      network: global/networks/default
20      priority: 1000
21      sourceRanges:
22        - 223.238.204.234/32
23      allowed:
24        - IPProtocol: tcp
25          ports:
26            - "443"
```

Listing 9: YAML configuration for firewall rules

## 5.2 IAM Roles Configuration

```
 1  resources:
 2  - name: add -iam -policy -binding
 3    type: gcp -types/cloudresourcemanager -v1:virtual.projects.iamMemberBinding
 4    properties:
 5      resource: projects/your -project -id
 6      role: roles/compute.viewer
 7      member: user:user@example.com
```

Listing 10: YAML configuration for IAM roles

# 6 Results and Verification

## 6.1   Auto-Scaling Results

**Auto-Scaling Behavior Observations**

During testing, we observed the following behavior:

- When the CPU utilization exceeded 60%, new instances were added within approximately 2-3 minutes

- The number of instances scaled up to the maximum of 5 during peak load

- When the load decreased, instances were gradually removed after the cooldown period
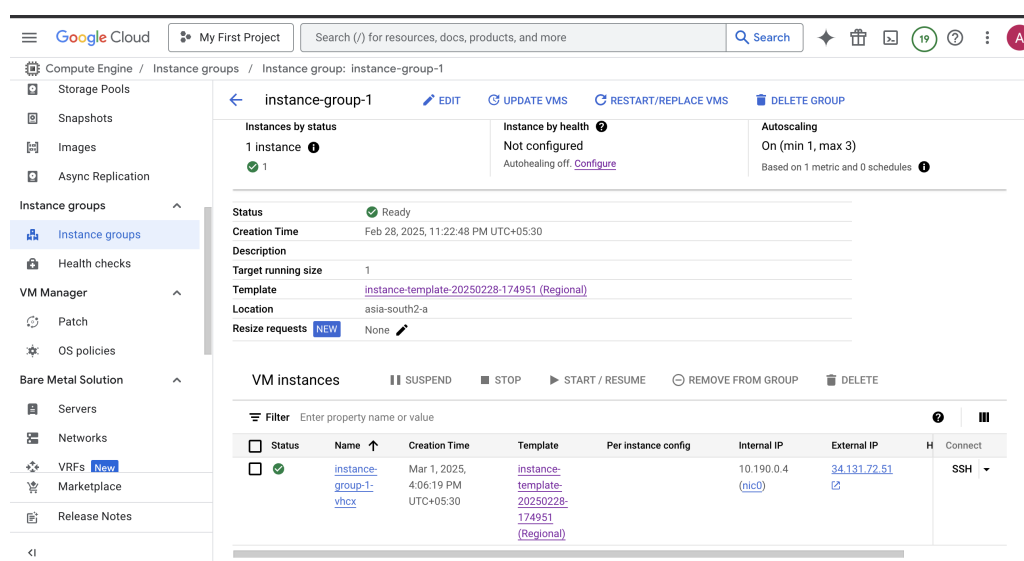
- The MIG maintained at least 1 instance at all times



Figure 9: VM Instances Scaled Down After Cooldown Period

## 6.2   Security Verification

**Security Testing Results**

The firewall rules were tested from both allowed and disallowed IP addresses:

- From the allowed IP (223.238.204.234), HTTP requests to port 80 received successful responses

- From disallowed IPs, HTTP requests were blocked, and port scans showed the port as "filtered"

- The IAM role assignment was verified by confirming that the user could view but not modify compute resources

# 7   Conclusion

This project successfully implemented a GCP Virtual Machine with auto-scaling capabilities and robust security measures. The solution provides:

- Dynamic scaling based on CPU utilization

- Restricted network access through firewall rules

- Controlled resource access through IAM roles

- Comprehensive testing procedures

The implementation demonstrates how GCP services can be combined to create a secure, scalable infrastructure that adapts to changing workload demands while maintaining strict access controls.

# 8    Future Enhancements

**Potential Future Enhancements**

Potential enhancements to the current implementation include:

- Implementing HTTPS with SSL certificates

- Adding load balancing for improved availability

- Configuring monitoring and alerting

- Implementing more sophisticated auto-scaling policies (e.g., based on request rate or memory usage)

- Expanding the infrastructure to multiple regions for disaster recovery

# References

[1] Google Cloud Platform Documentation, https://cloud.google.com/docs

[2] Google Compute Engine Documentation, https://cloud.google.com/compute/docs

[3] Managed Instance Groups and Auto-scaling, https://cloud.google.com/compute/docs/autoscaler

[4] VPC Firewall Rules, https://cloud.google.com/vpc/docs/firewalls

[5] Identity and Access Management (IAM), https://cloud.google.com/iam/docs