

VCC Assignment 3

Report

Gourav Sen

M24CSA011

Github: [Github link](#)

Video: [Video Link](#)

1. Introduction

Objective: **Autoscaling a web application from a local VM to GCP** when resource usage reaches 75%, followed by live migration back to the local machine when load decreases.

- Tools and technologies used:
 - **VirtualBox** for local VM
 - **Lubuntu** as the guest OS
 - **Docker** for containerizing the application
 - **htop** for resource monitoring
 - **GCP** for autoscaling and live migration
 - **Bash scripts** for automation
 - Architecture Diagram
-

2. Prerequisites and Setup

- **System Requirements**
 - Host Machine: OS, RAM, CPU, Disk space
 - Guest VM: Lubuntu configuration (RAM, CPU, Disk)
 - **Downloading Required Software**
 - **VirtualBox**
 - **Lubuntu ISO**
 - **GCP SDK**
 - **Installing VirtualBox**
 - **Creating Lubuntu VM**
 - **VM Configuration**
 - **Network Configuration (NAT with Port Forwarding)**
-

3. VM Configuration and Access

- **Starting the Lubuntu VM**
 - **Port Forwarding Rules**
 - SSH Access: Host Port 2222 → Guest Port 22
 - WebApp Access: Host Port 8080 → Guest Port 80
 - **SSH into the VM from the Host**
 - Verifying connectivity
 - Troubleshooting SSH issues
-

4. Docker Installation and Setup

- **Installing Docker on Lubuntu**

- **Configuring Docker**
 - **Verifying Docker Installation**
 - **Installing Required Packages**
 - Python
 - Pip
 - Docker dependencies
-

5. Web Application Development

- **Creating app.py:**
 - Simple Python Flask-based web application
 - Application will stress the CPU to simulate high load
 - **Creating requirements.txt:**
 - Define required Python packages
 - **Building and Running the Docker Container**
 - **Testing the Web Application**
-

6. Monitoring and Resource Management

- **Installing htop:**
 - Real-time resource usage monitoring
 - **Creating the Dockerfile:**
 - Define environment and dependencies
 - **Resource Monitoring Script:**
 - **Bash script** to monitor CPU usage
 - Trigger autoscaling when CPU > 75%
-

7. Autoscaling to GCP

- **Installing Google Cloud SDK**
 - Download and install GCP SDK
 - Authenticate using the **service account**
 - **Service Account Setup**
 - Create and configure service account
 - Assign necessary IAM roles
 - **GCP VM Instance Creation**
 - Launch VM on GCP
 - Assign static external IP
 - **Load Transfer to GCP**
 - Redirect traffic from local VM to GCP
-

8. Live Migration and Load Balancing

- **Live Migration:**
 - Trigger GCP instance when CPU > 75%
 - Migrate back to local VM when load decreases
 - **Automation Script:**
 - Bash script to automate the entire process
-

9. Results and Observations

- VM Creation and Deletion Logs
-

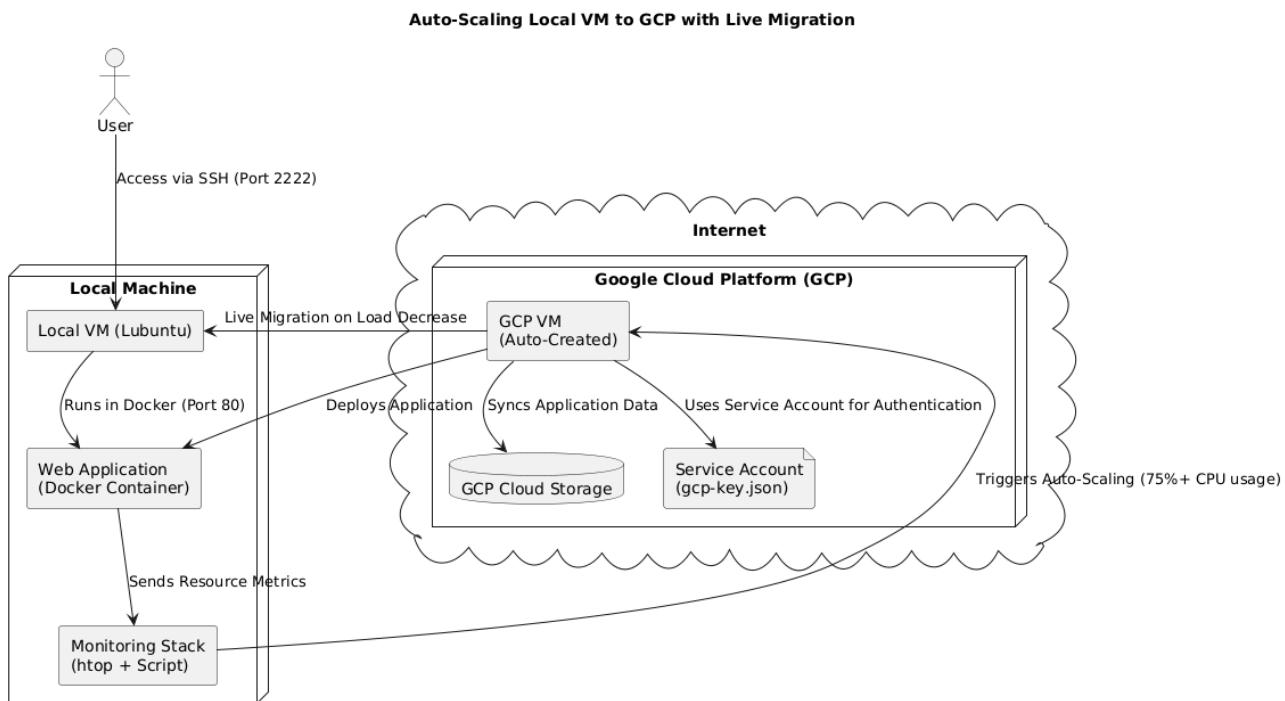
10. Conclusion

- Summary of the project
- Challenges faced
- Possible improvements

• 1. Introduction

The objective of this project is to demonstrate **autoscaling and live migration of a web application** from a local Virtual Machine (VM) to **Google Cloud Platform (GCP)** based on **resource usage**. The application, running on a **Lubuntu VM**, is monitored for CPU utilization. When the CPU usage exceeds **75%**, the application is automatically scaled to a **GCP instance**. When the load decreases, the application is migrated back to the local VM to **optimize resource usage**.

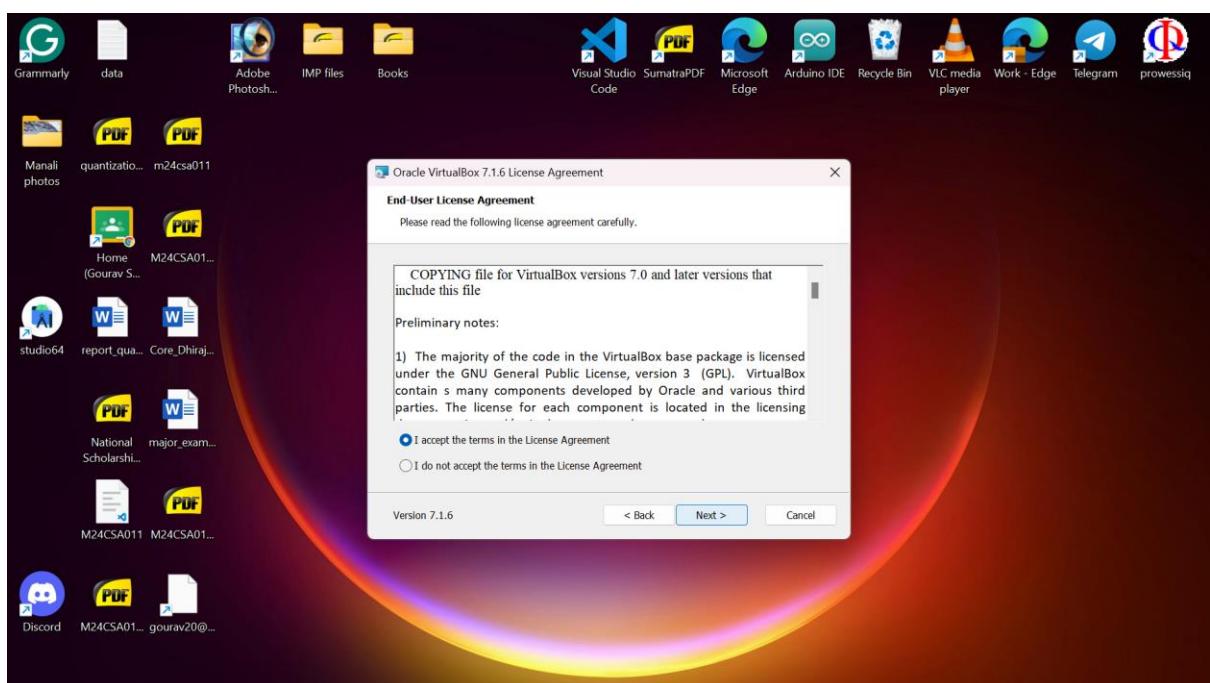
Architecture Diagram

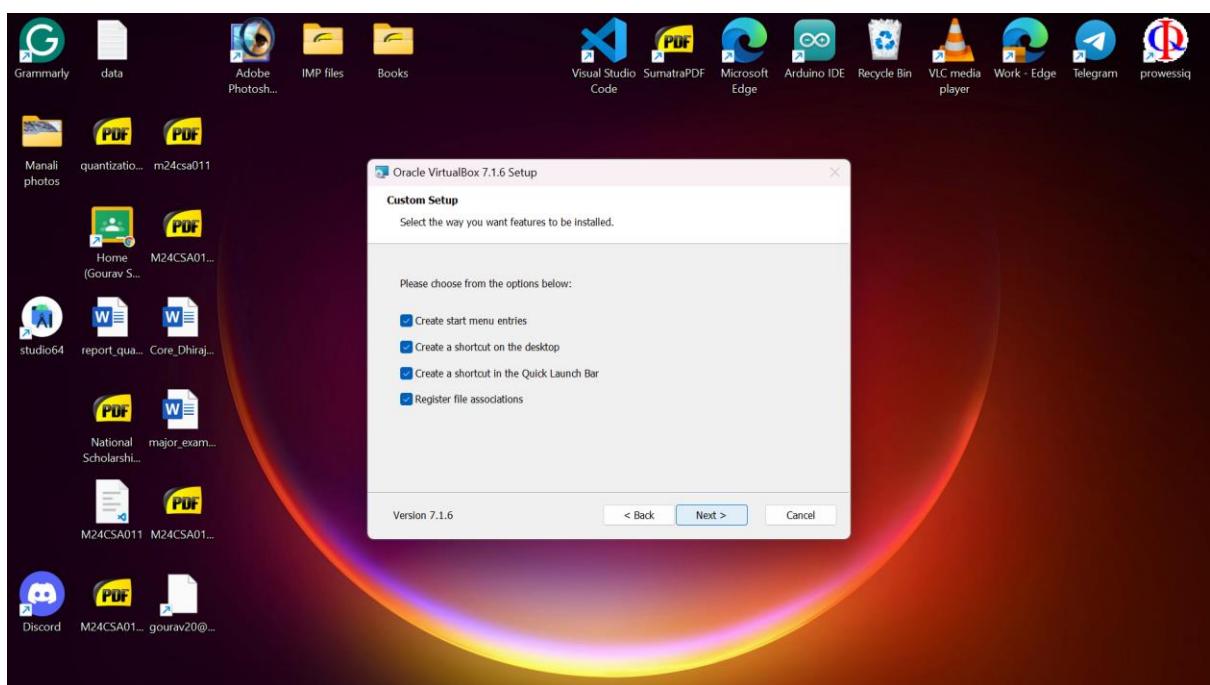
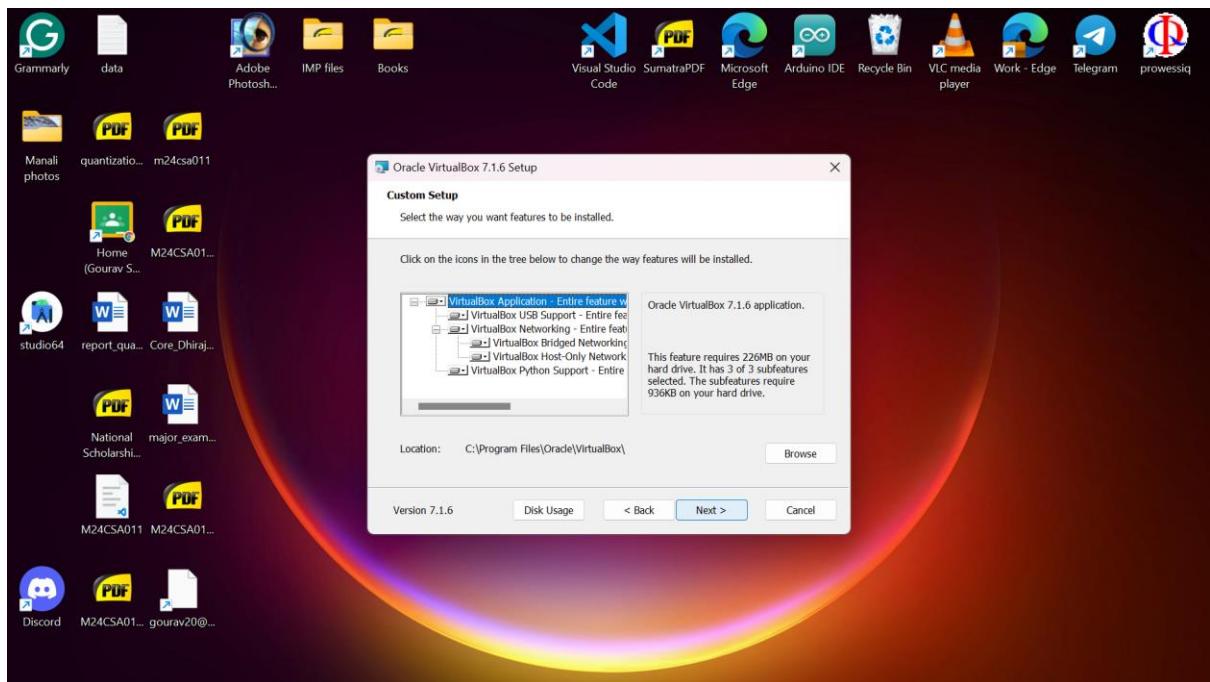


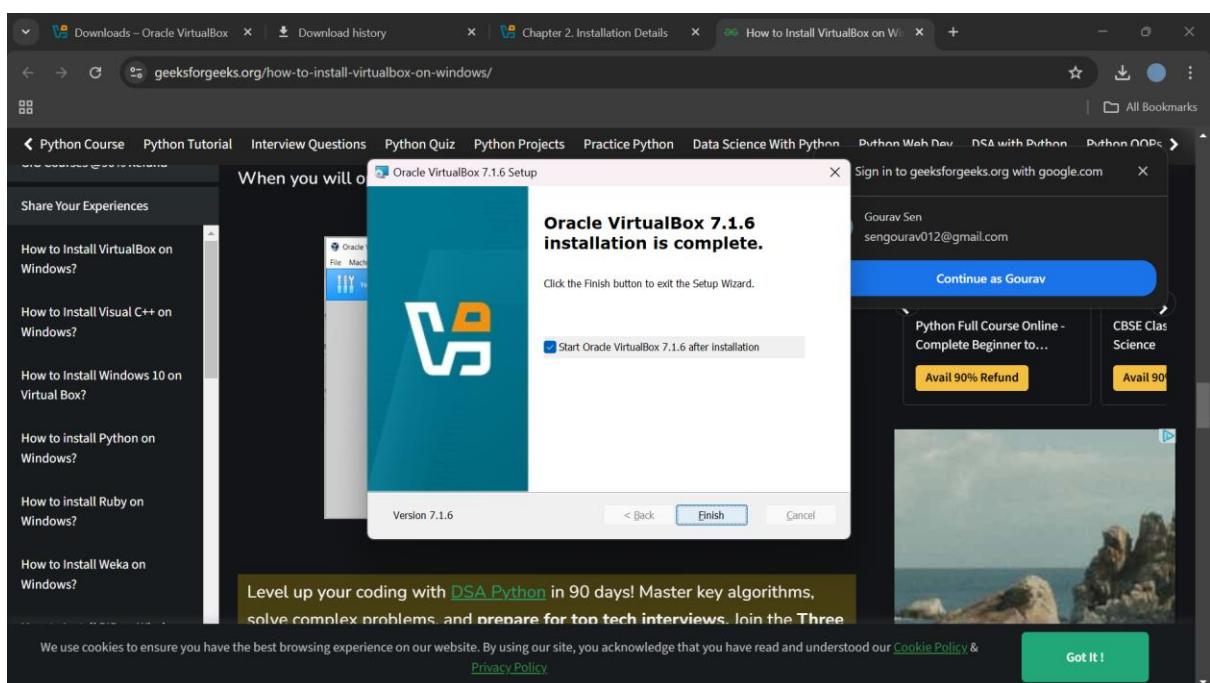
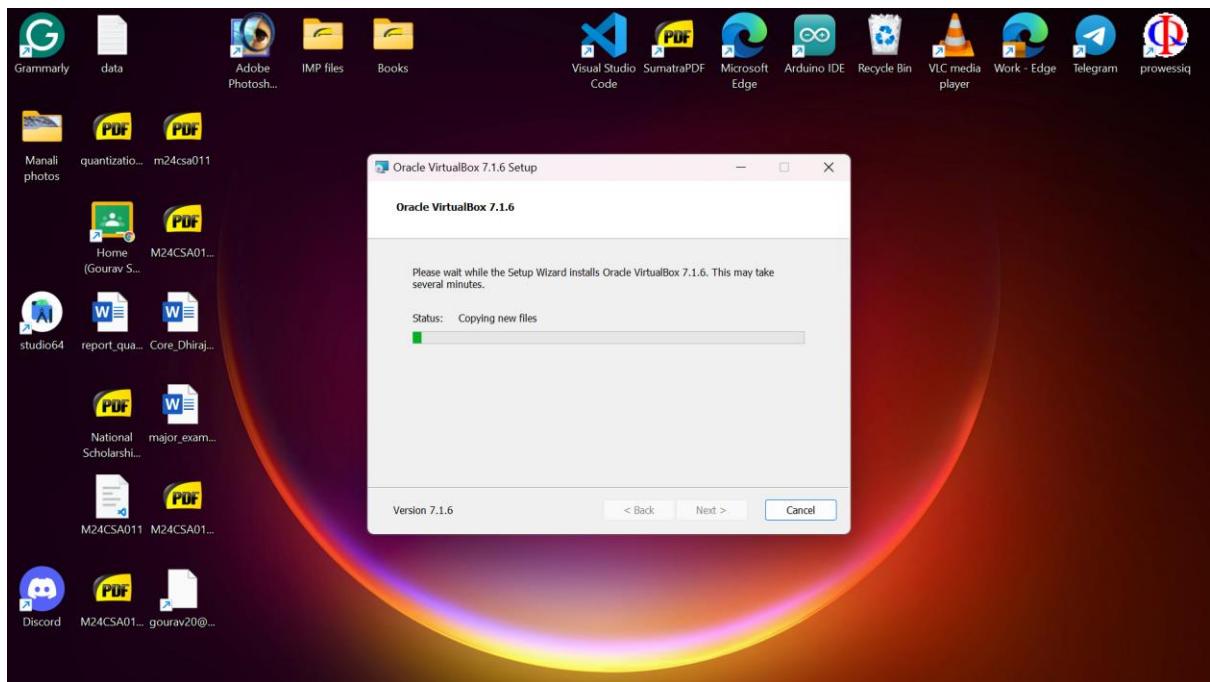
2. Prerequisites and Setup

2.1 Install VirtualBox

Download and install VirtualBox from the official website:
[VirtualBox Download](#)







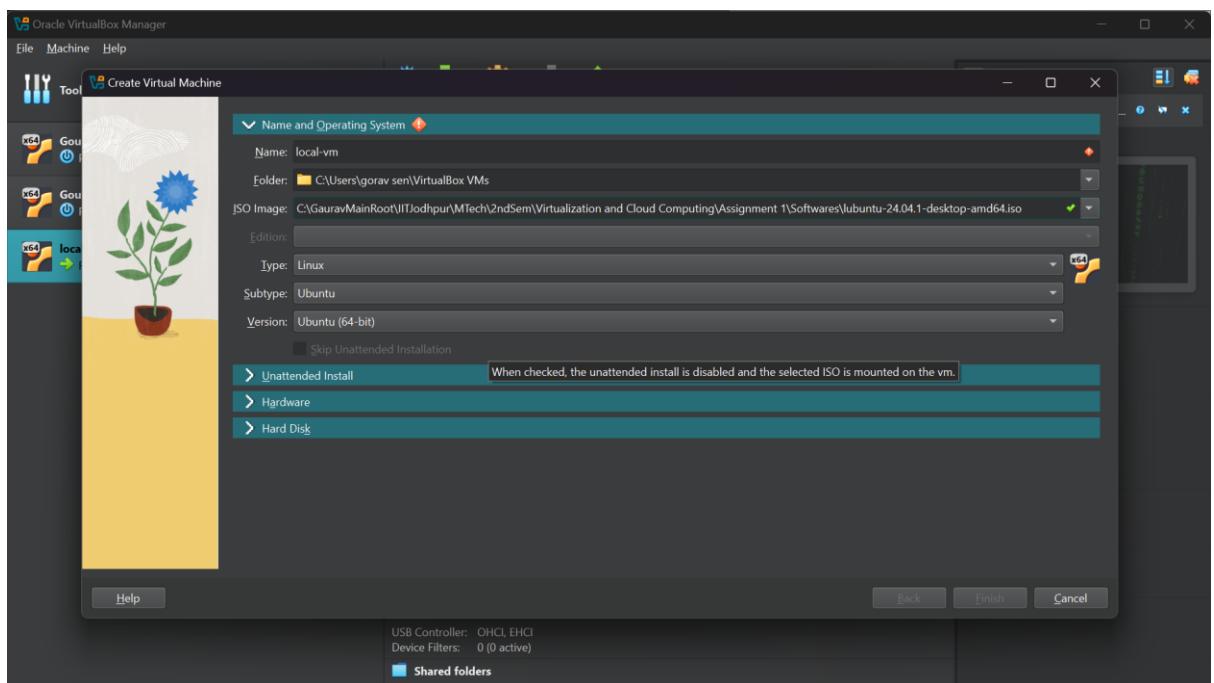
2.3 Download Lubuntu ISO

Since Ubuntu required more resources than available on our machine, we opted for Lubuntu, which is lightweight and suitable for systems with limited resources

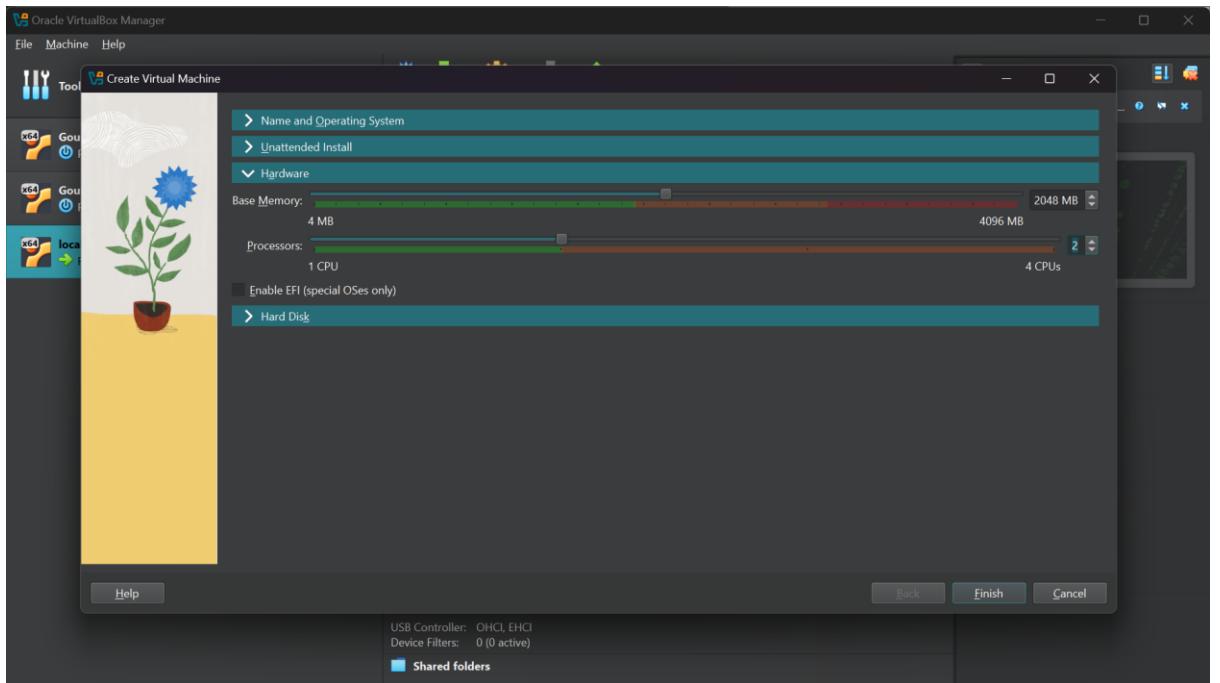
Download the **Lubuntu 22.04 ISO** from the official website:
[Lubuntu Download](https://lubuntu.download)

2.4 Create a New VM in VirtualBox

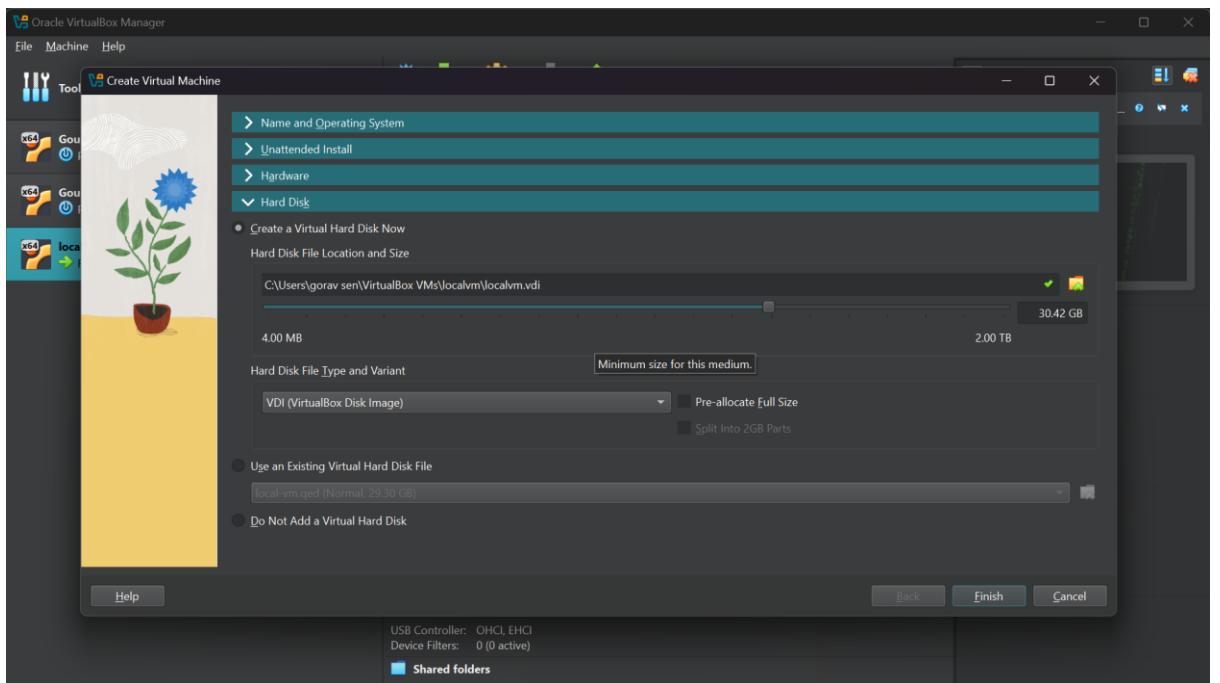
1. Open VirtualBox and click on New.
2. Set the following configurations:
 - o **Name:** Lubuntu-VM
 - o **Type:** Linux
 - o **Version:** Ubuntu (64-bit)
 - o **Memory Size:** 2 GB (2048 MB)
 - o **Hard Disk:** Create a virtual hard disk now
 - o **Disk Size:** 30 GB
3. Click **Create**.



Memory and processors:



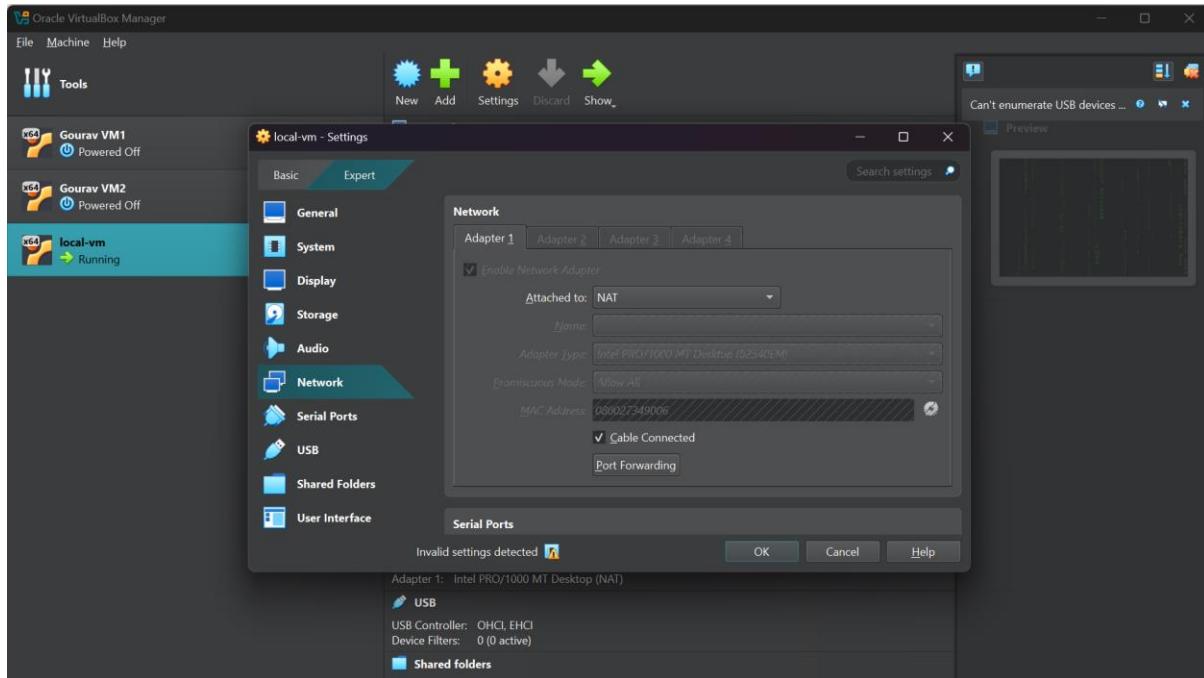
Allocate hard disk and Finish:



2.5 Configure VM Settings

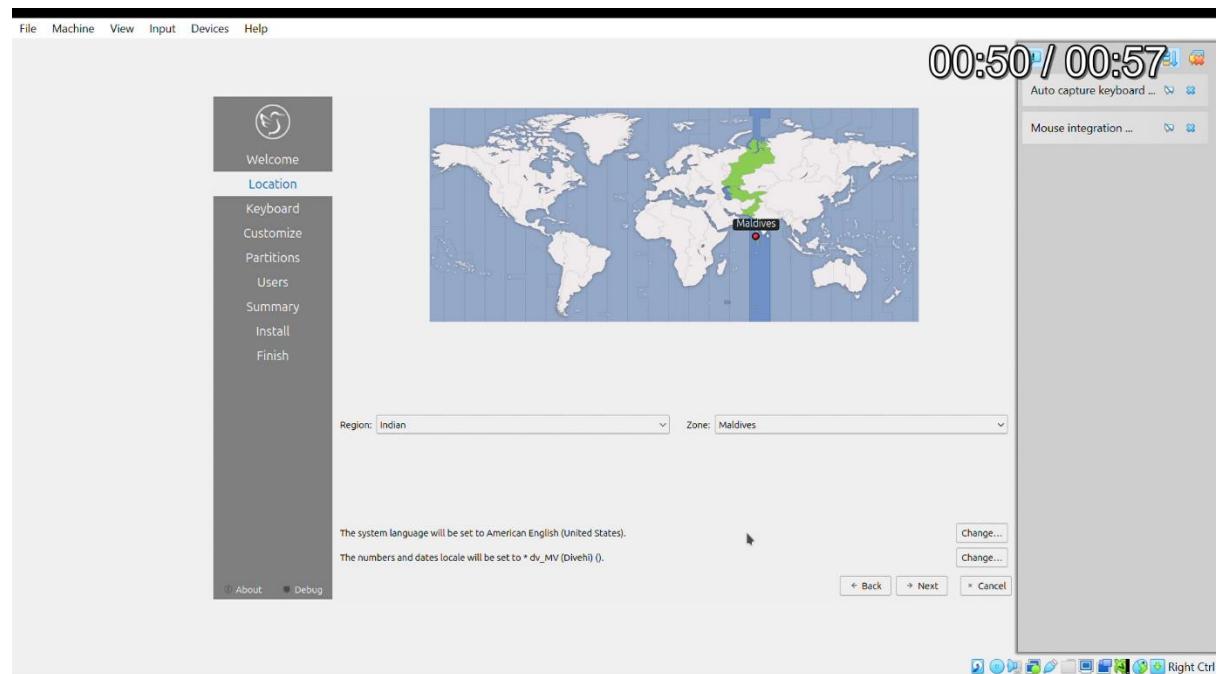
1. Select the newly created VM and click **Settings**.

2. Go to **Network > Adapter 1**:
 - o **Attached to:** NAT
 - o Click **Advanced > Port Forwarding**.
3. Add the following rules:
 - **SSH Rule:**
 - o **Name:** SSH
 - o **Protocol:** TCP
 - o **Host IP:** 127.0.0.1
 - o **Host Port:** 2222
 - o **Guest Port:** 22
 - **Web Application Rule:**
 - o **Name:** WebApp
 - o **Protocol:** TCP
 - o **Host IP:** 127.0.0.1
 - o **Host Port:** 8080
 - o **Guest Port:** 80

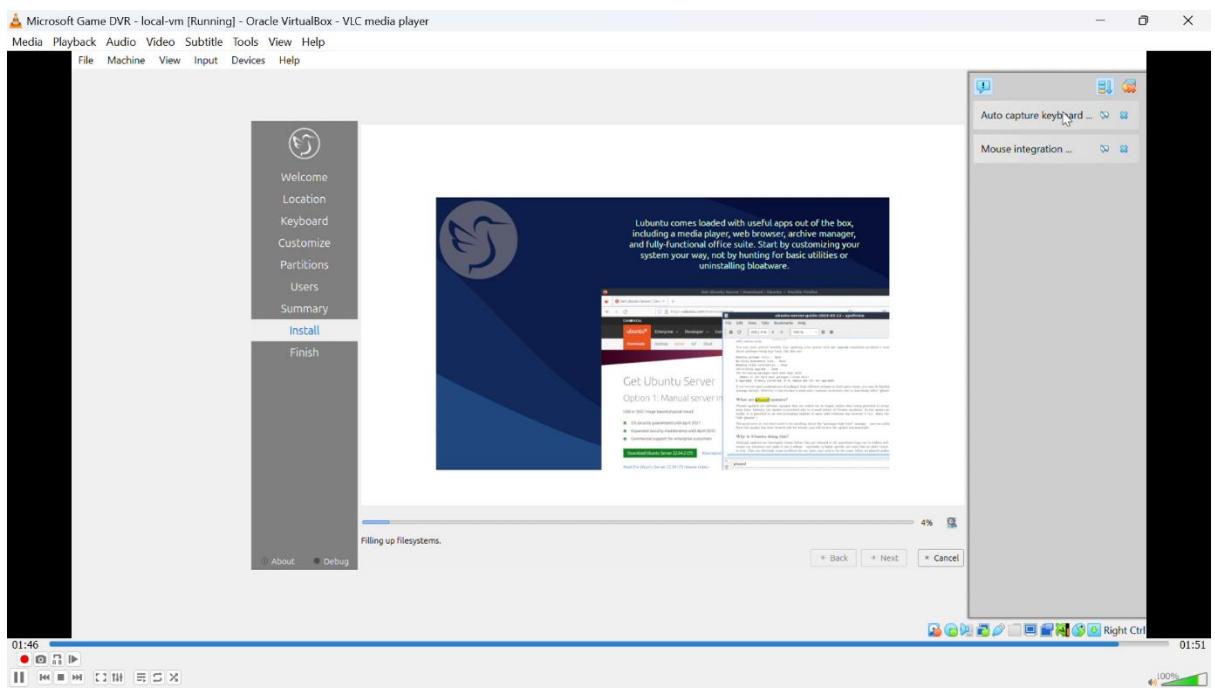
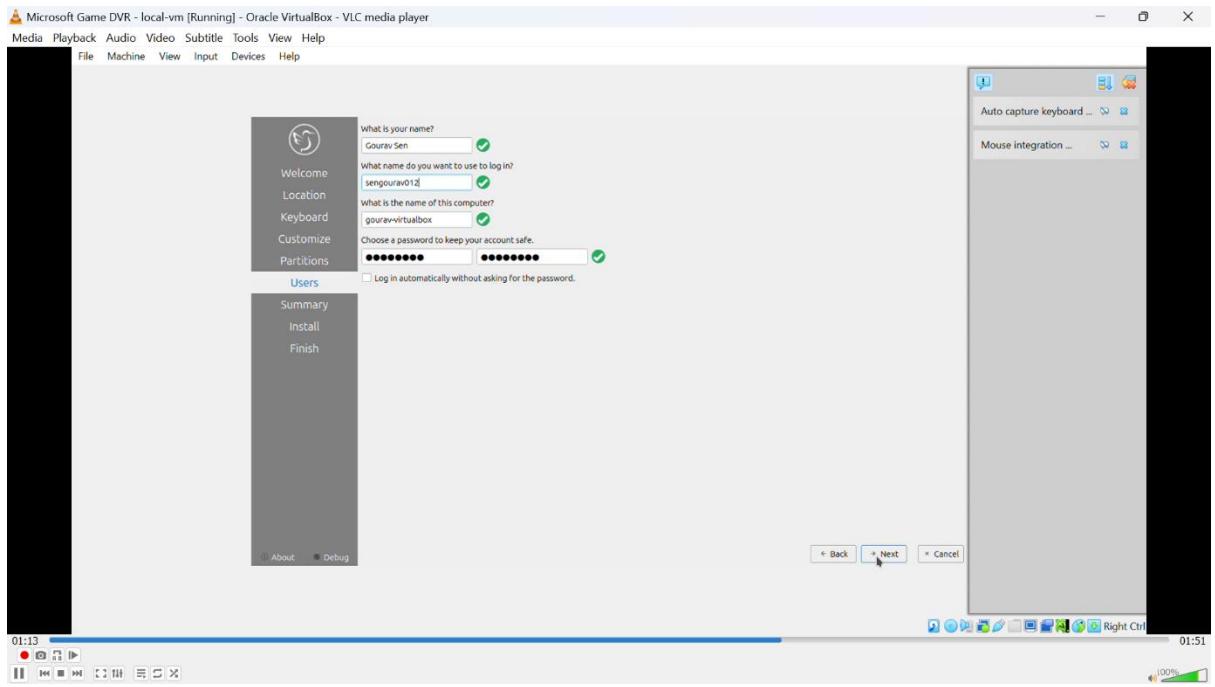


2.6 Install Lubuntu on the VM

1. Start the VM and select the **Lubuntu ISO** to boot.
2. Choose **Install Lubuntu**.
3. Follow the on-screen instructions:
 - o Select **Language, Timezone, and Keyboard Layout**.
 - o Set username and password.
4. Once the installation is complete, **remove the ISO** and restart the VM.



5.



2.7 SSH into the VM

After the VM is up and running, SSH into the VM from your local machine.

First, find the local VM IP by running:

```
ip a
```

From the local machine, SSH into the VM:

```
ssh username@127.0.0.1 -p 2222
```

3. Docker Installation and Setup

3.1 Install Docker on the VM

Once you have SSH access to the VM, install Docker using the following steps.

3.1.1 Update the Package Index

```
sudo apt update
```

3.1.2 Install Required Dependencies

```
sudo apt install -y apt-transport-https ca-certificates curl software-properties-common
```

3.1.3 Add Docker's Official GPG Key

```
curl -fsSL https://download.docker.com/linux/ubuntu/gpg | sudo gpg --dearmor -o /usr/share/keyrings/docker-archive-keyring.gpg
```

3.1.4 Add Docker Repository

```
echo "deb [arch=$(dpkg --print-architecture) signed-by=/usr/share/keyrings/docker-archive-keyring.gpg] \n https://download.docker.com/linux/ubuntu $(lsb_release -cs) stable" | sudo tee /etc/apt/sources.list.d/docker.list > /dev/null
```

3.1.5 Install Docker Engine

```
sudo apt update
```

```
sudo apt install -y docker-ce docker-ce-cli containerd.io
```

3.3 Add User to Docker Group

To run Docker without sudo, add your user to the Docker group:

```
sudo usermod -aG docker $USER
```

```
newgrp docker
```

4. Package Installation and Application Setup

containerize the web application

create app.py

```
from flask import Flask, jsonify
import threading
import time
import random
|
app = Flask(__name__)

# Global variable to control the load
is_loading = False

def cpu_load():
    """Function to simulate continuous CPU load while web app is open."""
    global is_loading
    while is_loading:
        # Perform heavy computation
        [x**2 for x in range(10**7)] # CPU intensive task
        time.sleep(0.1)

@app.route('/')
def home():
    return "Web App is running! Visit /load to increase CPU usage."

@app.route('/load')
def load():
    global is_loading
    if not is_loading:
        is_loading = True
        threading.Thread(target=cpu_load).start()
    return jsonify({"message": "CPU load started"})

@app.route('/stop')
def stop():
    global is_loading
    is_loading = False
    return jsonify({"message": "CPU load stopped"})

if __name__ == '__main__':
    app.run(host='0.0.0.0', port=80)
```

create Dockerfile with following content:

```
# Use lightweight Python image
FROM python:3.9-slim

# Set the working directory
WORKDIR /app

# Copy the app files
COPY app.py /app

# Install dependencies
RUN pip install Flask

# Expose port 80
EXPOSE 80

# Run the app
CMD ["python", "app.py"]
```

Build the image:

```
docker build webapp
```

Create the Docker network

```
docker network create my-network
```

Run the Docker Container with Custom Names

Use the following command to run the web application with:

- **Container Name:** container
- **Network:** my-network
- **Port Mapping:** 80:80

```
docker run -d --name container --network my-network -p 80:80 webapp
```

4.4 Test the Application

Access the application from the host machine:

```
curl http://127.0.0.1:8080
```



Step 8: Resource Monitoring Script

Objective:

- Create a script to monitor the CPU usage of the local VM.
- Trigger autoscaling by creating a GCP VM instance when the CPU usage exceeds **75%**.
- Perform live migration of the application to the GCP VM.
- Migrate the application back to the local VM when the CPU usage decreases below **40%**

Create the Resource Monitoring Script:

`nano autoscale.sh`

```
#!/bin/bash

# Variables
THRESHOLD=75          # Upper threshold to start the VM
MIN_THRESHOLD=30        # Lower threshold to stop the VM
INSTANCE_NAME="web-app-instance"
ZONE="asia-south1-a"
MACHINE_TYPE="e2-medium"
IMAGE_FAMILY="ubuntu-2204-lts"
IMAGE_PROJECT="ubuntu-os-cloud"

# Docker configuration
CONTAINER_NAME="container"
LOCAL_PORT=8080
DOCKER_PORT=80
DOCKER_IMAGE="webapp"

# --- Function to check if GCP VM is reachable ---
check_vm_ready() {
    echo "Checking VM status..."
    local retries=15  # Number of retries (30 seconds total)
    local delay=2      # Delay between retries

    for ((i=1; i<retries; i++)); do
        if gcloud compute ssh "$INSTANCE_NAME" --zone="$ZONE" --strict-host-key-checking=no --command "echo 'VM is ready'" &>/dev/null; then
            echo "VM is reachable!"
            return 0
        else
            echo "Waiting for VM to be reachable... Attempt $i/$retries"
            sleep "$delay"
        fi
    done

    echo "✗ Failed to connect to VM after $((retries * delay)) seconds."
    exit 1
}
```

```

# --- Check if the GCP VM exists ---
echo "Checking if GCP VM exists..."
if ! gcloud compute instances describe "$INSTANCE_NAME" --zone="$ZONE" &> /dev/null; then
    echo "GCP VM does not exist. Creating it..."

    # Create the GCP VM
    gcloud compute instances create "$INSTANCE_NAME" \
        --zone="$ZONE" \
        --machine-type="$MACHINE_TYPE" \
        --image-family="$IMAGE_FAMILY" \
        --image-project="$IMAGE_PROJECT" \
        --tags=http-server

    echo "Waiting for VM to start..."
    sleep 30

    # Ensure the VM is reachable before proceeding
    check_vm_ready

    # SSH into the VM and install Docker, pull image, and run container
    echo "Installing Docker and deploying the web app on the new VM..."

    gcloud compute ssh "$INSTANCE_NAME" --zone="$ZONE" --strict-host-key-checking=no --command "
        sudo apt update -y &&
        sudo apt install -y docker.io &&
        sudo systemctl start docker &&
        sudo systemctl enable docker &&
        sudo docker run -d --name $CONTAINER_NAME -p $DOCKER_PORT:$DOCKER_PORT $DOCKER_IMAGE
    "

    echo "✅ Web app deployed successfully on the new GCP VM!"
else
    echo "GCP VM already exists."
fi

```

```

# --- CPU Monitoring Section ---

# Get Docker container CPU usage
CONTAINER_CPU=$(docker stats --no-stream --format "{{.CPUPerc}}" "$CONTAINER_NAME" | sed 's/%//')

# Get VM CPU usage
VM_CPU=$(top -bn1 | grep "Cpu(s)" | awk '{print $2 + $4}')

# Combine CPU usage by adding container and VM CPU usage
TOTAL_CPU=$(echo "$CONTAINER_CPU + $VM_CPU" | bc -l)
echo "Container CPU: $CONTAINER_CPU%"
echo "VM CPU: $VM_CPU%"
echo "Total CPU Usage: $TOTAL_CPU%"

# --- Autoscaling Logic ---
if (( $(echo "$TOTAL_CPU > $THRESHOLD" | bc -l) )); then
    echo "⚠️ High CPU usage detected: $TOTAL_CPU%"

    if gcloud compute instances describe "$INSTANCE_NAME" --zone="$ZONE" &> /dev/null; then
        echo "Starting existing GCP VM: $INSTANCE_NAME"
        gcloud compute instances start "$INSTANCE_NAME" --zone="$ZONE"

        # Wait until the VM is ready
        check_vm_ready

    else
        echo "Creating and starting a new GCP VM: $INSTANCE_NAME"
        gcloud compute instances create "$INSTANCE_NAME" \
            --zone="$ZONE" \
            --machine-type="$MACHINE_TYPE" \
            --image-family="$IMAGE_FAMILY" \
            --image-project="$IMAGE_PROJECT" \
            --tags=http-server

        echo "Installing Docker and deploying web app..."
    fi
fi

```

```

echo "Installing Docker and deploying web app..."

gcloud compute ssh "$INSTANCE_NAME" --zone="$ZONE" --strict-host-key-checking=no --command "
    sudo apt update -y &&
    sudo apt install -y docker.io &&
    sudo systemctl start docker &&
    sudo systemctl enable docker &&
    sudo docker run -d --name $CONTAINER_NAME -p $DOCKER_PORT:$DOCKER_PORT $DOCKER_IMAGE
"

echo "✓ Web app deployed successfully!"
fi

elif (( $(echo "$TOTAL_CPU < $MIN_THRESHOLD" | bc -l) )); then
    echo "✓ Low CPU usage detected: $TOTAL_CPU%"
    echo "Stopping GCP VM: $INSTANCE_NAME"
    gcloud compute instances stop "$INSTANCE_NAME" --zone="$ZONE"

else
    echo "✓ CPU usage is within normal range: $TOTAL_CPU%"
fi

```

Make the Script Executable:

```
chmod +x autoscale.sh
```

Step 9: Google Cloud SDK Installation and Service Account Setup

Objective:

- Install the Google Cloud SDK to interact with GCP.
- Authenticate with GCP using a service account.
- Configure the SDK with the Asia-South1 (Mumbai) region.
- Verify GCP connectivity.

```
# Add the Cloud SDK distribution URI as a package source
```

```
echo "deb [signed-by=/usr/share/keyrings/cloud.google.gpg] http://packages.cloud.google.com/apt cloud-sdk
main" | sudo tee -a /etc/apt/sources.list.d/google-cloud-sdk.list
```

```
# Import the Google Cloud public key
```

```
curl https://packages.cloud.google.com/apt/doc/apt-key.gpg | sudo gpg --dearmor -o
/usr/share/keyrings/cloud.google.gpg
```

```
# Update and install Google Cloud SDK
```

```
sudo apt update && sudo apt install -y google-cloud-sdk
```

Verify the installation by running:

```
gcloud version
```

If this appears you are good to go.

```
Google Cloud SDK 465.0.0
bq 2.0.88
core 2024.03.01
gsutil 5.30
```

Create Service Account :

Go to IAM and Admin:

The screenshot shows the Google Cloud IAM & Admin interface. The left sidebar has 'IAM & Admin / IAM' selected. The main area is titled 'Search Results' and shows a list of items related to IAM and Admin. The first item is 'IAM' under 'IAM & Admin'. Other items include 'Service Accounts', 'Use IAM securely | IAM Documentation', 'PAM', 'Roles', and 'Gemini for Google Cloud'. A search bar at the top right contains the text 'iam & admin'. At the bottom, there's a table with columns 'Role', 'Security insights', and 'Actions'.

Go to Service account:

Service accounts for project "autoscale-project"

A service account represents a Google Cloud service identity, such as code running on Compute Engine VMs, App Engine apps, or systems running outside Google. [Learn more about service accounts](#).

Organization policies can be used to secure service accounts and block risky service account features, such as automatic IAM Grants, key creation/upload, or the creation of service accounts entirely. [Learn more about service account organization policies](#).

Email	Status	Name	Description	Key ID	Actions
1008372883176-compute@developer.gserviceaccount.com	Enabled	Compute Engine default service account		2372a0a3fd15386310b6e106e3	

Create service account here.

After creating service account we need to create a key as well, this will be used by our VM to authenticate with the GCP –

Service accounts for project "autoscale-project"

A service account represents a Google Cloud service identity, such as code running on Compute Engine VMs, App Engine apps, or systems running outside Google. [Learn more about service accounts](#).

Organization policies can be used to secure service accounts and block risky service account features, such as automatic IAM Grants, key creation/upload, or the creation of service accounts entirely. [Learn more about service account organization policies](#).

Email	Status	Name	Description	Key ID	Actions
1008372883176-compute@developer.gserviceaccount.com	Enabled	Compute Engine default service account		2372a0a3fd15386310b6e106e3	
new-service@autoscale-project-454519.iam.gserviceaccount.com	Enabled	new-service	No keys		

Service account created

The screenshot shows the Google Cloud IAM & Admin interface for creating a new service account. The left sidebar is titled 'Google Cloud' and has a 'Service Accounts' section selected. The main page title is 'new-service'. The 'Keys' tab is active. A prominent message at the top right says: 'Service account keys could pose a security risk if compromised. We recommend you avoid downloading service account keys and instead use the Workload Identity Federation. Learn more about the best way to authenticate service accounts on Google Cloud.' Below this, another message states: 'Google automatically disables service account keys detected in public repositories. You can customize this behavior by using the 'iam.serviceAccountKeyExposureResponse' organization policy. Learn more.' A button labeled 'Add key' is visible, with a dropdown menu showing 'Create new key' and 'Upload existing key'. A success message 'Service account created' is displayed in a black box at the bottom.

This screenshot shows a modal dialog box titled 'Create private key for "new-service"'. It contains instructions: 'Downloads a file that contains the private key. Store the file securely because this key can't be recovered if lost.' It offers two options for 'Key type': 'JSON' (selected) and 'P12'. A note below says: 'For backward compatibility with code using the P12 format'. At the bottom of the dialog are 'Cancel' and 'Create' buttons. The background of the main interface shows the same 'new-service' service account page as the previous screenshot, with the 'Keys' tab selected and a 'Service account created' message at the bottom.

This will download the json file with authentication key.

Authenticate with the service account using the key:

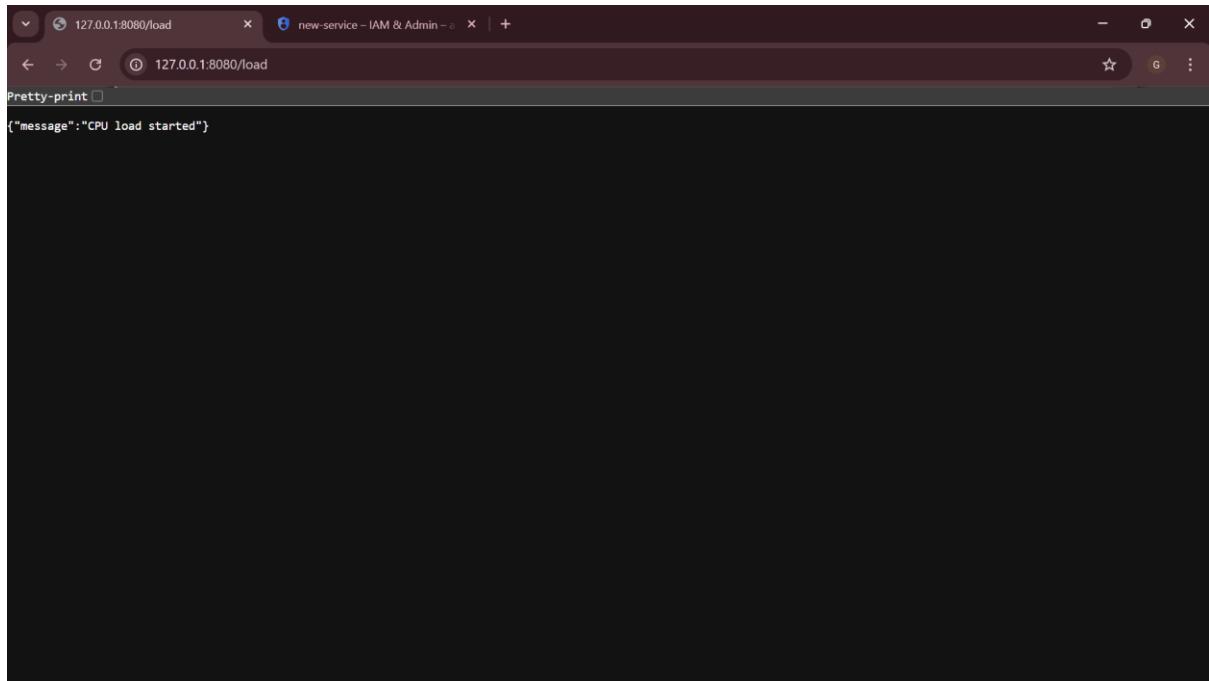
```
gcloud auth activate-service-account --key-file ~/gcp-key.json
```

Final Results and Output:

Now after doing all the setups we can check the autoscaling part:

Lets give load to site:

Run <http://127.0.0.1:8080/load> in the browser:



Run the monitoring script:

```
sengourav@12@gourav-virtualbox:~/autoscale-app$ bash autoscale4.sh
```

Result:

```

sengourav012@gourav-virtualbox:~/autoscale-app$ bash autoscale4.sh
Checking if GCP VM exists...
GCP VM does not exist. Creating it...
Created [https://www.googleapis.com/compute/v1/projects/autoscale-project-454519/zones/asia-south1-a/instances/web-app-instance].
WARNING: Some requests generated warnings:
- You are creating a global DNS VM. VM instances using global DNS are vulnerable to cross-regional outages. To reduce the risk of widespread service disruption, use zonal DNS instead. Learn more at https://cloud.google.com/compute/docs/networking/zonal-dns

NAME          ZONE      MACHINE_TYPE  PREEMPTIBLE INTERNAL_IP  EXTERNAL_IP   STATUS
web-app-instance  asia-south1-a  e2-medium        10.160.0.4  34.93.218.174  RUNNING

Waiting for VM to start...
Checking VM status...
VM is reachable!
Installing Docker and deploying the web app on the new VM...

WARNING: apt does not have a stable CLI interface. Use with caution in scripts.

Hit:1 http://asia-south1.gce.archive.ubuntu.com/ubuntu jammy InRelease
Get:2 http://asia-south1.gce.archive.ubuntu.com/ubuntu jammy-updates InRelease [128 kB]
Get:3 http://security.ubuntu.com/ubuntu jammy-security InRelease [129 kB]
Get:4 http://asia-south1.gce.archive.ubuntu.com/ubuntu jammy-backports InRelease [127 kB]
Get:5 http://asia-south1.gce.archive.ubuntu.com/ubuntu jammy/universe amd64 Packages [14.1 MB]
Get:6 http://security.ubuntu.com/ubuntu jammy-security/main amd64 Packages [2145 kB]
Get:7 http://asia-south1.gce.archive.ubuntu.com/ubuntu jammy/universe Translation-en [5652 kB]
Get:8 http://asia-south1.gce.archive.ubuntu.com/ubuntu jammy/universe amd64 c-n-f Metadata [286 kB]
Get:9 http://asia-south1.gce.archive.ubuntu.com/ubuntu jammy/multiverse amd64 Packages [217 kB]

```

We can see that in the case of load (above 75%) , a new vm instance is being created in google cloud platform.

Run monitoring script:

```

sengourav012@gourav-virtualbox:~/autoscale-app$ nano autoscale4.sh
sengourav012@gourav-virtualbox:~/autoscale-app$ bash autoscale4.sh
Checking if GCP VM exists...
GCP VM already exists.
Container CPU: 194.21%
VM CPU: 73.4%
Total CPU Usage: 267.61%
⚠️ High CPU usage detected: 267.61%
Starting existing GCP VM: web-app-instance
Starting instance(s) web-app-instance...done.
Updated [https://compute.googleapis.com/compute/v1/projects/autoscale-project-454519/zones/asia-south1-a/instances/web-app-instance].
Instance internal IP is 10.160.0.5
Instance external IP is 35.244.38.92
Checking VM status...
Waiting for VM to be reachable... Attempt 1/15
VM is reachable!
sengourav012@gourav-virtualbox:~/autoscale-app$ |

```

We can see that the load is being transferred to GCP VM instance.

Lets decrease the load.

Run <http://127.0.0.1:8080/stop> in browser:

Output:

```
Command Prompt: sengourav012@gourav-virtualbox: ~/autoscale-app
    ~/autoscale-app          MV"
    ctr+alt+t                dv_MV",
    LC_NUMERIC = "en_IN",
    LC_PAPER = "dv_MV",
    LANG = "C.UTF-8"
    are supported and installed on your system.
perl: warning: Falling back to a fallback locale ("C.UTF-8").
VirtualBox: Running kernel seems to be up-to-date.
VirtualBox: No services need to be restarted.
VirtualBox: No containers need to be restarted.
VirtualBox: No user sessions are running outdated binaries.
VirtualBox: No VM guests are running outdated hypervisor (qemu) binaries on this host.
VirtualBox: Unable to find image 'webapp:latest' locally
VirtualBox: docker: Error response from daemon: pull access denied for webapp, repository does not exist or may require 'docker login': denied: requested access to the resource is denied.
VirtualBox: See 'docker run --help'.
VirtualBox:  Web app deployed successfully on the new GCP VM!
VirtualBox: Container CPU: 0.04%
VirtualBox: VM CPU: 5.6%
VirtualBox: Total CPU Usage: 5.64%
VirtualBox:  Low CPU usage detected: 5.64%
VirtualBox: Stopping GCP VM: web-app-instance
VirtualBox: Stopping instance(s) web-app-instance...done.
VirtualBox: Updated [https://compute.googleapis.com/compute/v1/projects/autoscale-project-454519/zones/asia-south1-a/instances/web-app-instance].
VirtualBox: https://compute.googleapis.com/compute/v1/projects/autoscale-project-454519/zones/asia-south1-a/instances/web-app-instance
```

Clearly the the VM instance is stopped, when the load is below 40% (min threshold).

Conclusion:

The project successfully demonstrated **auto-scaling and live migration** of a web application from a **local VM to GCP** based on CPU usage thresholds. The **monitoring script** effectively triggered the scaling event when CPU usage exceeded **75%**, deploying the application to a GCP VM. Upon load reduction, the application automatically migrated back to the local VM. This solution ensures **efficient resource utilization**, seamless scaling, and reduces infrastructure costs by dynamically managing cloud resources.