# Placement Empowerment Program

## *Cloud Computing and DevOps Centre*

Implement Auto-scaling in the Cloud

Set up an auto-scaling group for your cloud VMs to handle variable workloads

Name : Harshana Perianayaki B          Department : IT

# *INTRODUCTION*

In today's cloud-driven environment, businesses require scalable and cost-effective solutions to manage fluctuating workloads efficiently. Auto-scaling is a crucial feature in cloud computing that enables applications to dynamically adjust resources based on demand. Cloud platforms, such as Microsoft Azure, AWS, and Google Cloud, provide robust auto-scaling capabilities to ensure optimal performance and cost efficiency.

# *OVERVIEW*

Auto-scaling allows organizations to automatically increase or decrease the number of virtual machines (VMs) within an auto-scaling group based on predefined rules and real-time performance metrics. This ensures that applications can handle varying traffic loads while maintaining high availability and cost-effectiveness. Cloud-based auto-scaling solutions integrate with monitoring tools and load balancers to provide seamless scaling experiences.

# *OBJECTIVES*

- **Efficiently Manage Workloads**: Automatically adjust resources based on real-time demand.
- **Improve Application Performance**: Maintain optimal response times by scaling resources up or down.

- **Reduce Operational Costs**: Avoid over-provisioning and only use necessary resources.
- **Enhance System Reliability**: Minimize downtime by ensuring the availability of sufficient resources during high traffic periods.
- **Automate Scaling Policies**: Utilize predefined metrics such as CPU utilization, memory usage, or custom application triggers for auto-scaling decisions.
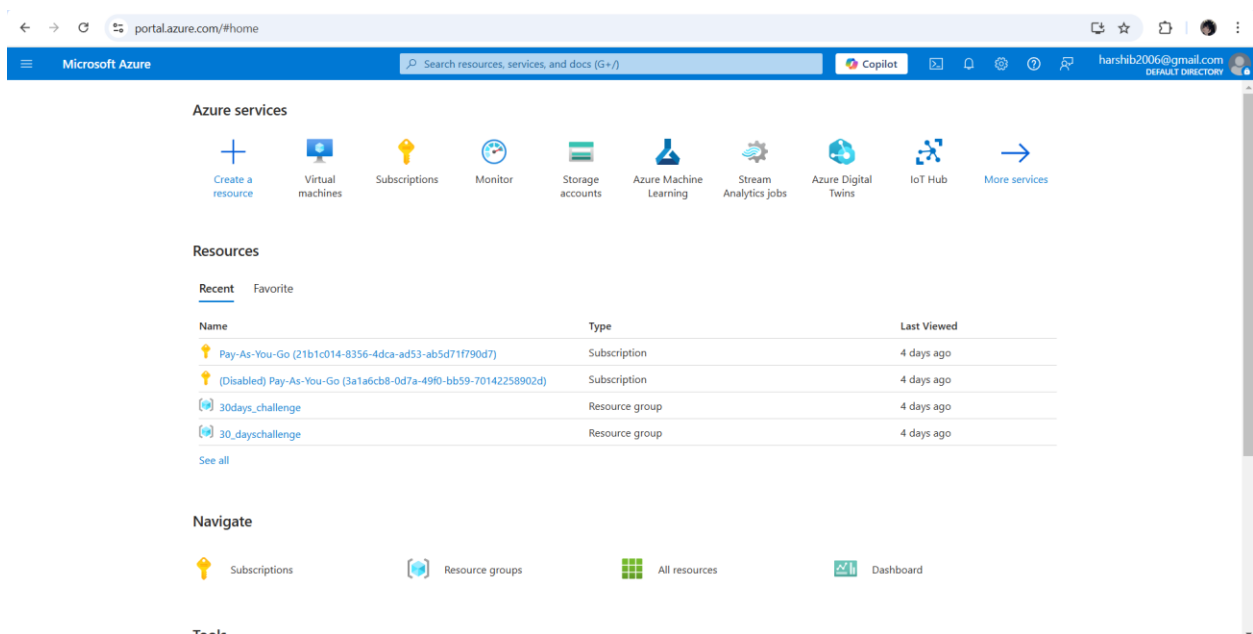
## *IMPORTANCE*

1. **Cost Efficiency** – Auto-scaling ensures that you only pay for the resources you use, reducing unnecessary infrastructure costs.
2. **High Availability** – Ensures applications remain available and responsive by scaling resources during peak loads.
3. **Performance Optimization** – Maintains optimal application performance by adding or removing instances based on demand.
4. **Resource Automation** – Reduces the need for manual intervention, enabling a more automated cloud infrastructure.
5. **Fault Tolerance** – Enhances system reliability by distributing workloads efficiently across multiple instances.

# *STEP BY STEP OVERVIEW*

Setting up an **Auto-scaling Group** in **Microsoft Azure** is done through **Virtual Machine Scale Sets (VMSS)**. This allows you to automatically adjust the number of virtual machines based on demand. Below are the steps to implement auto-scaling in **Azure**:

## Step 1: Create a Virtual Machine Scale Set (VMSS)

1. **Log in to Azure Portal**:
   Go to **Azure Portal**

2. **Navigate to "Virtual Machine Scale Sets"**: Search for **Virtual Machine Scale Sets** in the top search bar and click **Create**.



3. **Configure Basic Settings**:
   a. Select your **Subscription** and **Resource Group**.
   b. Choose a **Region** where you want the scale set to be deployed.
   c. Set the **Instance details**:
      i. **Virtual Machine Scale Set name**
      ii. **Orchestration mode**: Choose **Uniform** for standard VM scaling or **Flexible** for greater instance independence.
      iii. Select an appropriate **Image** (e.g., Ubuntu, Windows Server).
      iv. Choose a **VM Size** (e.g., Standard_DS2_v2).
4. **Configure Scaling Policy**:
   a. Choose **Scaling Policy**: Manual or Automatic.
   b. Set the **Minimum**, **Maximum**, and **Default** number of instances.

c. Select **Enable Autoscale**.

5. **Networking and Security**:

   a. Choose **Virtual Network (VNet)** and **Subnet**.

   b. Enable **Public IP** if needed.

   c. Configure **Load Balancer** (optional but recommended).

6. **Review and Create**:

   Click **Review + Create** and then **Create** to deploy the scale set.



# Step 2: Configure Auto-scaling Rules

1. **Go to the Scale Set**:

   a. In the **Azure Portal**, navigate to your **Virtual Machine Scale Set**.

2. **Configure Autoscale Settings**:

    a. Click on **Scaling**.

    b. Select **Custom Autoscale**.

    c. Click **+ Add a rule** to create scaling conditions.

3. **Define Scaling Triggers**:
    a. **Metric-based Scaling** (Recommended):
        i. Choose **CPU Percentage**, **Memory Usage**, or **Disk IOPS**.
        ii. Set **Thresholds** (e.g., scale out when CPU > 75% for 5 minutes).
        iii. Define **Action**: Increase instances by **1**.
        iv. Similarly, create a **Scale-in rule** (e.g., remove instances when CPU < 30%).
4. **Set Instance Limits**:
    a. **Minimum**: 1 (or as required)
    b. **Maximum**: Define based on expected workload
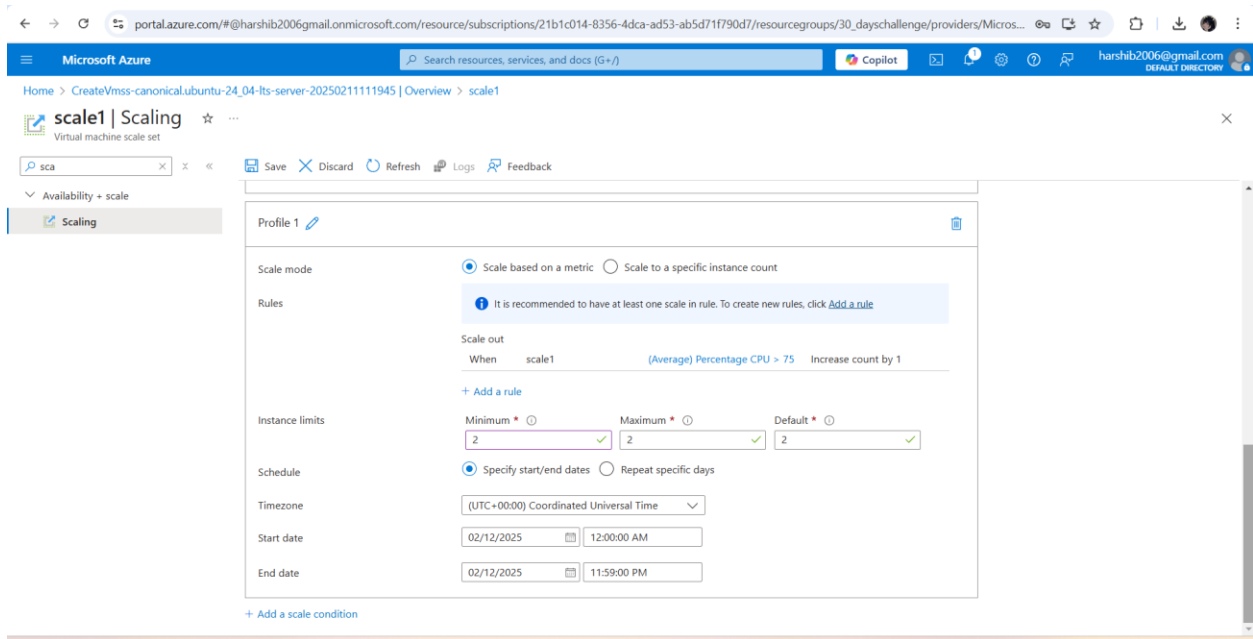    c. **Default**: Initial number of instances
5. **Apply and Save**.

## Step 3: Test the Auto-scaling Setup

1. **Generate Load**: Use a load-testing tool like **Apache JMeter**, **Azure Load Testing**, or **locust.io**.
2. **Monitor Scaling**:
   a. Go to **Azure Monitor** > **Metrics**.
   b. Check if new instances are created when demand increases.

3. **Verify Scale-in**:
   a. Reduce load and check if instances are removed.

## OUTCOME

By implementing auto-scaling in the cloud, businesses can achieve increased efficiency, optimized resource utilization, and reduced costs. Applications remain highly available and responsive to fluctuating workloads without human intervention. This approach ensures sustainable infrastructure management, paving the way for enhanced business continuity and scalability in a cloud-native environment.