# 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:** CHANDRU S

**Department:** INFORMATION TECHNOLOGY

# Introduction

Modern applications often experience fluctuating workloads, making optimal performance and availability crucial. Auto Scaling, a feature provided by cloud platforms like AWS, dynamically adjusts computing resources based on demand changes. This Proof of Concept (PoC) demonstrates how to set up an Auto Scaling Group (ASG) for virtual machines (VMs) to efficiently handle varying workloads. The process includes defining launch configurations, setting scaling policies, and testing automatic scaling based on CPU usage.

# Overview

This PoC focuses on building a scalable architecture using AWS Auto Scaling Groups. The key steps include:

1. **Defining a Launch Template:** Configuring VMs with specifications such as instance type, AMI, key pairs, and security groups.
2. **Creating an Auto Scaling Group:** Setting the initial group size and linking it to the launch template for dynamic instance management.
3. **Configuring Scaling Policies:** Establishing metrics like CPU utilization to trigger scaling actions, such as scaling up during high CPU usage.
4. **Testing Auto Scaling:** Simulating high CPU load to verify that the ASG launches additional instances when demand increases.

This PoC demonstrates the reliability, flexibility, and cost-efficiency of dynamic scaling in a cloud environment.

# Objectives

The primary objectives of this PoC are to:

1. Implement an Auto Scaling Group (ASG) to manage workloads effectively.
2. Define and configure a Launch Template for virtual machines.
3. Set up and test scaling policies based on predefined metrics, such as CPU utilization.
4. Validate the scaling process by simulating real-world scenarios (e.g., high CPU usage).

By completing this PoC, users will gain hands-on experience with Auto Scaling and understand its importance in ensuring application availability and cost management.

# Importance

Key benefits of implementing AWS Auto Scaling include:

- **Improved Application Availability:** Auto Scaling ensures applications remain available during traffic spikes by automatically adding more VMs.
- **Cost Optimization:** Resources are scaled down during low traffic periods, reducing unnecessary costs.
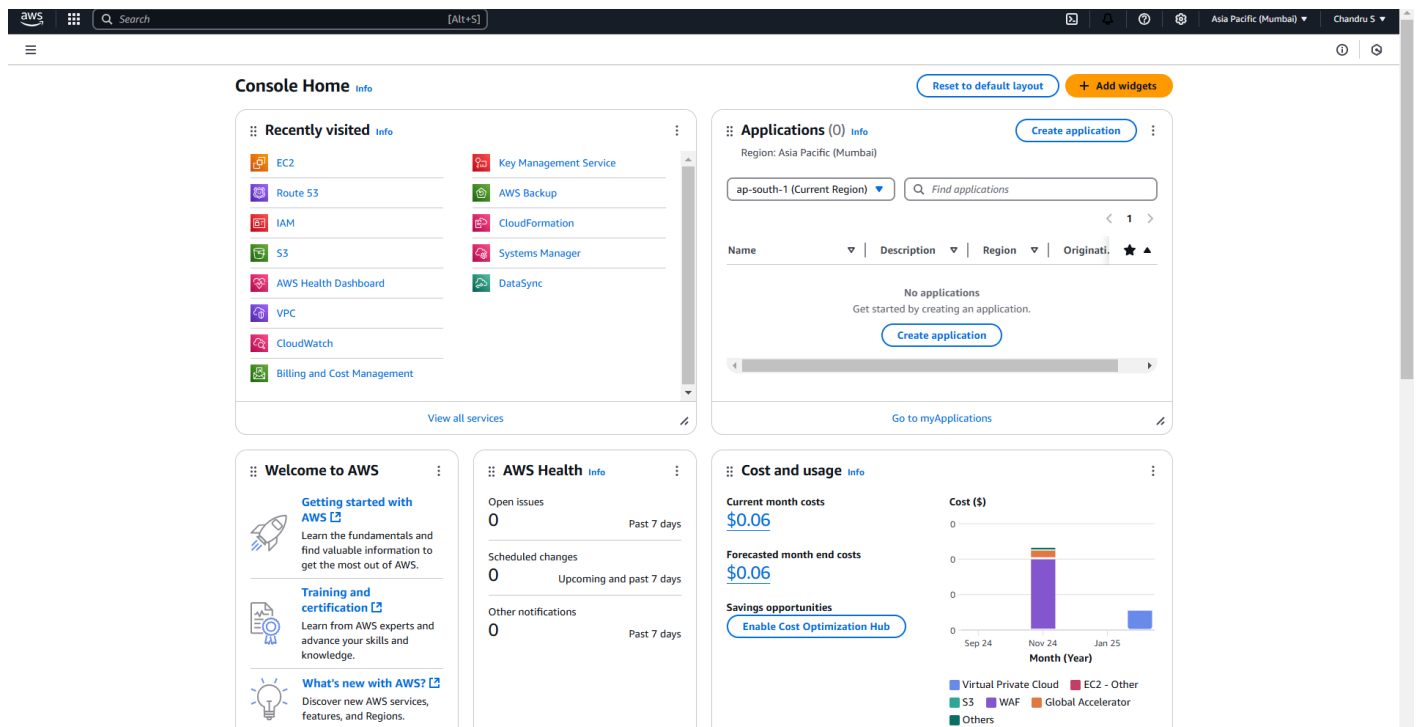
- **Efficient Resource Utilization:** Resources are provisioned based on actual demand, preventing over-provisioning and underutilization.
- **Resilience to Failures:** Unhealthy instances are automatically replaced, ensuring consistent performance.
- **Real-World Relevance:** Managing variable workloads is a critical cloud computing skill aligned with industry practices.

## Step-by-Step Overview
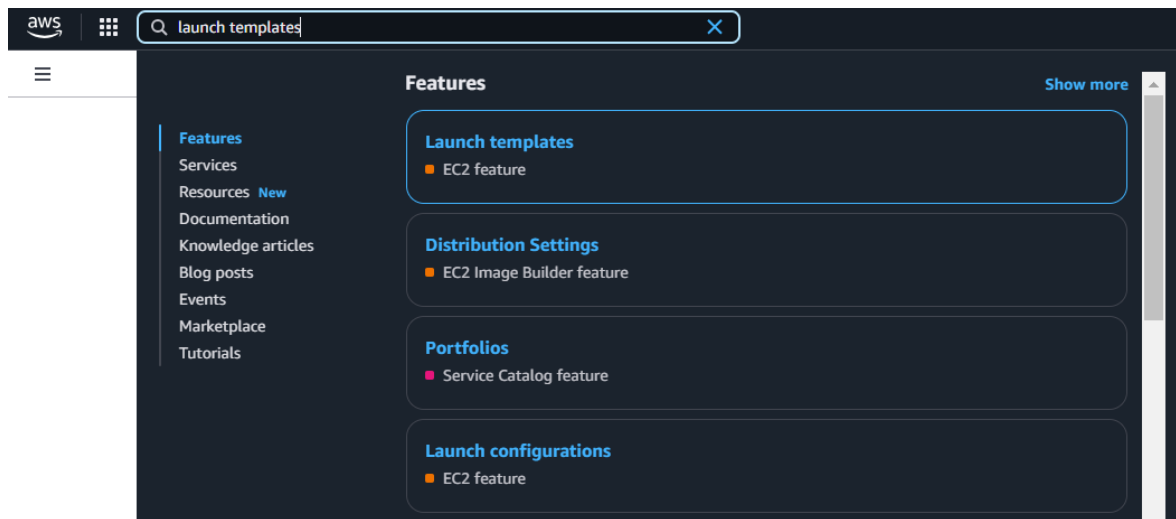
## Step 1:

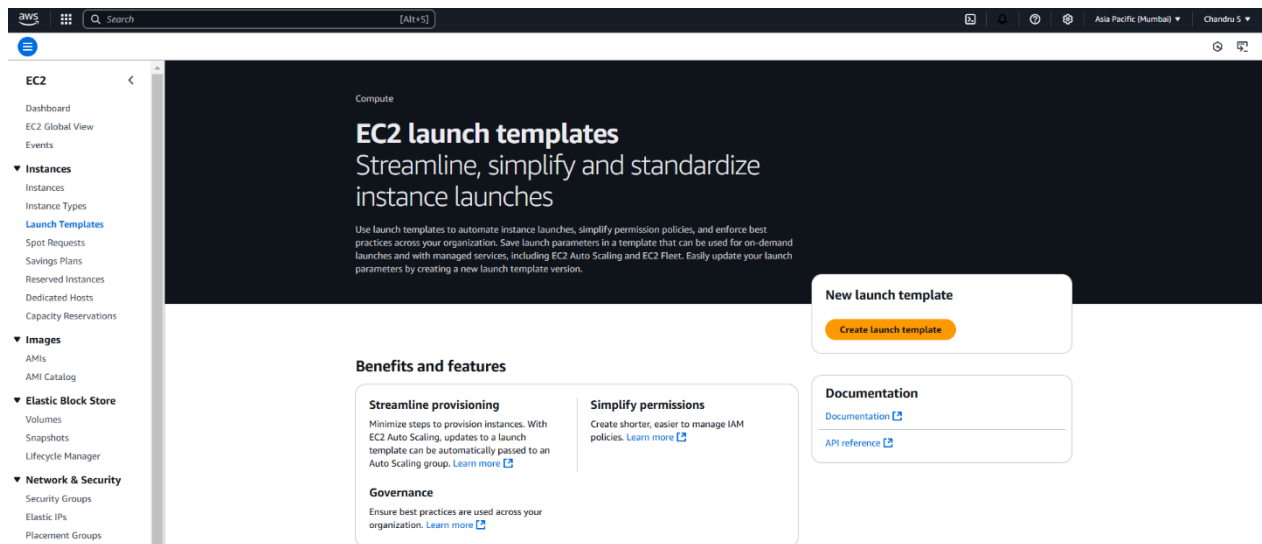**Access AWS Management Console**

- Log in to the AWS Management Console with your credentials.



## Step 2:
**Create Launch Template**

1. Search for **Launch Templates** in the AWS Console.

2. Click Create Launch Template.



3. Configure the template:

   o **Template Name:** AutoScalingTemplate

   o **AMI:** Amazon Linux 2 (or any default image)

   o **Instance Type:** t2.micro (Free-tier eligible)

   o **Key Pair:** Choose an existing key pair or create a new one for SSH access.

   o **Security Group:** Allow HTTP (port 80) and SSH (port 22).

4. Review the configuration and click Create Launch Template.



## Step 3:
**Create Auto Scaling Group**

1. Go to the **EC2 Dashboard**.

2. In the left sidebar, click **Auto Scaling Groups**.

3. Click Create Auto Scaling Group.



4. Configure the group:
   - **Group Name:** MyAutoScalingGroup
   - **Launch Template:** Select AutoScalingTemplate.

o **VPC and Subnets:** Use the default VPC and select at least two subnets in different Availability Zones for high availability.

5. Leave other settings as **default** and click **Next**.

6. Review the configuration and click **Create Auto Scaling Group**.

## Step 4:
**Testing Auto Scaling**
*Important Note: Avoid this test if you wish to prevent additional AWS costs.*

1. **Simulate High CPU Usage**
    - Connect to an EC2 instance in the Auto Scaling Group using SSH.
    - Install the stress package and simulate CPU load:

**sudo yum install -y stress**

**stress --cpu 2 --timeout 300**
    - This command utilizes 2 CPU cores for 5 minutes to simulate high usage.

2. **Monitor Scaling Activities**
    - Go to **AWS Management Console > EC2 Dashboard > Auto Scaling Groups**.
    - Select your Auto Scaling Group and navigate to the Activity History tab.
    - Check if new instances are launched based on the scaling policy (e.g., CPU utilization exceeding 50%).

3. **Terminate the Stress Test**
    - Stop the CPU load by pressing Ctrl+C in the terminal or terminating the stress process.

4. **Verify Scaling Down**
    - After CPU usage drops, check the Auto Scaling Group to confirm that unnecessary instances are terminated, returning the group to its desired capacity.

---

## Outcome

This Proof of Concept successfully demonstrated how AWS Auto Scaling dynamically manages EC2 instances based on workload demand, ensuring efficient resource utilization and cost-effectiveness. Key outcomes include:

1. **Launch Template and ASG Setup:** Successfully created a launch template and Auto Scaling Group with scaling policies.

2. **Dynamic Scaling and Monitoring:** Implemented scaling policies triggered by CPU utilization and verified scaling actions using the Activity History.

3. **Cost Awareness:** Highlighted potential costs of running additional instances beyond the AWS Free Tier and ensured optimal resource usage.

By completing this PoC, users gain hands-on experience with AWS Auto Scaling, enhancing their cloud computing skills for real-world applications.