# 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: NIDHISHA A DHAS

DEPARTMENT: AML

# INTRODUCTION:

**Auto Scaling** is an AWS services that ensure high availability, fault tolerance, and cost efficiency in cloud applications.

### Auto Scaling (ASG - Auto Scaling Group):

Auto Scaling automatically adjusts the number of EC2 instances in response to demand.

### What Auto Scaling Does:

- Increases Instances (Scale Out): When traffic or CPU usage increases, Auto Scaling launches new EC2 instances.
- Decreases Instances (Scale In): When traffic decreases, Auto Scaling terminates extra instances to reduce costs.
- Ensures Availability: If an instance fails (due to crash, hardware failure, etc.), Auto Scaling replaces it.
- Optimizes Costs: Ensures you only use the required number of instances at any given time.
- Elastic Load Balancer distributes incoming traffic across multiple EC2 instances to improve performance and availability.
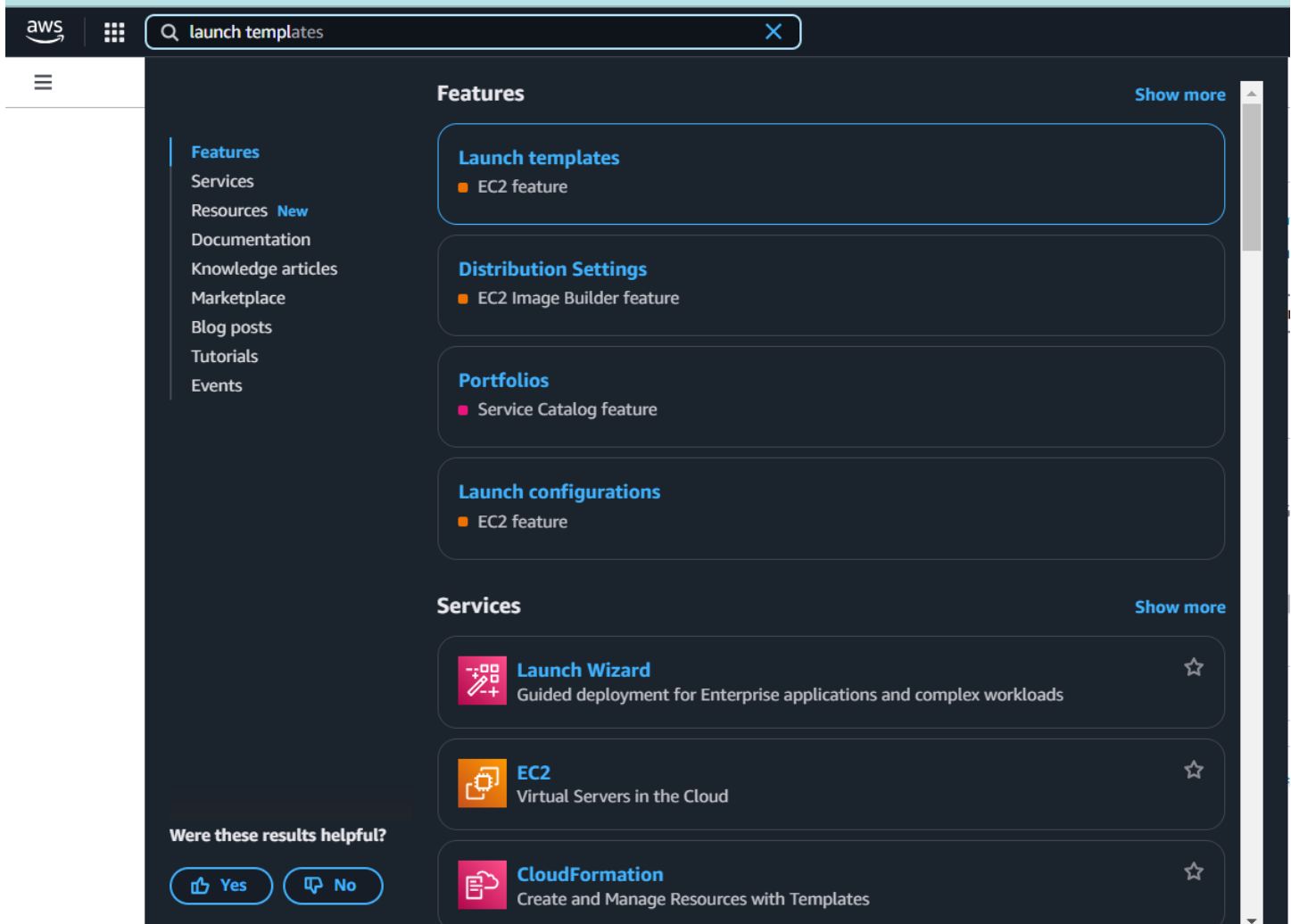
# IMPORTANCE:

- Auto Scaling Group monitors load (CPU, traffic) and adds/removes instances as needed.
- If traffic increases, ASG launches more instances, and ELB distributes traffic evenly.
- If traffic decreases, ASG removes extra instances, reducing costs.
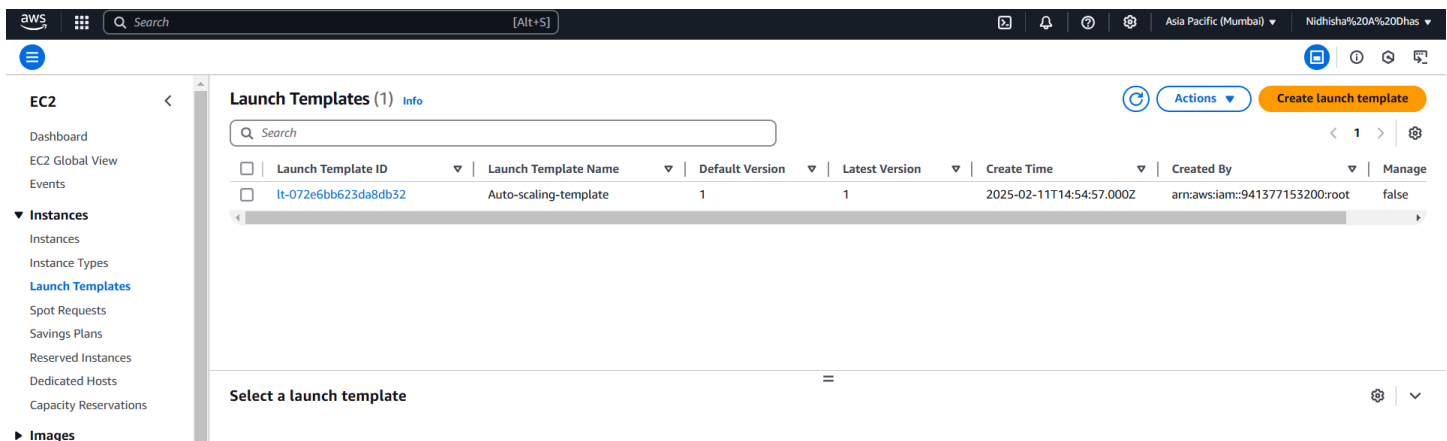- Result: High availability, better performance, and optimized cloud costs!

# STEP BY STEP OVERVIEW:

Step 1: LAUNCH TEMPLATES

- Login into your AWS console and search for launch templates.

- Navigate to it and click on launch template.
- Specify the template name, AMI image (Amazon Linux), Instance type (t2.micro), Key pair.
- Mainly, configure the security groups and add the security group role(such as: HTTP, ssh)
- Now launch the template.

## Step 2: AUTO SCALING GROUPS

- Go to your EC2 dashboard, on the left sidebar you will find 'Auto Scaling groups'. Click on create.
- Give the name to the ASG and select the template you have just created.



➕ Choose your VPC or the default VPC and select two subnets of two different availability zone.



➕ Leave the next setting as default. Now, review and click on create Auto Scaling group.

Step 3: TESTING AUTO SCALING

⚠️ Important Note: Do Not Perform This Test If You Want to Avoid Costs:

1. Launching and running additional EC2 instances will incur charges beyond the AWS Free Tier.

2. Simulating high CPU usage and triggering scaling may increase costs temporarily due to additional resource allocation.

TESTING:

1. Simulate High CPU Usage on an EC2 Instance: Connect to one of your EC2 instances in the Auto Scaling Group using SSH.
2. Run a command to create artificial CPU load. For example: sudo yum install -y stress, stress--cpu 2--timeout 300.

This command will utilize 2 CPU cores for 5 minutes, simulating high CPU usage.

MONITOR SCALING ACTIVITIES:

Navigate to the AWS Management Console > EC2 Dashboard > Auto Scaling Groups.

Select your Auto Scaling Group and go to the Activity History tab.

Check if a new instance is being launched based on your scaling policy (e.g., CPU utilization exceeding 50%).

 TERMINATE THE STRESS TEST:

Once testing is done, stop the CPU load by pressing Ctrl+C in the terminal or by terminating the stress process.

 VERIFY SCALING DOWN:

After the CPU usage drops, monitor the Auto Scaling Group again to confirm that unnecessary instances are terminated, returning to the desired capacity.

## **CONCLUSION:**

By completing this PoC, we will be able to:

- Implement Auto Scaling group and launch the Template.
-  Attain dynamic scaling and monitoring.
- Most importantly, we should be aware of Costs.