# SCALING AND LOAD BALANCING AN ARCHITECTURE

AWS provides the feature of auto scaling and load balancing to make your architecture more robust.

ELB automatically distributes incoming application traffic across multiple Amazon Elastic Compute Cloud (Amazon EC2) instances. ELB provides the amount of load balancing capacity needed to route application traffic to help you achieve fault tolerance in your applications.

Auto scaling can automatically increase the number of EC2 instances during spikes in demand to maintain performance and can decrease capacity during lulls to reduce costs. Auto scaling is well suited to applications that have stable demand patterns or that experience hourly, daily, or weekly variability in usage.

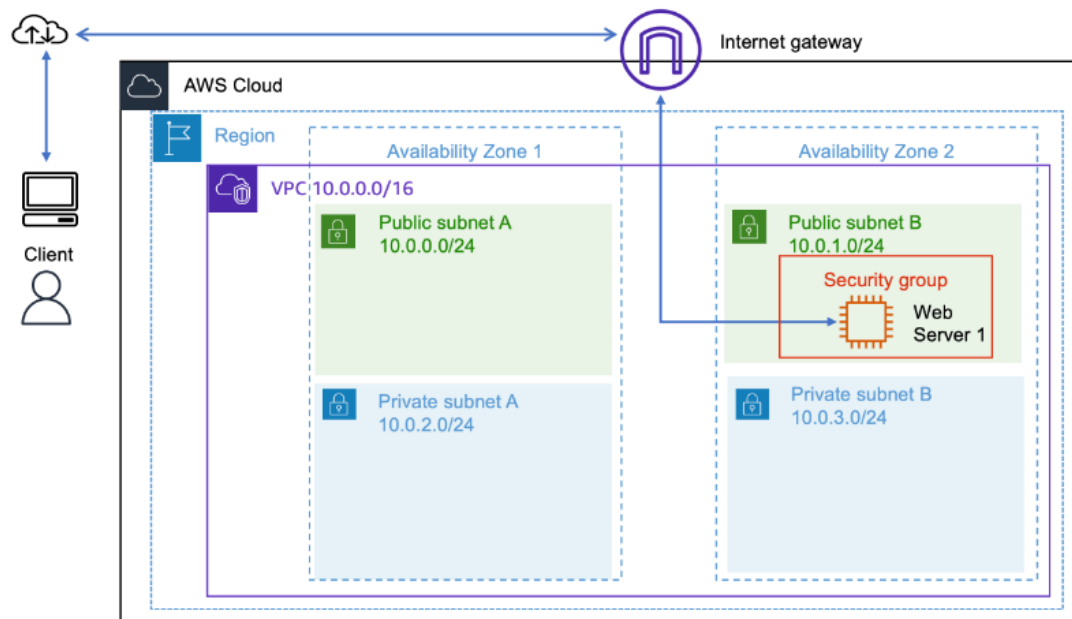With that said, let's make it happen!!

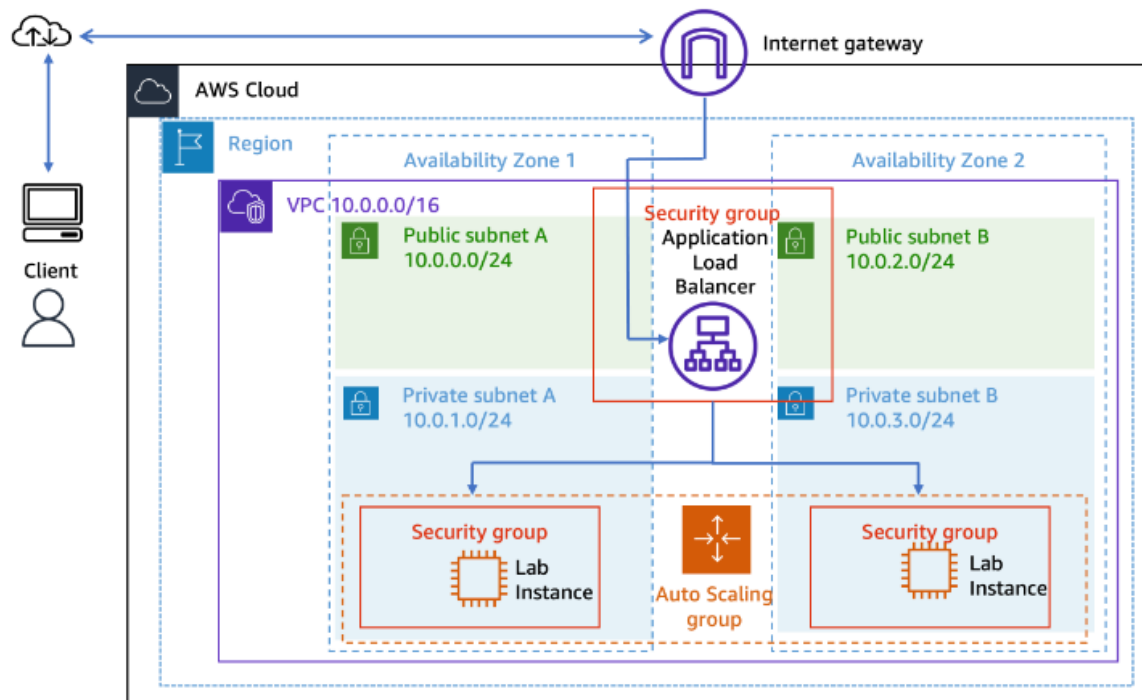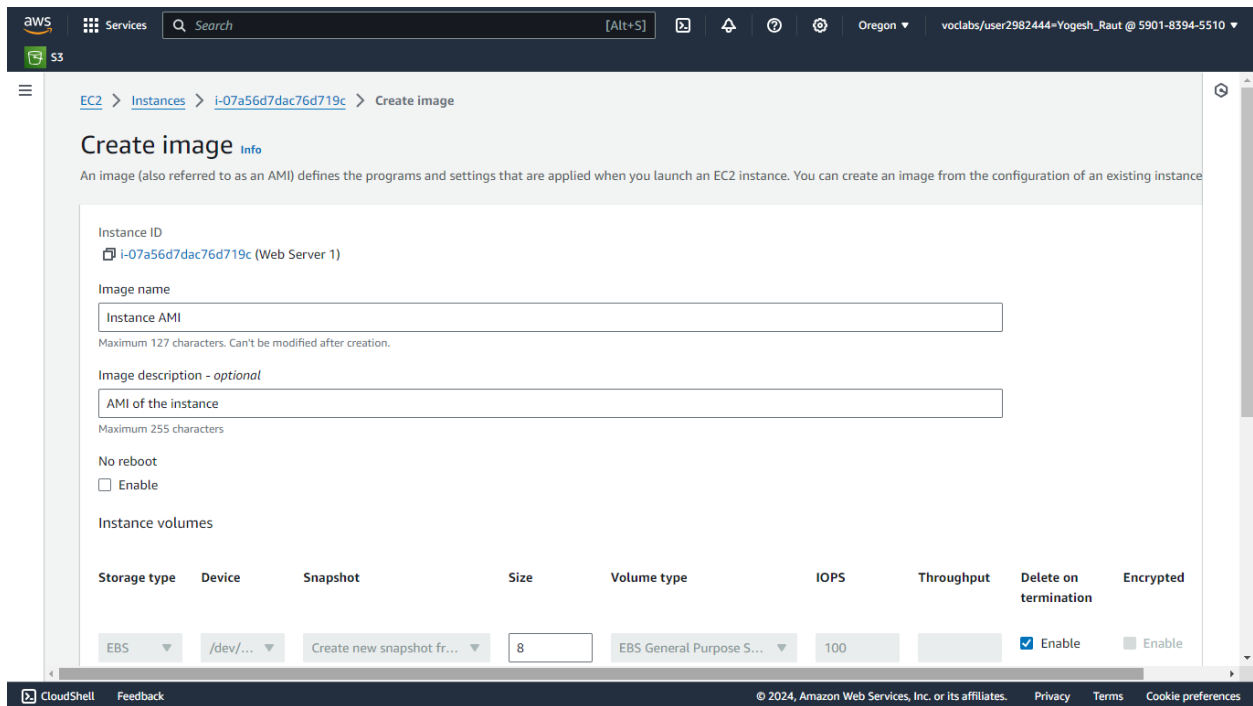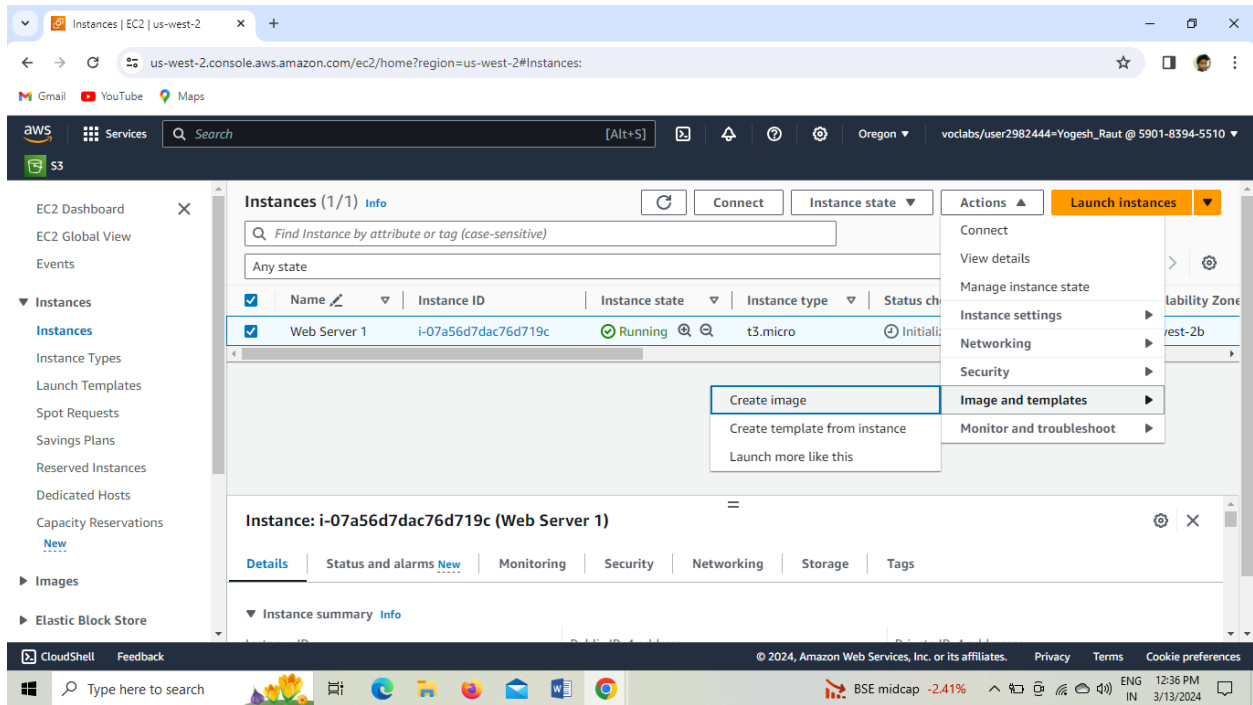## REFERANCE DIAGRAMS:



**Fig 1. Initial Architecture**

**Fig 2. Final Architecture**

STEPS TO FOLLOW:

A. Creating an Amazon Machine Image (AMI) of an EC2 Instance

1. In **AWS Management Console**, search for **EC2** after which you select **Instances** from the left navigation pane to list the instances.
2. Select your instance, and from **Actions > Image and templates**, select **Create Image**. Now, configure the image as follows.
   I.   **Image name:** Instance AMI
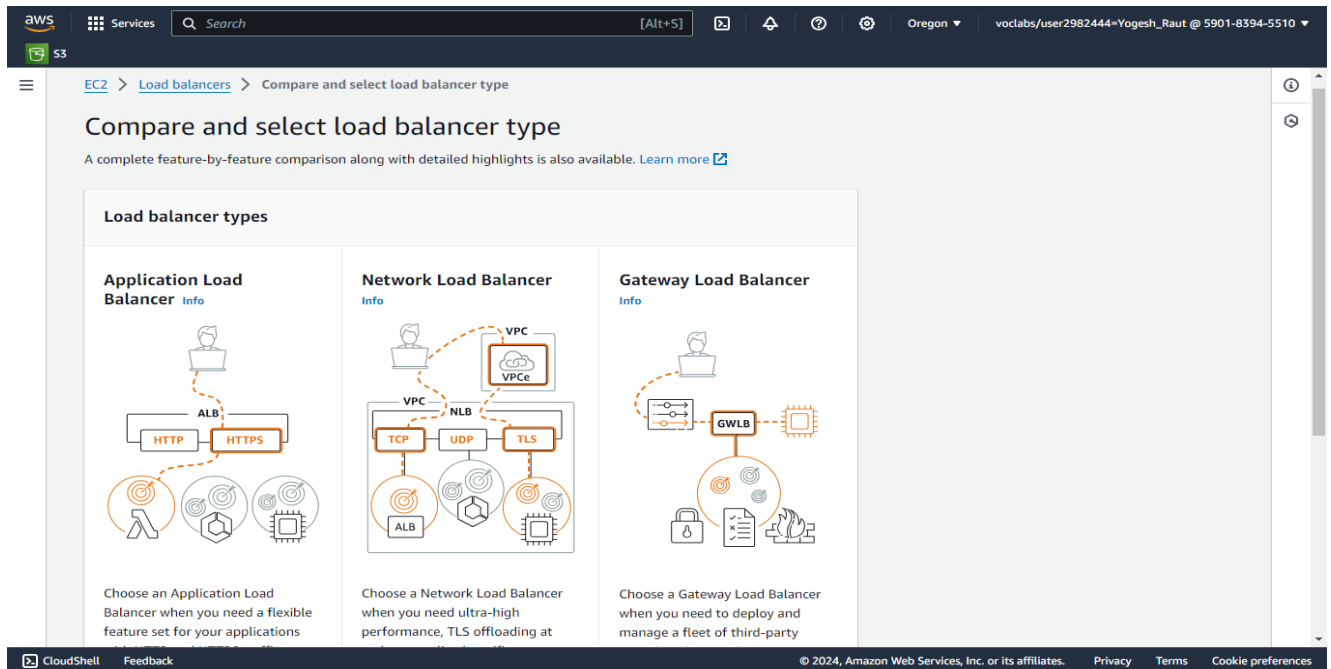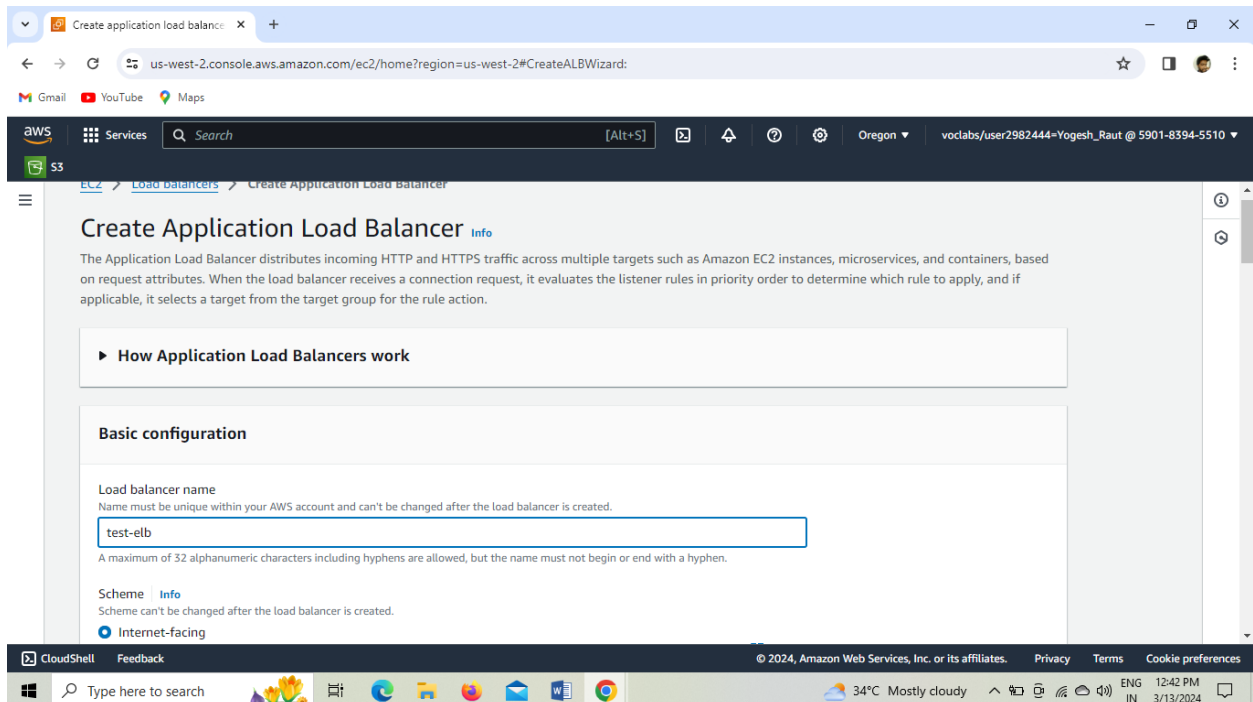   II.  **Description:** AMI of the instance

**3.** Select **Create Image.**

## B. Creating a load balancer

**1.** Select **Load Balancing** section, where you'll find **Load balancer** from left navigation pane.

**2.** Choose **Create Load balancer.** If asked for **Load balancer type**, select **Application Load Balancer**.



**3.** In **Basic Configuration** section, give **Load balancer name** as **test-elb**.



**4.** In **Network mapping** section, configure the settings as follows:
  **I.** **VPC:** Select your VPC.

**II.** **Mappings:** Choose both Availability Zones. For **AZ1** and **AZ2,** choose **Public Subnets 1 and 2** respectively.



**5.** **Security Groups:** Web Security Group (**HTTP permitted**)



**6.** In the **Listeners and routing** section, choose the **Create target group** link.

7. On the new **Target groups browser tab,** in the **Basic configuration** section, configure the following:

       I.    **Target type:** Instances

      II.    **Target Group Name:** test-target-group

Click **Next.**



8. On the **Register targets** page, choose **Create target group.** Once the target group has been created successfully, close the **Target groups browser tab.**

**9.** Return to the **Load balancers browser tab**. In the **Listeners and routing** section, choose **Refresh.** Now from the **Forward to dropdown list,** choose **test-target-group** and choose **Create Load Balancer.**

**10.** To view the load balancer that you created, choose **View load balancer** and copy **DNS Name.**
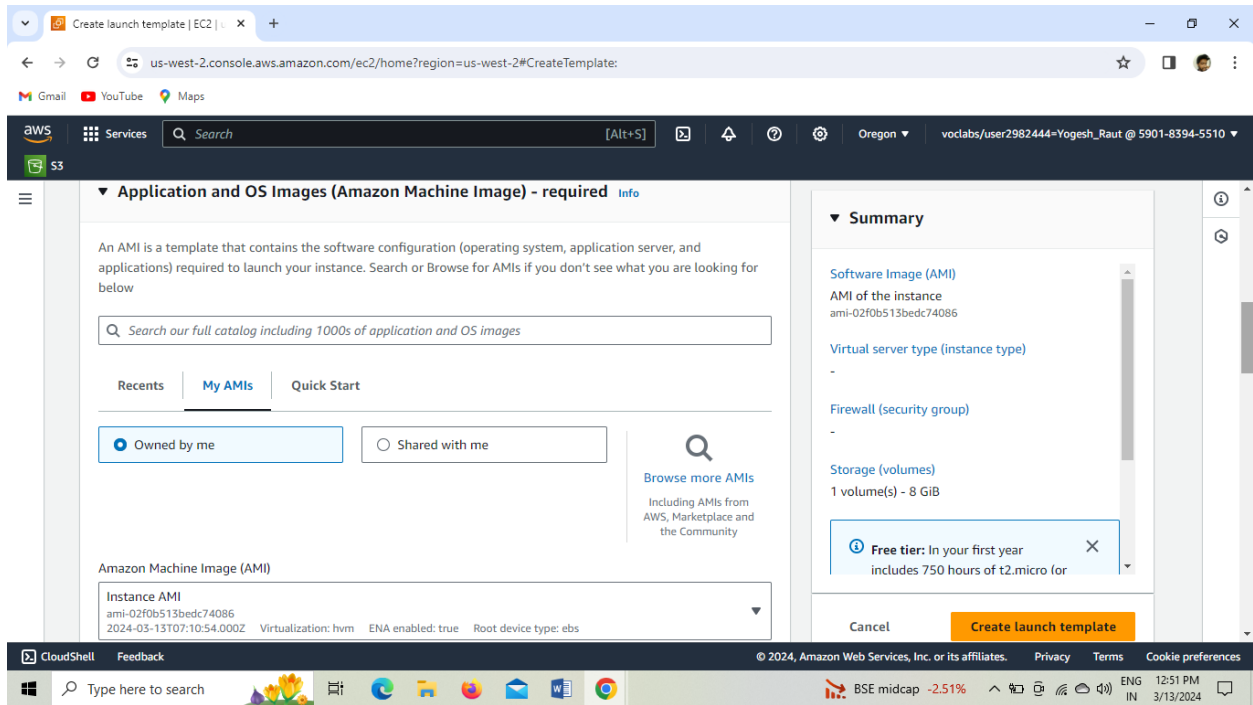
## C. Create a launch template

1. At the top of the **AWS Management Console**, in the search bar, enter and choose **EC2**.
2. In the left navigation pane, locate the **Instances** section, and choose **Launch Templates** and select **Create launch template**.
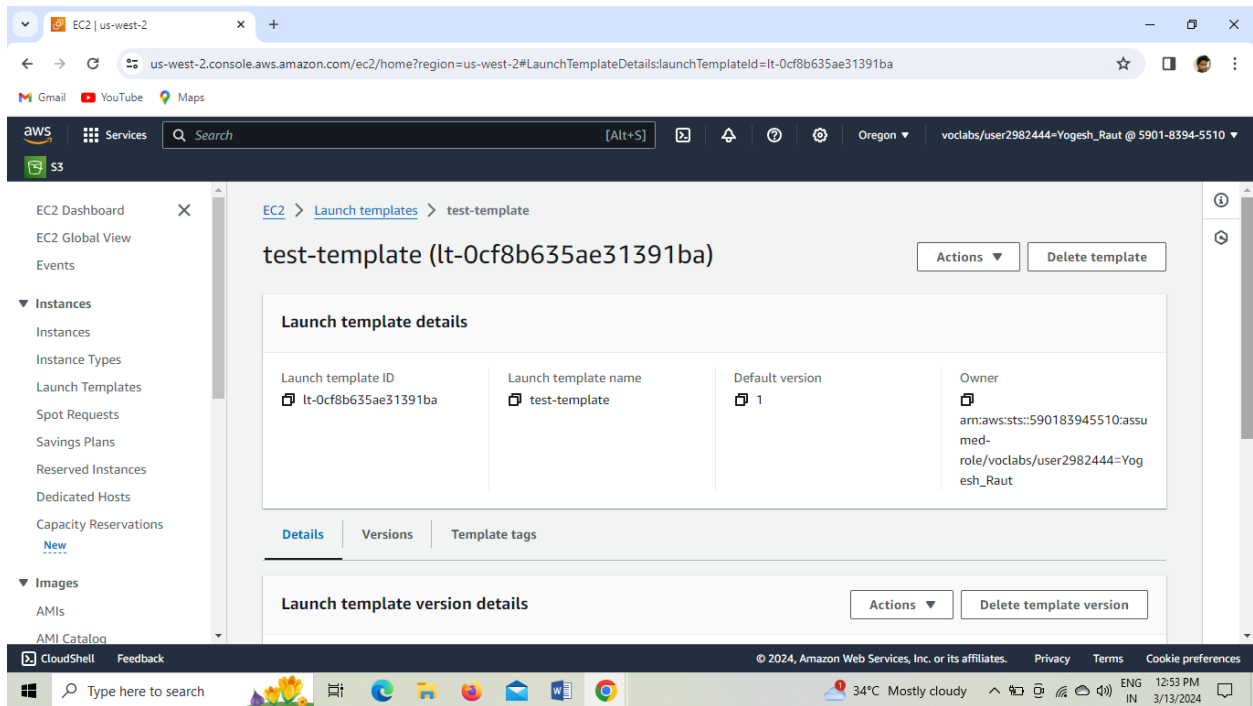


3. Now **configure the launch template** as follows:
      I. **Launch template name:** test-template
      II. **Template version description:** A web server for the load test app
      III. **Auto Scaling guidance:** Provide guidance to help me set up a template that I can use with EC2 Auto Scaling.
      IV. **Application and OS Images (Amazon Machine Image):** Instance AMI
      V. **Instance Type:** t3.micro
      VI. **Key pair name:** Don't include in launch template**.**
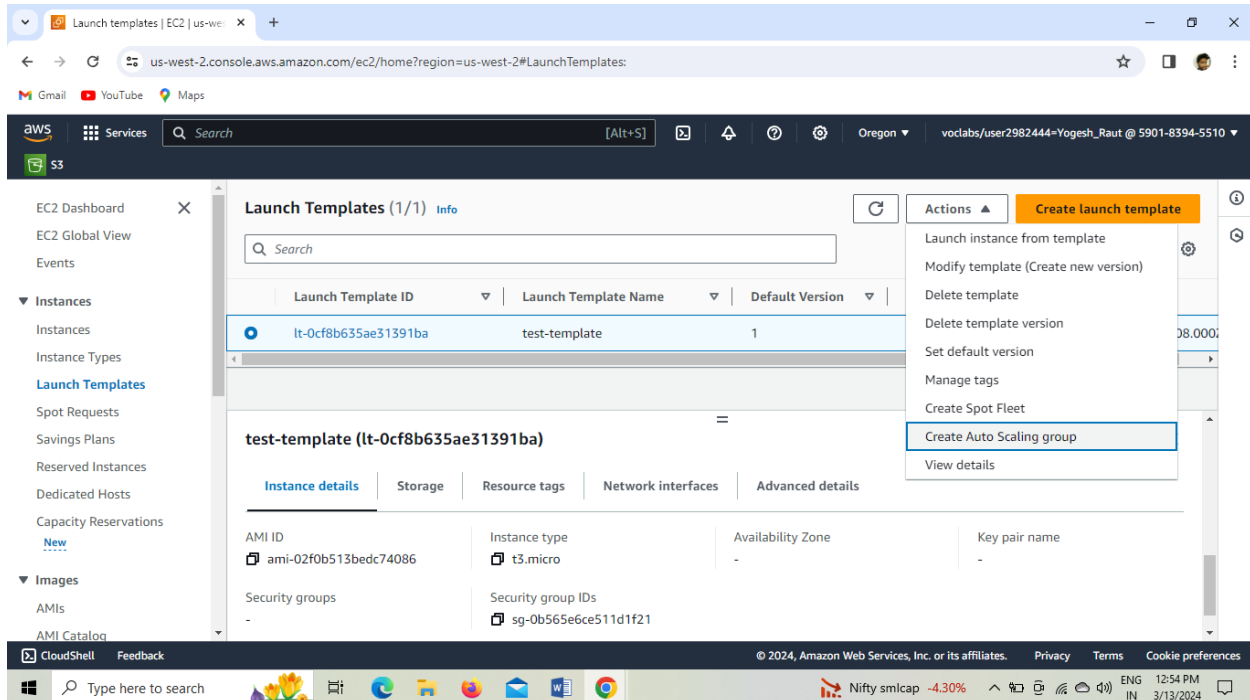      VII. **Network Settings > Security Group:** Web Security Group.

4. Select **Create launch template.**

## D. Creating an Auto Scaling group

1.  Select the **launch template** you created, and choose **Create Auto Scaling group** from **Actions** tab.



2.  Configure the group as follows:

I.      **Name:** test-auto-scaling-group
II.     **Network:**
        a.   **VPC:** Lab VPC
        b.   **Availability Zones and Subnets:** Private Subnets 1 and 2.

III.    **Configure advanced options:**
        a.   **Load balancing:** Attach to an existing load balancer.
             • **Choose from your load balancer target groups:** test-target-group
IV.     **Health check type:** ELB
V.      **Group size:**
        a.   **Desired capacity:** 2
        b.   **Minimum capacity:** 2
        c.   **Maximum capacity:** 4

VI.     **Scaling policies:** Target tracking scaling policy
VII.    **Metric type:** Average CPU utilization.
VIII.   **Target Value:** 50

**3.** Choose **Create Auto Scaling group.**



## E. Verifying if Load Balancer works:

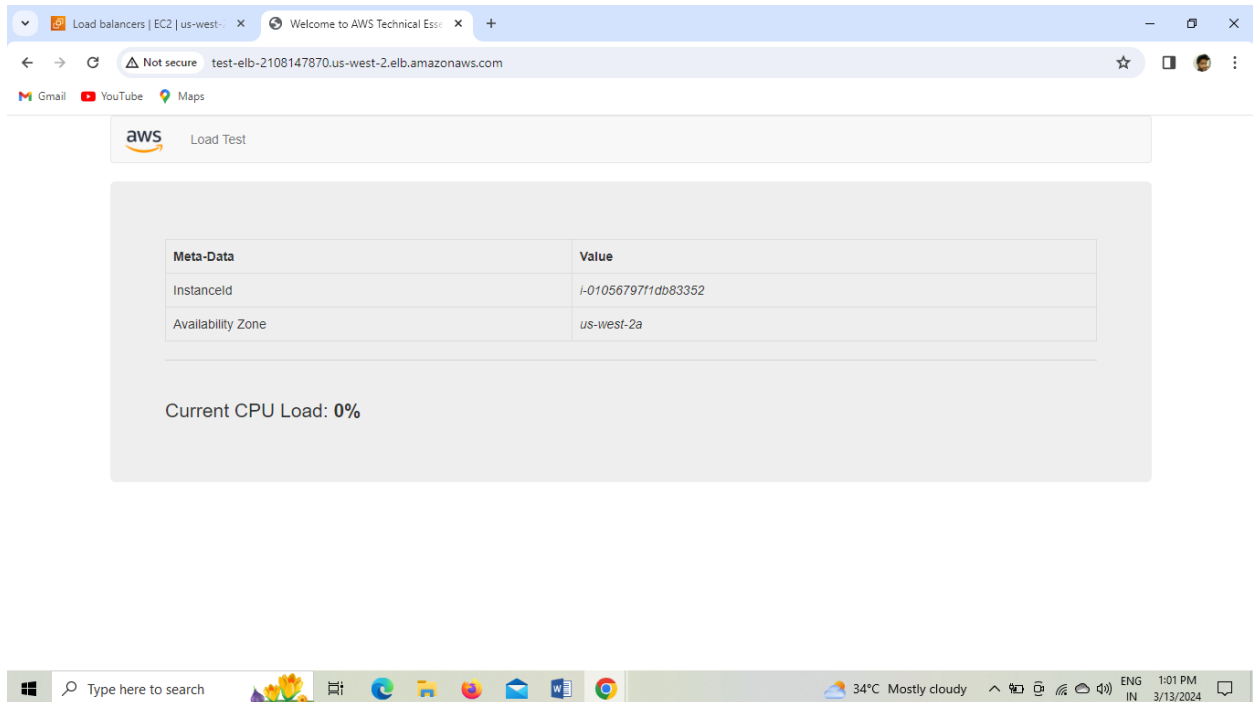**1.** To verify load balancer, **check Instances** if additional instances are added.

2. In left navigation pane, in **Load Balancing** section, check **if new targets are registered in your target group** and that they are **healthy.**



3. Open a new web browser tab, paste the **DNS name** that you copied before, and press Enter.

**The Load Test application should appear in your browser**, which means that the load balancer received the request, sent it to one of the EC2 instances, and then passed back the result.



## F. Testing Auto Scaling

1. Return to the **AWS Management Console**, but keep the Load Test application tab open.
2. In the AWS Management Console, in the search bar, enter and choose **CloudWatch**
3. In left navigation pane, choose **Alarms.** In **All Alarms**, check if **AlarmHigh** state is **OK.**
4. Choose **Load Test** from Load Test application tab to **increase system load.**
5. Go back to **CloudWatch,** and check **if Alarm Status changes.**
6. Once the **Alarm Status** changes to **In alarm,** new **instances will be launched**.

Congrats!! Your system is now automatically scalable and balancing the load too.