

## **Placement Empowerment Program**

### ***Cloud Computing and DevOps Centre***

Set Up a Load Balancer in the Cloud Configure a load balancer to distribute traffic across multiple VMs hosting your web application.

Name: SAMRAJ K

Department: CSE

# Introduction

In this Proof of Concept (POC), the focus is on setting up a cloud-based Load Balancer using AWS to distribute traffic across multiple virtual machines (EC2 instances). Load Balancers play a crucial role in modern cloud architectures by ensuring high availability, fault tolerance, and scalability for web applications. This POC demonstrates the basic setup of an AWS Load Balancer, allowing traffic to be distributed between two EC2 instances running simple web servers.

## Overview

The POC covers the following:

- 1. Creating EC2 Instances:** Setting up two virtual machines (WebServer1 and WebServer2) in the AWS Free Tier.
- 2. Configuring Web Servers:** Installing and configuring Apache HTTP Server on each instance to host simple HTML web pages.
- 3. Setting Up a Load Balancer:** Creating an Application Load Balancer (ALB) to distribute incoming traffic evenly between the two EC2 instances.
- 4. Testing the Load Balancer:** Verifying that the Load Balancer works by checking the DNS name and ensuring it alternates traffic between the two servers.

# Objectives

1. To understand the process of creating and configuring EC2 instances in AWS.
2. To install and configure a web server (Apache HTTP Server) on Linux-based EC2 instances.
3. To set up an Application Load Balancer to distribute traffic across multiple servers.
4. To validate that the Load Balancer works as intended by testing it with unique responses from each server.
5. To build a foundational understanding of cloud-based load balancing for real-world use cases.

# Importance

- 1. Scalability:** Demonstrates how load balancing allows scaling applications by adding or removing servers as traffic demands change.
- 2. Fault Tolerance:** Ensures that if one server goes down, the Load Balancer redirects traffic to the healthy server, improving reliability.
- 3. Cost Efficiency:** Explores how to leverage AWS Free Tier services to test and deploy cloud-based solutions with minimal cost.
- 4. Hands-On Experience:** Provides practical experience in configuring essential AWS services, an important skill for cloud computing professionals.
- 5. Foundation for Advanced Concepts:** Sets the stage for more complex setups, such as auto-scaling, secure traffic distribution, and monitoring solutions.

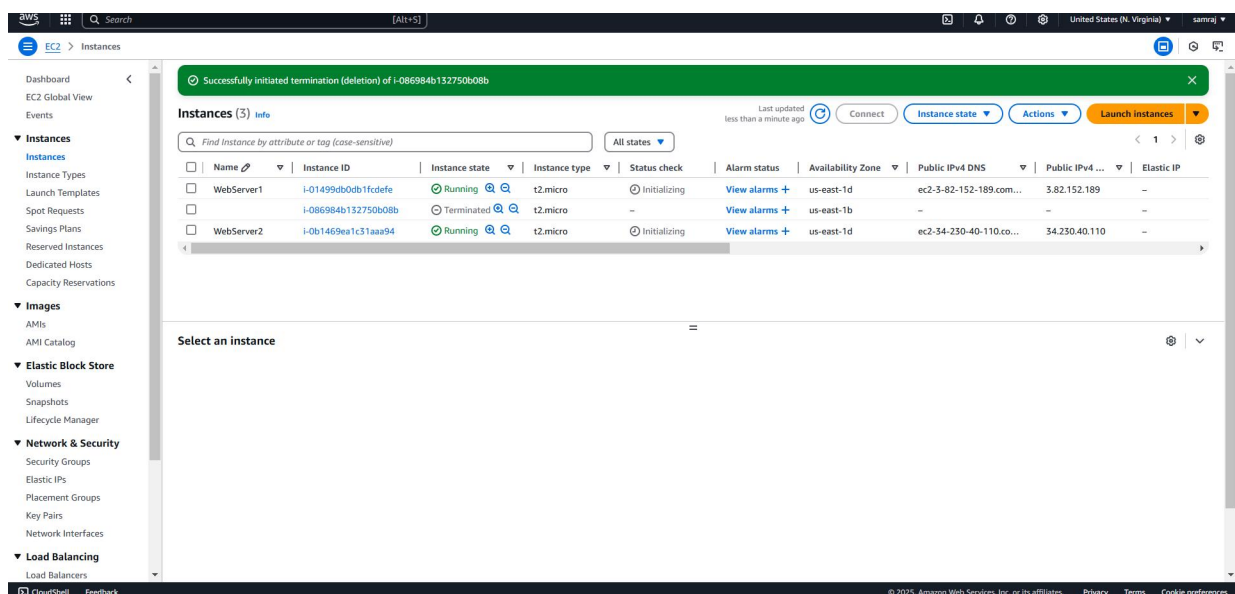
# Step-by-Step Overview

## Step 1:

1. Go to [AWS Management Console](#).
2. Enter your username and password to log in.

## Step 2:

To create your instances, click **Launch Instance** and fill in the details: name the first instance "WebServer1," select **Amazon Linux 2 AMI (Free Tier eligible)** as the OS, and choose the **t2.micro** instance type. For the Key Pair, either select an existing one or create a new key pair to use for SSH access. Under **Network Settings**, click "Edit" and ensure "Allow HTTP traffic from the internet" is checked to enable web traffic. Keep the storage size at the default 8 GB, then click **Launch Instance**. Repeat the same steps for the second instance, naming it "WebServer2."



## Step 3:

Click on **WebServer1**, then click **Connect**.

Use the instructions under **SSH client** to connect to your instance via terminal.

The screenshot shows the AWS Management Console interface for connecting to an EC2 instance. The breadcrumb trail is **EC2 > Instances > i-01499db0db1fcdefe > Connect to instance**. The page title is **Connect to instance** with an info icon. Below the title, it says "Connect to your instance i-01499db0db1fcdefe (WebServer1) using any of these options". There are four tabs: **EC2 Instance Connect**, **Session Manager**, **SSH client** (which is selected), and **EC2 serial console**.

Under the **SSH client** tab, the **Instance ID** is `i-01499db0db1fcdefe (WebServer1)`. The instructions are:

1. Open an SSH client.
2. Locate your private key file. The key used to launch this instance is DEV2.pem
3. Run this command, if necessary, to ensure your key is not publicly viewable.  
`chmod 400 "DEV2.pem"`
4. Connect to your instance using its Public DNS:  
`ec2-3-82-152-189.compute-1.amazonaws.com`

A green tooltip "Command copied" points to the terminal command:  
`ssh -i "DEV2.pem" ec2-user@ec2-3-82-152-189.compute-1.amazonaws.com`

A blue note box contains the text: **Note:** In most cases, the guessed username is correct. However, read your AMI usage instructions to check if the AMI owner has changed the default AMI username.

A "Cancel" button is located at the bottom right of the panel.

The footer of the console shows "CloudShell Feedback" on the left, and "© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences" on the right.

## Step 3:

Run the following commands to install and start a web server

```
PS C:\Users\samni> cd downloads
PS C:\Users\samni\downloads> ssh -i "DEV2.pem" ec2-user@ec2-3-82-152-189.compute-1.amazonaws.com
[ec2-user@ip-172-31-91-108 ~]$ echo "Hello from Webserver1" | sudo tee /var/www/html/index.html
Hello from Webserver1

[ec2-user@ip-172-31-91-108 ~]$ sudo yum update -y
[ec2-user@ip-172-31-91-108 ~]$ sudo yum install httpd -y
[ec2-user@ip-172-31-91-108 ~]$ sudo systemctl start httpd
[ec2-user@ip-172-31-91-108 ~]$ sudo systemctl enable httpd

[ec2-user@ip-172-31-91-108 ~]$ echo "Hello from Webserver1" | sudo tee /var/www/html/index.html
Hello from Webserver1

[ec2-user@ip-172-31-91-108 ~]$ exit
logout
Connection to ec2-3-82-152-189.compute-1.amazonaws.com closed.
```

## Step 4:

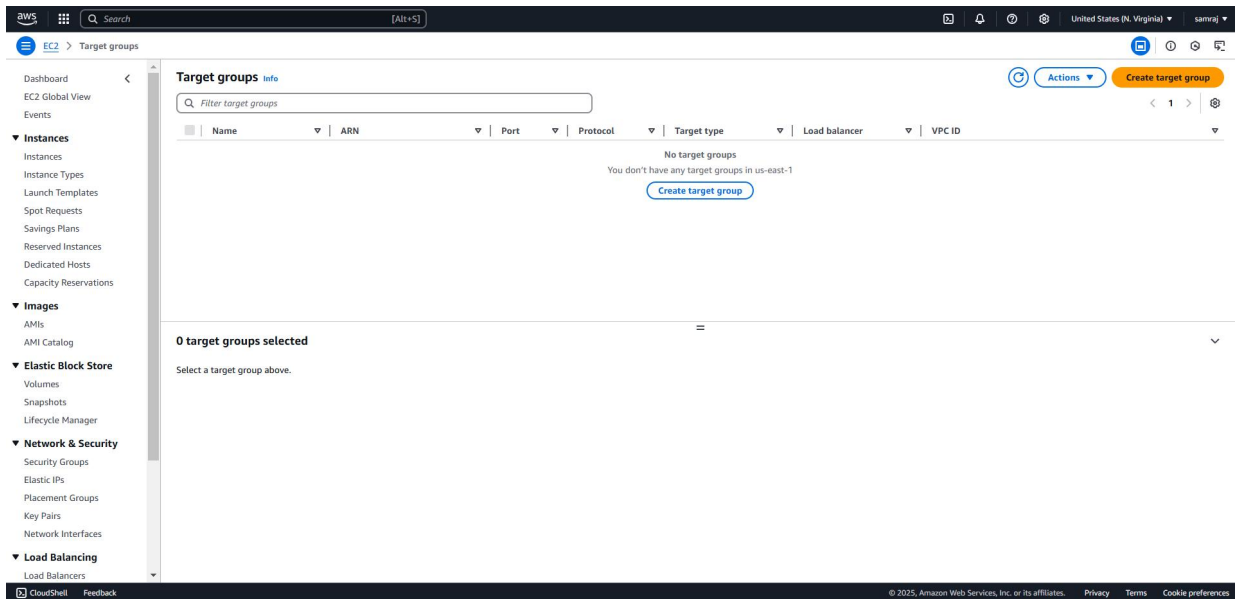
Repeat these steps for **WebServer2** but change the message in the last command to:

```
PS C:\Users\samni> cd downloads
PS C:\Users\samni\downloads> ssh -i "DEV2.pem" ec2-user@ec2-34-230-40-110.compute-1.amazonaws.com

[ec2-user@ip-172-31-93-44 ~]$ echo "Hello from webserver2" | sudo tee /var/www/html/index.html
Hello from webserver2
```

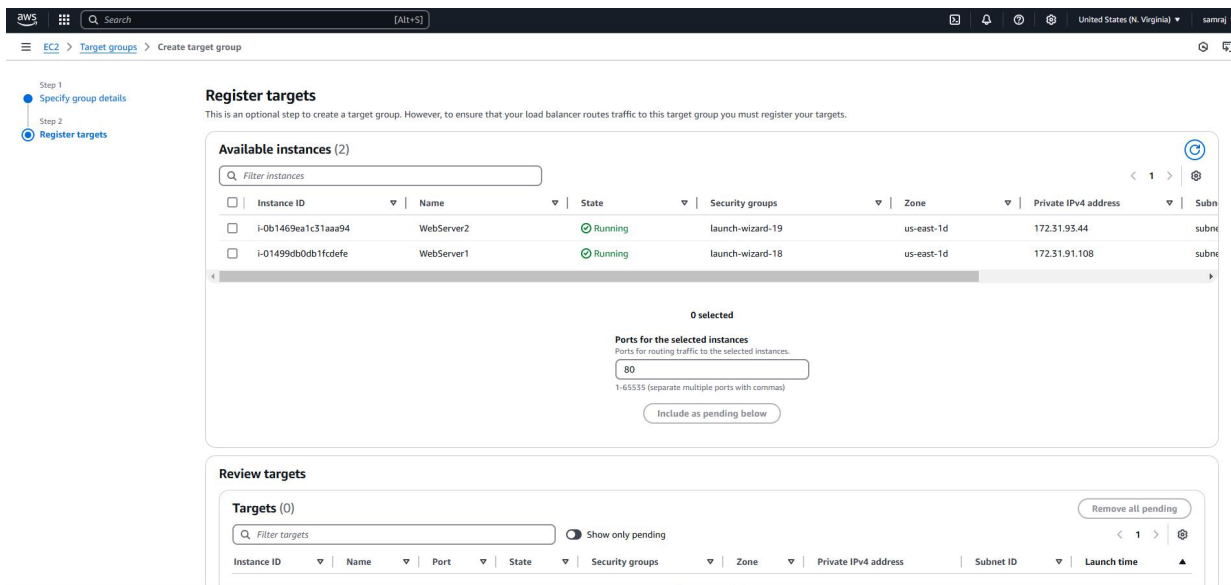
## Step 5:

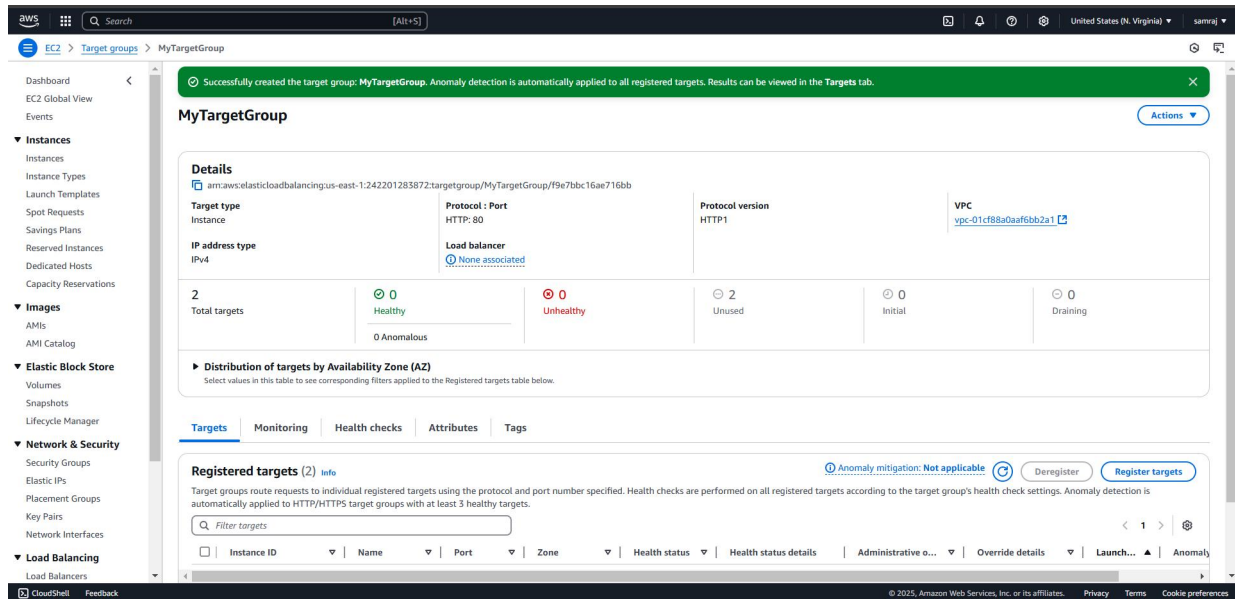
1. In the **AWS Management Console**, go to the **EC2 Dashboard**.
2. Scroll down and click on **Target Groups** under "Load Balancing."
3. Click **Create Target Group**.



## Step 6:

To create a target group, select **Instances** as the target type, name it (e.g., "MyTargetGroup"), set the **Protocol** to HTTP and **Port** to 80, and choose the same VPC as your EC2 instances (usually the default VPC). Keep the **Health Check Path** as / to verify the web server's status. Click **Next**, select both WebServer1 and WebServer2 under "Register Targets," click **Include as pending below**, and then create the target group.





## Step 7:

In the EC2 Dashboard, go to **Load Balancers** under "Load Balancing" and click **Create Load Balancer**. Select **Application Load Balancer (free tier eligible)** and configure it: name it (e.g., "MyALB"), set the **Scheme** to Internet-facing, **IP Address Type** to IPv4, and ensure the listener is HTTP on port 80. Select the VPC and at least two subnets for high availability. Skip the security settings since this is HTTP. On the **Security Groups** page, choose or create a security group that allows HTTP traffic. On the **Routing** page, select



the previously created target group (e.g., "MyTargetGroup") and click **Create Load Balancer**.

This screenshot shows the AWS Management Console 'Load balancers' page. The left-hand navigation pane is expanded to 'Load Balancing' > 'Load Balancers'. The main content area, titled 'Load balancers', includes a description: 'Elastic Load Balancing scales your load balancer capacity automatically in response to changes in incoming traffic.' Below this is a search bar labeled 'Filter load balancers'. A table with columns 'Name', 'DNS name', 'State', 'VPC ID', 'Availability Zones', 'Type', and 'Date created' is present. A message states: 'No load balancers. You don't have any load balancers in us-east-1.' with a 'Create load balancer' button. At the bottom, a section '0 load balancers selected' prompts the user to 'Select a load balancer above.' The top of the console shows the user is in the 'United States (N. Virginia)' region.

This screenshot shows the AWS Management Console 'Load balancers' page after a load balancer has been created. The left-hand navigation pane remains the same. The main content area now shows 'Load balancers (1)'. The table lists one load balancer: 'MyALB' with a DNS name 'MyALB-557583686.us-east-1...', a state of 'Provisioning...', VPC ID 'vpc-01cf889a0aaf6bb2a1', 2 Availability Zones, and type 'application'. It was created on 'February 12, 2025, 21:16 (UTC+05:30)'. The '0 load balancers selected' section at the bottom remains unchanged. The top of the console shows the user is in the 'United States (N. Virginia)' region.

## Step 8:

To verify the functionality of your Load Balancer:

1. Go to the **Load Balancers** section in the AWS Management Console.
2. Select your Load Balancer and find its **DNS name** under the **Description** tab.
3. Copy the DNS name and open it in your browser.
4. Refresh the page to confirm that traffic is being alternated between the two EC2 instances. You should see the messages **"Hello from WebServer1!"** and **"Hello from WebServer2!"** displayed alternately.

This confirms that the Load Balancer is correctly distributing traffic and ensuring high availability.

## Outcome

By completing this POC of setting up an Application Load Balancer in AWS, you will:

1. Launch and configure two EC2 instances with Amazon Linux 2, each hosting a simple web server with unique content.
2. Create and configure an Application Load Balancer to distribute incoming traffic between the two EC2 instances.
3. Verify the functionality of the Load Balancer by accessing the DNS name and observing traffic alternation between the two web servers.
4. Understand the importance of Load Balancers in ensuring high availability and fault tolerance for web applications.