# AWS Auto Scaling with Application Load Balancer (ALB)

## Step 1: Host a Simple PHP Application in EC2

1. Launch an EC2 instance (Amazon Linux 2 preferred).
2. Connect via SSH:
   ssh -i test.pem ec2-user@65.2.6.16
3. Install Apache and PHP:
   sudo yum install -y httpd php
4. Start and enable Apache:
   sudo systemctl start httpd
   sudo systemctl enable httpd
5. Create a PHP test page:
   echo "<?php phpinfo(); ?>" | sudo tee /var/www/html/index.php
6. Open HTTP (port 80) in the instance's Security Group.
7. Test by visiting http://65.2.6.16/ in a browser.

## Step 2: Configure an Application Load Balancer (ALB)

1. Go to EC2 → Load Balancers → Create Load Balancer → Application Load Balancer.
2. Select scheme: Internet-facing, and IP type: IPv4.
3. Choose at least two subnets from different Availability Zones.
4. Select a Security Group allowing inbound HTTP (port 80).
5. Create a Target Group:
   - Target type: Instance
   - Protocol: HTTP
   - Port: 80
   - Health check path: /
6. Register your running EC2 instance with the target group.
7. Finish ALB setup and note the DNS name (e.g., my-alb-12345.ap-south-

# 1.elb.amazonaws.com).



Load balancer: **my-app-lb**

EC2 > Load balancers > my-app-lb

## my-app-lb

Actions ▼

### ▼ Details

| Load balancer type | Status | VPC | Load balancer IP address type |
|---|---|---|---|
| Application | ⊖ Provisioning | vpc-0bc290fdc970e066e | IPv4 |

| Scheme | Hosted zone | Availability Zones | Date created |
|---|---|---|---|
| Internet-facing | ZP97RAFLXTNZK | subnet-09d8518b56d59e08a ap-south-1b (aps1-az3) | October 29, 2025, 21:32 (UTC+05:30) |
| | | subnet-0ab4c48c0e13d3293 ap-south-1a (aps1-az1) | |
| | | subnet-05991032cd04c7d2b ap-south-1-az2) | |

**Load balancer ARN**
arn:aws:elasticloadbalancing:ap-south-1:182399719928:loadbalancer/app/my-app-lb/da71acc35723d06d

✓ DNS name copied

my-app-lb-291126939.ap-south-1.elb.amazonaws.com (A Record)

Listeners and rules | Network mapping | Resource map | Security | Monitoring | Integrations | Attributes | Ca

---



EC2 > Load balancers > my-app-lb

-az2)

✓ DNS name copied

**Load balancer ARN**
arn:aws:elasticloadbalancing:ap-south-1:182399719928:loadbalancer/app/my-app-lb/da71acc35723d06d

my-app-lb-291126939.ap-south-1.elb.amazonaws.com (A Record)

Listeners and rules | Network mapping | Resource map | Security | Monitoring | Integrations | Attributes | Ca

### Listeners and rules (1) Info

Manage rules ▼ | Manage listener ▼ | Add listener

A listener checks for connection requests on its configured protocol and port. Traffic received by the listener is routed according to the default action and any additional rules.

Filter listeners

1

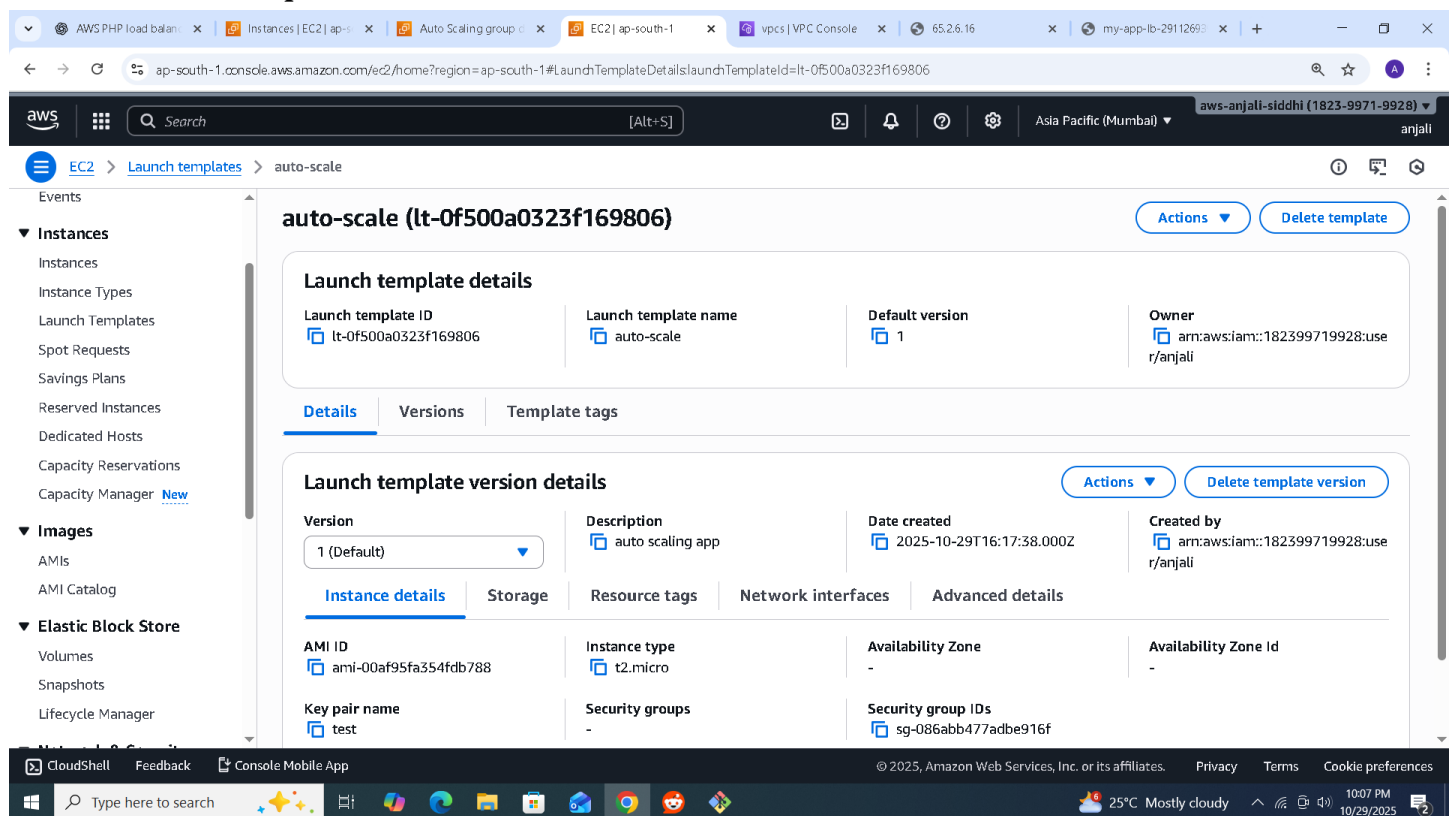| | Protocol:Port ▼ | Default action ▼ | Rules ▼ | ARN ▼ | Security policy ▼ |
|---|---|---|---|---|---|
| | HTTP:80 | • Forward to target group<br>new-tg : 1 (100%)<br>Target group stickiness: Off | 1 rule | ARN | Not applicable |

## Step 3: Create Launch Template for Auto Scaling
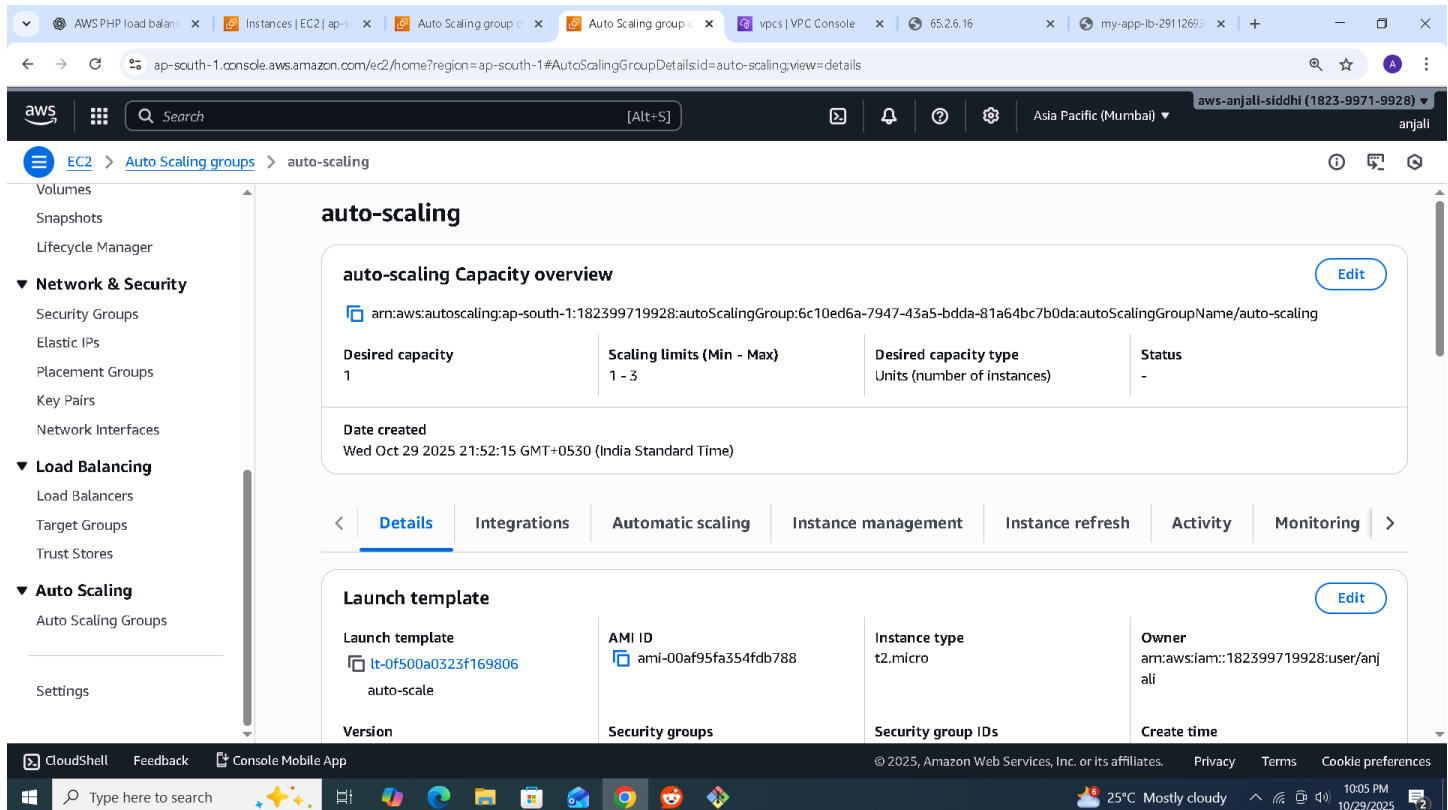
1. Go to EC2 → Launch Templates → Create Launch Template.
2. Configure:
   - Name: auto-scale
   - AMI: same as your instance
   - Instance type: t2.micro
   - Key pair: test
   - Security Group: same as your instance (allow port 80)
3. Leave subnet blank (ASG will handle multi-AZ placement).
4. Save the template.
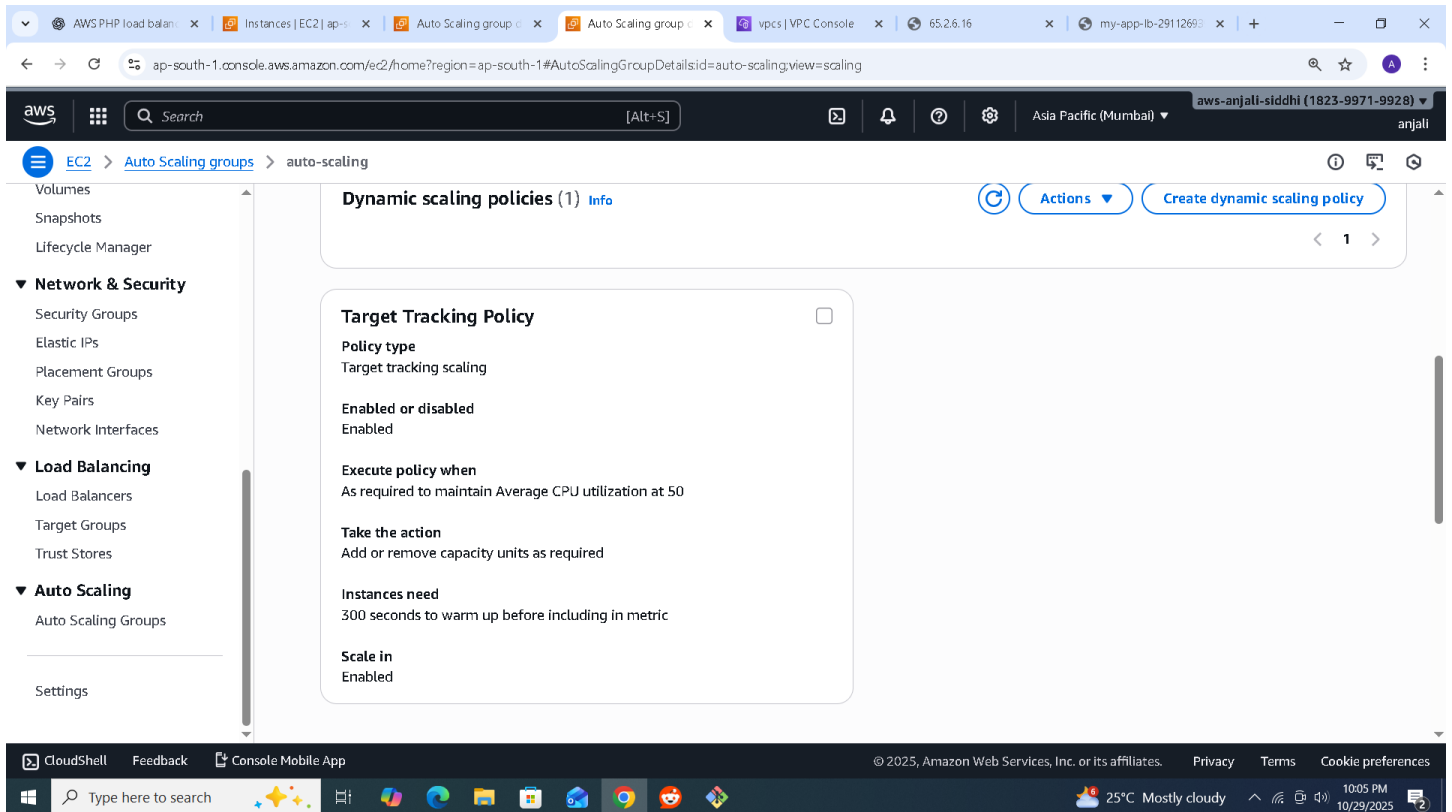


## Step 4: Create an Auto Scaling Group (ASG)

1. Go to EC2 → Auto Scaling Groups → Create Auto Scaling Group.
2. Select the Launch Template created earlier.
3. Set:
   - Desired capacity: 1
   - Minimum capacity: 1

- Maximum capacity: 3

4. Choose the same VPC and subnets as the ALB.

5. Attach the previously created Target Group.

6. Enable EC2 and ELB health checks (recommended).

7. Set Health check grace period: 300 seconds.

8. Skip or disable VPC Lattice integration unless required.

9. Finish creation.



## Step 5: Configure Scaling Policy (CPU-based)

1. Go to your Auto Scaling Group → Automatic Scaling tab.

2. Click "Add policy" → Choose "Target tracking scaling policy."

3. Configure as:
   - Policy name: CPU50-TargetTracking
   - Metric type: Average CPU utilization
   - Target value: 50%
   - Cooldown period: 300 seconds (default)

4. Save the policy.

# Step 6: Test Scaling Behavior

1. SSH into your instance:

   ssh -i test.pem ec2-user@<instance-public-ip>

2. Install the stress tool:

   sudo yum install -y stress

3. Run the load test:

   stress --cpu 2 --timeout 300

4. Monitor scaling:

   - In EC2 → Auto Scaling Group → Activity tab, watch for new instance launches.

   - In EC2 → Instances, new instances will appear.

   - In Load Balancer → Target Groups, all instances will show as healthy.

5. Once CPU usage drops below 50%, the ASG will terminate extra instances.

# Step 7: Verification and Monitoring

1. Verify in CloudWatch:
    - Go to CloudWatch → Metrics → EC2 → CPUUtilization.
    - Observe the spike during the stress test and the automatic scaling response.
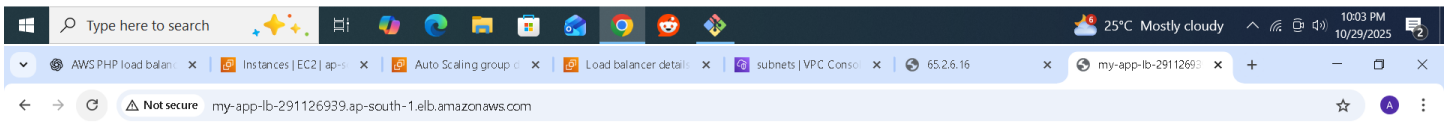2. Visit your ALB DNS URL and ensure the load is distributed across instances.

✓Successfully configured AWS Auto Scaling with an Application Load Balancer for a PHP web application.

**Hello from EC2 instance: ip-172-31-12-61.ap-south-1.compute.internal**

**Hello from Auto-Scaled instance: ip-172-31-37-121.ap-south-1.compute.internal**

# Hello from EC2 instance: ip-172-31-12-61.ap-south-1.compute.internal

```
Available Versions:

Version 2023.9.20251027:
    Run the following command to upgrade to 2023.9.20251027:

        dnf upgrade --releasever=2023.9.20251027

    Release notes:
        https://docs.aws.amazon.com/linux/al2023/release-notes/relnotes-2023.9.2025
1027.html

=======================================================================

Installed:
    stress-1.0.7-2.amzn2023.0.1.x86_64

Complete!
[ec2-user@ip-172-31-12-61 ~]$ stress --cpu 2 --timeout 300
stress: info: [28229] dispatching hogs: 2 cpu, 0 io, 0 vm, 0 hdd
^C
[ec2-user@ip-172-31-12-61 ~]$ stress --cpu 2 --timeout 300
stress: info: [28283] dispatching hogs: 2 cpu, 0 io, 0 vm, 0 hdd
```

EC2 > Auto Scaling groups > auto-scaling

## Activity notifications (0)

Actions ▼   Create notification

Filter notifications

☐  Send to                    ▼        On instance action                    ▼

No notifications are currently specified

Create notification

## Activity history (2)

Filter activity history

| Status | Description | Cause |
|--------|-------------|-------|
| ⊘ Successful | Launching a new EC2 instance: i-01b208faabdc5c1a0 | At 2025-10-29T16:22:30Z an instance was started in response to a difference between desired and actual capacity, increasing the capacity from 0 to 1. |

CloudShell   Feedback   Console Mobile App          © 2025, Amazon Web Services, Inc. or its affiliates.   Privacy   Terms   Cookie preferences