# Elastic Load Balancing & Auto Scaling

## 1) Create Launch Template

- Go to **EC2 → Launch Templates**
- Click **"Create launch template"**
- Configure the following:

**Template name: web-server-template**

**AMI: Amazon Linux 2 (ami-0abcdef1234567890)**

**Instance type: t2.micro**

**Key pair: Select existing or create new**

**Security groups: Create new with HTTP (80) and SSH (22)**

**User Data Script**

**Add this script in Advanced details → User data:**
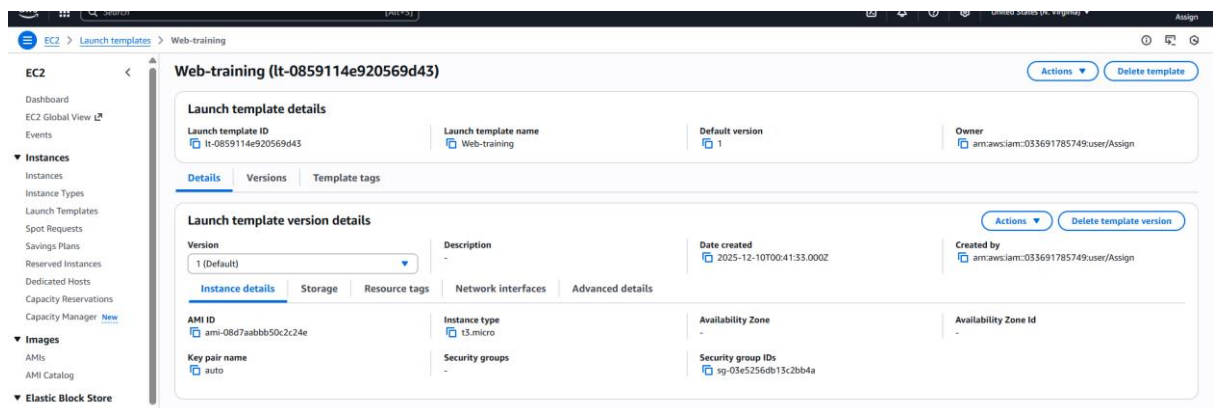
```
#!/bin/bash
yum update -y
yum install -y httpd
systemctl start httpd
systemctl enable httpd
echo "<h1>Web Server $(hostname -f)</h1>" >
/var/www/html/index.html
```



## 2) Create Auto Scaling Group

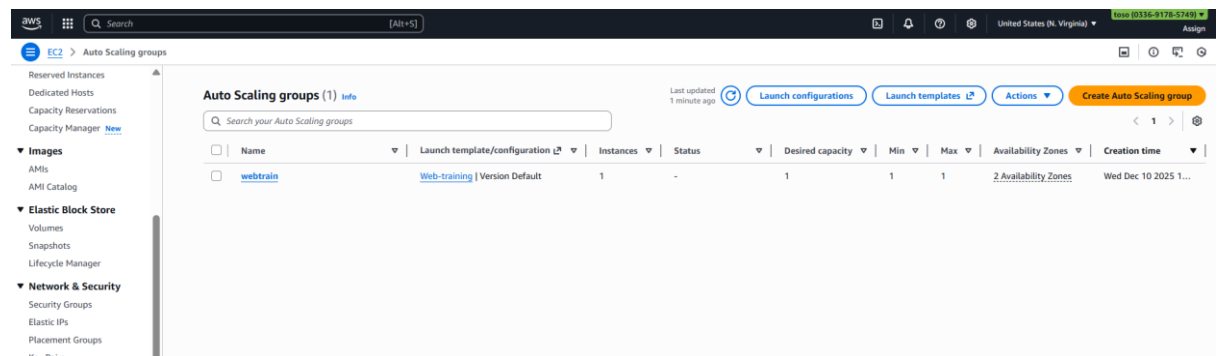1. Go to EC2 → Auto Scaling Groups

2. **Click "Create Auto Scaling group"**
3. **Select your launch template**
    **Auto Scaling group name: webtrain**
    **Launch template: web-server-template**
    **VPC: Default VPC**
    **Subnets: Select 2+ subnets in different AZs**



# 3) Create Application Load Balancer

- Go to **EC2 → Load Balancers**
- Click **"Create Load Balancer"**
- Choose **"Application Load Balancer"**

**Name: webtrain**
**Scheme: Internet-facing**
**IP address type: IPv4**
**VPC: Default VPC**
**Mappings: Select 2+ AZs with public subnets**

**Target Group Configuration**

- **Target type: Instances**
- **Protocol: HTTP**
- **Port: 80**
- **Health check path: /**
- **Health check interval: 30 seconds**

webtrain

**Details**
arn:aws:elasticloadbalancing:us-east-1:033691785749:targetgroup/webtrain/2be5ca203a909861

| Target type | Protocol : Port | Protocol version | VPC |
|---|---|---|---|
| Instance | HTTP: 80 | HTTP1 | vpc-06iaeced855313192b |
| **IP address type** | **Load balancer** | | |
| IPv4 | webtrain | | |

| 1 | ⊘ 1 | ⊗ 0 | ⊖ 0 | ⊘ 0 | ⊖ 0 |
|---|---|---|---|---|---|
| Total targets | Healthy | Unhealthy | Unused | Initial | Draining |
| | 0 Anomalous | | | | |

▸ **Distribution of targets by Availability Zone (AZ)**
Select values in this table to see corresponding filters applied to the Registered targets table below.

Targets | Monitoring | Health checks | Attributes | Tags

**Registered targets (1)** Info

Target groups route requests to individual registered targets using the protocol and port number specified. Health checks are performed on all registered targets according to the target group's health check settings. Anomaly detection is automatically applied to HTTP/HTTPS target groups with at least 3 healthy targets.

ⓘ Anomaly mitigation: Not applicable    Deregister    Register targets

| ☐ | Instance ID ▽ | Name ▽ | Port ▽ | Zone ▽ | Health status ▽ | Health status details | Administrative override ▽ | Override details ▽ | Launch ti |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | i-09c4745703d9f7545 | | 80 | us-east-1a (us... | ⊘ Healthy | - | ⊖ No override | No override is currently active ... | December |

---

aws    Q Search    [Alt+S]    United States (N. Virginia) ▾    Assign

EC2 > Load balancers > webtrain

**EC2**
Dashboard
EC2 Global View
Events
▾ Instances
  Instances
  Instance Types
  Launch Templates
  Spot Requests
  Savings Plans
  Reserved Instances
  Dedicated Hosts
  Capacity Reservations
  Capacity Manager New
▾ Images
  AMIs
  AMI Catalog
▾ Elastic Block Store
  Volumes
  Snapshots
  Lifecycle Manager
▾ Network & Security
  Security Groups
  Elastic IPs
  Placement Groups
  Key Pairs
  Network Interfaces
▾ Load Balancing
  Load Balancers
  Target Groups
  Trust Stores

ⓘ **Introducing ALB target optimizer**
Target optimizer lets you enforce a maximum number of requests per target using an ALB-provided agent, improving success rates, latency, and efficiency. Learn more    ✕

webtrain    Actions ▾

▾ **Details**

| Load balancer type | Status | VPC | Load balancer IP address type |
|---|---|---|---|
| Application | ⊘ Active | vpc-06aeced855313192b | IPv4 |
| **Scheme** | **Hosted zone** | **Availability Zones** | **Date created** |
| Internet-facing | Z355XDOTRQ7X7K | subnet-052577fe7dbc78a2e us-east-1e (use1-az3) | December 14, 2025, 22:34 (UTC+05:30) |
| | | subnet-027576533e131255a us-east-1a (use1-az1) | |
| **Load balancer ARN** | | **DNS name** Info | |
| arn:aws:elasticloadbalancing:us-east-1:033691785749:loadbalancer/app/webtrain/2da7b8336783Oc28 | | webtrain-825822773.us-east-1.elb.amazonaws.com (A Record) | |

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

**Listeners and rules (1)** Info
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.
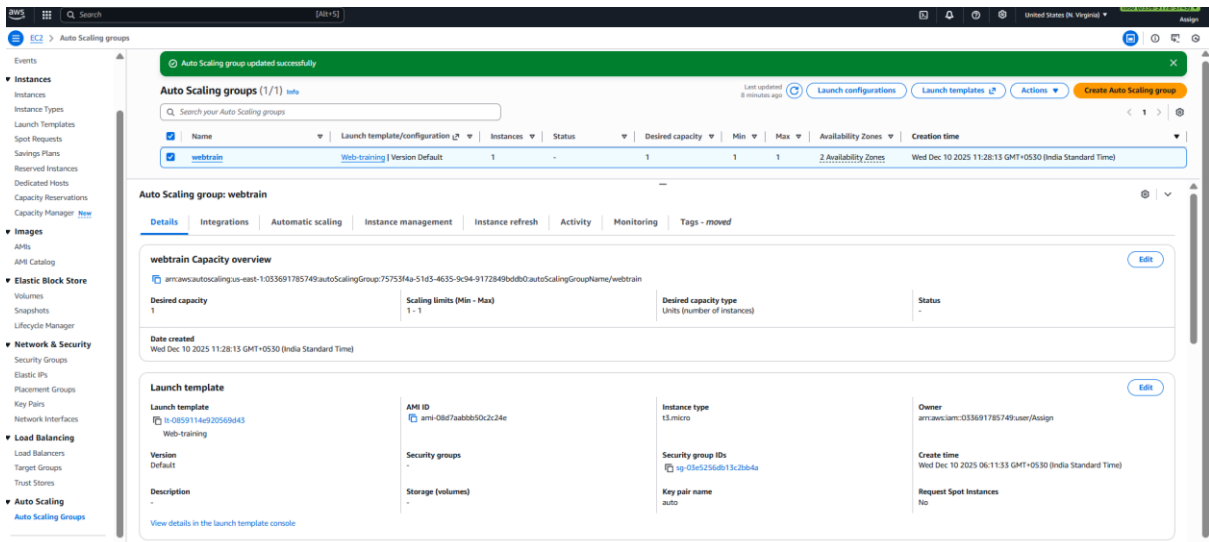
Manage rules ▾    Manage listener ▾    Add listener

| ☐ | Protocol:Port | Default action | Rules ▽ | ARN ▽ | Security policy ▽ | Default SSL/TLS certificate ▽ | mTLS ▽ | Trust store |
|---|---|---|---|---|---|---|---|---|
| ☐ | HTTP:80 | + Forward to target group<br>webtrain : 1 (100%)<br>Target group stickiness: Off | 1 rule | ARN | Not applicable | Not applicable | Not applicable | Not applica |

---

## 4) Attach Load Balancer to Auto Scaling Group

- Go back to **Auto Scaling Groups**
- Select your ASG: **web-server-asg**
- Go to **Details** tab
- Click **"Edit"** in Load balancing section

**Load balancing: Enable**

**Target groups: Select your target group**

**Health check type: ELB**

**Health check grace period: 300 seconds**

# 5) Configure Dynamic Scaling Policies
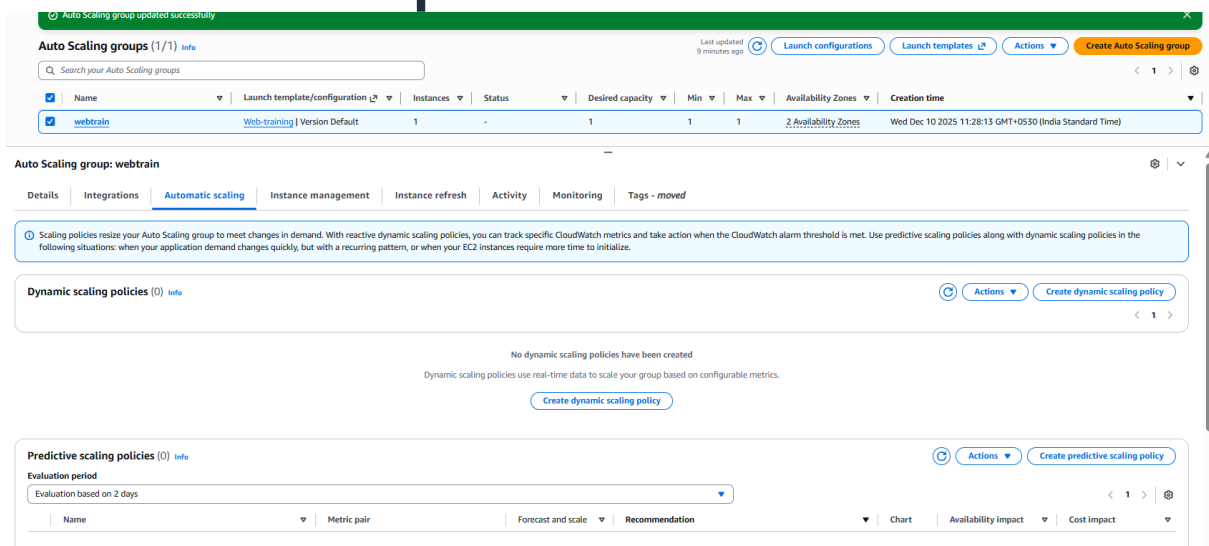
6) n ASG, go to **Automatic scaling** tab

7) Click **"Create dynamic scaling policy"**

## Policy type: Target tracking scaling

## Metric type: Average CPU Utilization

## Target value: 70%

## Instance warmup: 300 seconds

# Create dynamic scaling policy

**Policy type**

Target tracking scaling ▼

**Scaling policy name**

Target Tracking Policy

**Metric type** | Info

Monitored metric that determines if resource utilization is too low or high. If using EC2 metrics, consider enabling detailed monitoring for better scaling performance.

Average CPU utilization ▼

**Target value**

70

**Instance warmup** | Info

300 seconds

☐ Disable scale in to create only a scale-out policy

---



EC2 > Auto Scaling groups

⊘ Auto Scaling group updated successfully. ✕

⊘ Dynamic scaling policy created or edited successfully. ✕

**Auto Scaling groups** (1/1) Info

Launch configurations | Launch templates | Actions ▼ | **Create Auto Scaling group**

🔍 Search your Auto Scaling groups

| ☑ | Name ▼ | Launch template/configuration ▼ | Instances ▼ | Status ▼ | Desired capacity ▼ | Min ▼ | Max ▼ | Availability Zones ▼ | Creation time |
|---|---|---|---|---|---|---|---|---|---|

**Auto Scaling group: webtrain**

**Dynamic scaling policies** (1) Info

Actions ▼ | Create dynamic scaling policy

**Target Tracking Policy**

**Policy type**
Target tracking scaling

**Enabled or disabled**
Enabled

**Execute policy when**
As required to maintain Average CPU utilization at 70

**Take the action**
Add or remove capacity units as required

**Instances need**
300 seconds to warm up before including in metric

**Scale in**
Enabled

**Predictive scaling policies** (0) Info

Actions ▼ | Create predictive scaling policy

**Evaluation period**
Evaluation based on 2 days