# Scalable Web Application using **NLB + EC2 Auto Scaling**

## 1. Create Launch Template (EC2)

## 2. Create Target Group (Critical for NLB)

☰ EC2 > Target groups > Create target group ⬚ ⊘

Select the VPC with the instances that you want to include in the target group. Only VPCs that support the IP address type selected above are available in this list.

vpc-05073fe61b42c061b (default) ▼ ⟳ Create VPC ⬈
172.31.0.0/16

### Health checks
The associated load balancer periodically sends requests, per the settings below, to the registered targets to test their status.

**Health check protocol**

TCP ▼

▼ Advanced health check settings
Restore defaults

**Health check port**
The port the load balancer uses when performing health checks on targets. By default, the health check port is the same as the target group's traffic port. However, you can specify a different port as an override.
● Traffic port
○ Override

**Healthy threshold**
The number of consecutive health checks successes required before considering an unhealthy target healthy.

5

2-10

**Unhealthy threshold**
The number of consecutive health check failures required before considering a target unhealthy.

2

2-10

---

☰ EC2 > Target groups > mynlbtg ⬚ ⊘

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

▼ Auto Scaling
    Auto Scaling Groups

    Settings

# mynlbtg
Actions ▼

## Details
⎘ arn:aws:elasticloadbalancing:ap-southeast-1:368314293908:targetgroup/mynlbtg/7f808b2a2c8424ce

| **Target type** | **Protocol : Port** | **VPC** | **IP address type** |
|---|---|---|---|
| Instance | TCP : 80 | vpc-05073fe61b42c061b ⬈ | IPv4 |

**Load balancer**
ⓘ None associated

| **Total targets** | **Healthy** | **Unhealthy** | **Unused** | **Initial** | **Draining** |
|---|---|---|---|---|---|
| 0 | ✅ 0 | ⊗ 0 | ⊖ 0 | ⊘ 0 | ⊖ 0 |

**Targets** | Monitoring | Health checks | Attributes | Tags

### Registered targets (0)
⟳  Deregister  Register targets

🔍 Filter targets                                                    < 1 > ⚙

| ☐ | Instance ID ▽ | Name ▽ | Port ▽ | Zone ▽ | Health status ▽ | Health status details | Administrative ... ▽ | Override details ▽ | Launch time ▲ |
|---|---|---|---|---|---|---|---|---|---|

**No registered targets**
You have not registered targets to this group yet
Register targets

# 3. Create Network Load Balancer

☰ EC2 > Load balancers > Compare and select load balancer type ⓘ ⬚ ⊘

**Application Load Balancer** Info

**Network Load Balancer** Info

**Gateway Load Balancer** Info

Choose an Application Load Balancer when you need a flexible feature set for your applications with HTTP and HTTPS traffic. Operating at the request level, Application Load Balancers provide advanced routing and visibility features targeted at application architectures, including microservices and containers.
Create

Choose a Network Load Balancer when you need ultra-high performance, TLS offloading at scale, centralized certificate deployment, support for UDP, and static IP addresses for your applications. Operating at the connection level, Network Load Balancers are capable of handling millions of requests per second securely while maintaining ultra-low latencies.
Create

Choose a Gateway Load Balancer when you need to deploy and manage a fleet of third-party virtual appliances that support GENEVE. These appliances enable you to improve security, compliance, and policy controls.
Create

▶ **Classic Load Balancer** - *previous generation*

## Basic configuration

**Load balancer name**
Name must be unique within your AWS account and can't be changed after the load balancer is created.

mynlb

A maximum of 32 alphanumeric characters including hyphens are allowed, but the name must not begin or end with a hyphen.

**Scheme**
Scheme can't be changed after the load balancer is created.

○ **Internet-facing**
- Serves internet-facing traffic.
- Has public IP addresses.
- DNS name resolves to public IPs.
- Requires a public subnet.

○ **Internal**
- Serves internal traffic.
- Has private IP addresses.
- DNS name resolves to private IPs.

**Load balancer IP address type** | Info
Select the front-end IP address type to assign to the load balancer. The VPC and subnets mapped to this load balancer must include the selected IP address types.

● IPv4
Includes only IPv4 addresses.

○ Dualstack
Includes IPv4 and IPv6 addresses.

## Network mapping Info
The load balancer routes traffic to targets in the selected subnets, and in accordance with your IP address settings.

**VPC**
The load balancer will exist and scale within the selected VPC. The selected VPC is also where the load balancer targets must be hosted unless routing to on-premises targets or if using VPC peering. To confirm the VPC for your targets, view target groups ↗.

vpc-05073fe61b42c061b
172.31.0.0/16
(default) ▾    Create VPC ↗

## Network mapping Info
The load balancer routes traffic to targets in the selected subnets, and in accordance with your IP address settings.

**VPC**
The load balancer will exist and scale within the selected VPC. The selected VPC is also where the load balancer targets must be hosted unless routing to on-premises targets or if using VPC peering. To confirm the VPC for your targets, view target groups ↗.

vpc-05073fe61b42c061b
172.31.0.0/16
(default) ▾    Create VPC ↗

**Availability Zones and subnets**
Select one or more Availability Zones and corresponding subnets. Enabling multiple Availability Zones increases the fault tolerance of your applications. The load balancer routes traffic to targets in the selected Availability Zones only. Availability Zones that are not supported by the load balancer or the VPC are not available for selection.

☑ **ap-southeast-1a (apse1-az1)**
Subnet
Only CIDR blocks corresponding to the load balancer IP address type are used. At least 8 available IP addresses are required for your load balancer to scale efficiently.

subnet-04bf478a0e7c6ded7
IPv4 subnet CIDR: 172.31.16.0/20 ▾

IPv4 address
The front-end IPv4 address of the load balancer in the selected Availability Zone.

● Assigned by AWS    ○ Use an Elastic IP address

☑ **ap-southeast-1b (apse1-az2)**
Subnet
Only CIDR blocks corresponding to the load balancer IP address type are used. At least 8 available IP addresses are required for your load balancer to scale efficiently.

subnet-06db806dacbeeb4d7
IPv4 subnet CIDR: 172.31.32.0/20 ▾

IPv4 address
The front-end IPv4 address of the load balancer in the selected Availability Zone.

● Assigned by AWS    ○ Use an Elastic IP address

☐ ap-southeast-1c (apse1-az3)

## Security groups Info
A security group is a set of firewall rules that control the traffic to your load balancer. Select an existing security group, or you can create a new security group ↗.

**Security groups – recommended**
Security groups support on Network Load Balancers can only be enabled at creation by including at least one security group. You can change security groups after creation. The security groups for your load balancer must allow it to communicate with registered targets on both the listener port and the health check port. For PrivateLink Network Load Balancers, security group rules are enforced on PrivateLink traffic; however, you can turn off inbound rule evaluation after creation within the load balancer's Security tab or using the API.

Select up to 5 security groups ▾

mylbsecuritygroup ✕
sg-0e2998ce1115d491e1    VPC: vpc-05073fe61b42c061b

## Listeners and routing Info
A listener is a process that checks for connection requests using the port and protocol you configure. The rules that you define for a listener determine how the load balancer routes requests to its registered targets.

▼ Listener TCP:80                                                    Remove

**Protocol**              **Port**
TCP ▾                     80
                         1-65535

**Forward to target group** | Info
Choose a target group and specify routing weight or create target group ↗.

**Target group**                          **Weight**    **Percent**
mynlbtg                          TCP ▾    1             100%
Target type: Instance, IPv4 | Target stickiness: Off   0-999

[Add target group]
You can add up to 4 more target groups.

**Target group stickiness** | Info
Enables the load balancer to bind a user's session to a specific target group. If you want to bind a user's session to a specific target, turn on the Target Group attribute Stickiness.
☐ Turn on target group stickiness

**Weighted routing evaluation – recommended**
Your Network Load Balancer's ability to distribute requests according to assigned target group weights can be hindered if your target groups don't follow best practices.
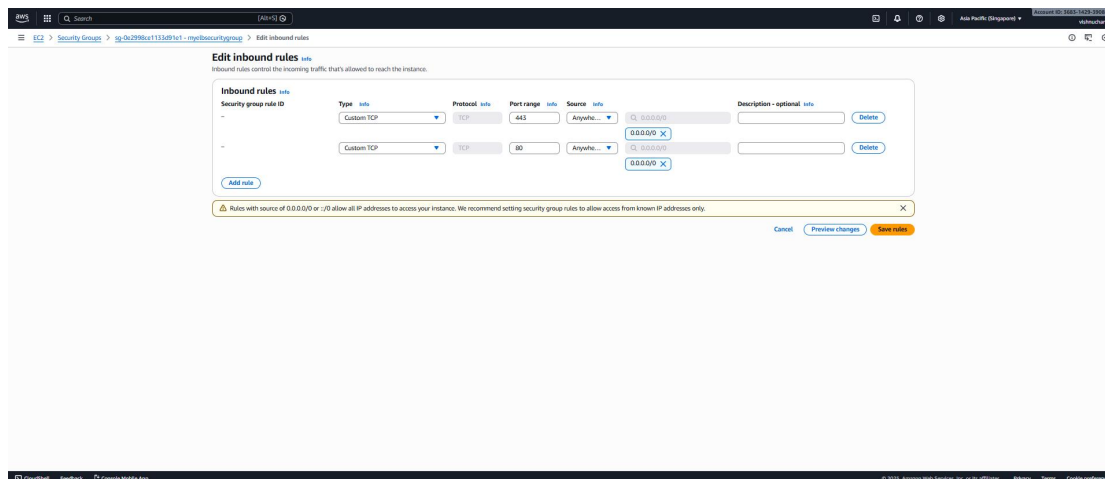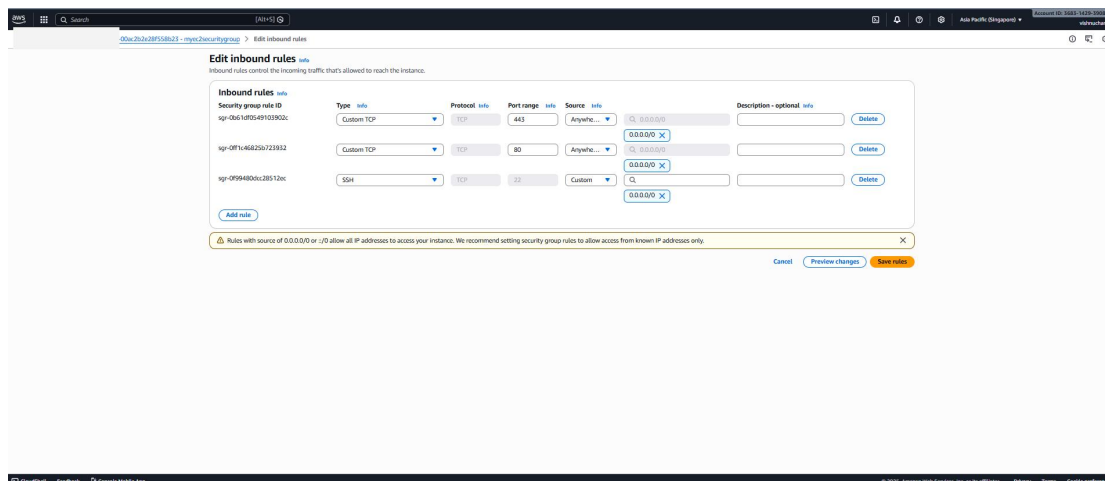[💡 Evaluate]

**Listener tags – optional**
Consider adding tags to your listener. Tags enable you to categorize your AWS resources so you can more easily manage them.

NLB Security group(Inbound rules):-



EC2 Security group(Inbound rules):-

# 4. Create Auto Scaling Group



## Choose launch template or configuration

Specify a launch template that contains settings common to all EC2 instances that are launched by this Auto Scaling group. If you currently use launch configurations, you might consider migrating to launch templates.

### Name

**Auto Scaling group name**
Enter a name to identify the group.

mynlbasg

Must be unique to this account in the current region and no more than 255 characters.

### Launch template Info

Switch to launch configuration

**Launch template**
Choose a launch template that contains the instance-level settings, such as the Amazon Machine Image (AMI), instance type, key pair, and security groups.

myec2lt

Create a launch template

**Version**

Latest (1)

Create a launch template version

| Description | Launch template | Instance type |
|---|---|---|
| - | myec2lt | t2.micro |
| | lt-0ce1857b268ee0ec6 | |
| **AMI ID** | **Security groups** | **Request Spot Instances** |
| ami-05f071c65e32875a8 | - | No |
| **Key pair name** | **Security group IDs** | |
| mydebian13 | sg-00ac2b2e28f558b23 | |

### Additional details

| Storage (volumes) | Date created |
|---|---|
| - | Mon Dec 22 2025 23:16:49 GMT+0530 (India Standard Time) |

Cancel    Next



## Choose instance launch options

Choose the VPC network environment that your instances are launched into, and customize the instance types and purchase options.

### Instance type requirements Info

Override launch template

You can keep the same instance attributes or instance type from your launch template, or you can choose to override the launch template by specifying different instance attributes or manually adding instance types.

| Launch template | Version | Description |
|---|---|---|
| myec2lt | Latest | |
| lt-0ce1857b268ee0ec6 | | |

**Instance type**
t2.micro

### Network Info

For most applications, you can use multiple Availability Zones and let EC2 Auto Scaling balance your instances across the zones. The default VPC and default subnets are suitable for getting started quickly.

**VPC**
Choose the VPC that defines the virtual network for your Auto Scaling group.

vpc-05073fe61b42c061b
172.31.0.0/16   Default

Create a VPC

**Availability Zones and subnets**
Define which Availability Zones and subnets your Auto Scaling group can use in the chosen VPC.

Select Availability Zones and subnets

apse1-az1 (ap-southeast-1a) | subnet-04bf478a0e7c6ded7  X
172.31.16.0/20   Default

apse1-az2 (ap-southeast-1b) | subnet-06db806dacbeeb4d7  X
172.31.32.0/20   Default

Create a subnet

**Availability Zone distribution - new**
Auto Scaling automatically balances instances across zones. If launch failures occur in a zone, select a strategy.

- ○ Balanced best effort
  If launches fail in one Availability Zone, Auto Scaling will attempt to launch in another healthy Availability Zone.
- ○ Balanced only
  If launches fail in one Availability Zone, Auto Scaling will continue to attempt to launch in the unhealthy Availability Zone to preserve balanced distribution.

## Integrate with other services - *optional* Info

Use a load balancer to distribute network traffic across multiple servers. Enable service-to-service communications with VPC Lattice. Shift resources away from impaired Availability Zones with zonal shift. You can also customize health check replacements and monitoring.

### Load balancing Info

Use the options below to attach your Auto Scaling group to an existing load balancer, or to a new load balancer that you define.

**Select Load balancing options**

- ○ No load balancer
  Traffic to your Auto Scaling group will not be fronted by a load balancer.
- ● Attach to an existing load balancer
  Choose from your existing load balancers.
- ○ Attach to a new load balancer
  Quickly create a basic load balancer to attach to your Auto Scaling group.

### Attach to an existing load balancer

**Select the load balancers to attach**

- ● Choose from your load balancer target groups
  This option allows you to attach Application, Network, or Gateway Load Balancers.
- ○ Choose from Classic Load Balancers

**Existing load balancer target groups**
Only instance target groups that belong to the same VPC as your Auto Scaling group are available for selection.

Select target groups

mynlbtg | TCP  X
Network Load Balancer: mynlb

### VPC Lattice integration options Info

To improve networking capabilities and scalability, integrate your Auto Scaling group with VPC Lattice. VPC Lattice facilitates communications between AWS services and helps you connect and manage your applications across compute services in AWS.

**Select VPC Lattice service to attach**

- ● No VPC Lattice service
  VPC Lattice will not manage your Auto Scaling group's network access and connectivity with other services.
- ○ Attach to VPC Lattice service
  Incoming requests associated with specified VPC Lattice target groups will be routed to your Auto Scaling group.

Create new VPC Lattice service

### Application Recovery Controller (ARC) zonal shift - *new* Info

During an Availability Zone impairment, target instance launches towards other healthy Availability Zones.

EC2 > Auto Scaling groups > Create Auto Scaling group

≡ EC2 > Auto Scaling groups > Create Auto Scaling group

**Step 1** Choose launch template or configuration

**Step 2** Choose instance launch options

**Step 3 - optional** Integrate with other services

**Step 4 - optional** Configure group size and scaling

**Step 5 - optional** Add notifications

**Step 6 - optional** Add tags

**Step 7** Review

### Configure group size and scaling - *optional* Info
Define your group's desired capacity and scaling limits. You can optionally add automatic scaling to adjust the size of your group.

#### Group size Info
Set the initial size of the Auto Scaling group. After creating the group, you can change its size to meet demand, either manually or by using automatic scaling.

**Desired capacity type**
Choose the unit of measurement for the desired capacity value. vCPUs and Memory(GiB) are only supported for mixed instances groups configured with a set of instance attributes.

Units (number of instances)

**Desired capacity**
Specify your group size.

2

#### Scaling Info
You can resize your Auto Scaling group manually or automatically to meet changes in demand.

**Scaling limits**
Set limits on how much your desired capacity can be increased or decreased.

**Min desired capacity** | **Max desired capacity**
1 | 3
Equal or less than desired capacity | Equal or greater than desired capacity

**Automatic scaling - *optional***
**Choose whether to use a target tracking policy** Info
You can set up other metric-based scaling policies and scheduled scaling after creating your Auto Scaling group.

○ No scaling policies
Your Auto Scaling group will remain at its initial size and will not dynamically resize to meet demand.

○ Target tracking scaling policy
Choose a CloudWatch metric and target value and let the scaling policy adjust the desired capacity in proportion to the metric's value.

#### Instance maintenance policy Info
Control your Auto Scaling group's availability during instance replacement events. This includes health checks, instance refreshes, maximum instance lifetime features and events that happen automatically to keep your group balanced, called rebalancing events.

**Choose a replacement behavior depending on your availability requirements**

**Mixed behavior** - No policy — For rebalancing events, new instances will
**Prioritize availability** - Launch before terminating — Launch new instances and wait for them to be
**Control costs** - Terminate and launch — Terminate and launch instances at the same
**Flexible** - Custom behavior — Set custom values for the minimum and

---

### Auto Scaling groups (1/1) Info
Last updated less than a minute ago | Launch configurations | Launch templates ⊘ | Actions ▾ | **Create Auto Scaling group**

Q Search your Auto Scaling groups

| Name | Launch template/configuration ⊘ | Instances | Status | Desired capacity | Min | Max | Availability Zones | Creation time |
|---|---|---|---|---|---|---|---|---|
| **mynlbasg** | myec2lt | Version Latest | 2 | - | 2 | 1 | 3 | 2 Availability Zones | Mon Dec 22 2025 23:41:16 GMT+0530 (India Standard Time) |

**Auto Scaling group: mynlbasg**

**Details** | Integrations | Automatic scaling | Instance management | Instance refresh | Activity | Monitoring | Tags – *moved*

#### mynlbasg Capacity overview
arn:aws:autoscaling:ap-southeast-1:368314293908:autoScalingGroup:a44a9b63-772a-45c0-a2cb-c01706a2bfac:autoScalingGroupName/mynlbasg

**Desired capacity** | **Scaling limits** | **Desired capacity type** | **Status**
2 | 1 - 3 | Units (number of instances) | -

**Date created**
Mon Dec 22 2025 23:41:16 GMT+0530 (India Standard Time)

#### Launch template

**Launch template**
lt-0ce18576268ae0ec6
myec2lt

**Version**
Latest

**Description**
-

**AMI ID**
ami-05f071c65e32875a8

**Security group**
-

**Storage (volumes)**
-

**Instance type**
t2.micro

**Security group IDs**
sg-00ac2b2e28f558b23

**Key pair name**
mydebian13

**Owner**
arn:aws:iam::368314293908:root

**Create time**
Mon Dec 22 2025 23:16:49 GMT+0530 (India Standard Time)

**Request Spot Instances**
No

View details in the launch template console

---

### Instances (2) Info
⟳ | Connect | Instance state ▾ | Actions ▾ | **Launch instances** ▾

Q Find Instance by attribute or tag (case-sensitive) | All states ▾

running ✕ | Clear filters

| Name | Instance ID | Instance state | Instance type | Status check | Alarm status | Availability Zone | Public IPv4 DNS | Public IPv4 ... | Elastic IP | IPv6 IPs | Monitoring | Security group name | Key name |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | i-095041782b6f6c89a1 | ⊘ Running | t2.micro | ⊘ Initializing | View alarms + | ap-southeast-1b | ec2-13-212-144-114.ap... | 13.212.144.114 | – | – | disabled | myec2securitygroup | mydebian13 |
| | i-0caa5b5a29ac2189d | ⊘ Running | t2.micro | ⊘ Initializing | View alarms + | ap-southeast-1a | ec2-3-1-222-65.ap-sout... | 3.1.222.65 | – | – | disabled | myec2securitygroup | mydebian13 |

**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

**Auto Scaling**
- Auto Scaling Groups

Settings

**Select an instance**