# 46°
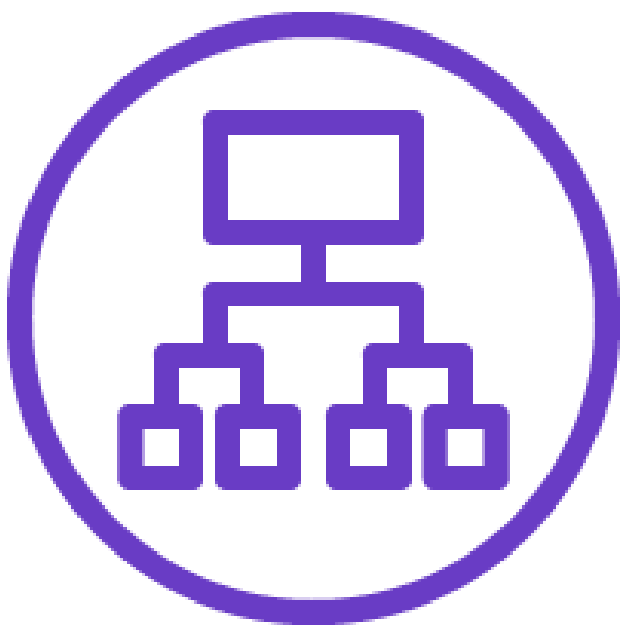
## Lab - AWS re/Start
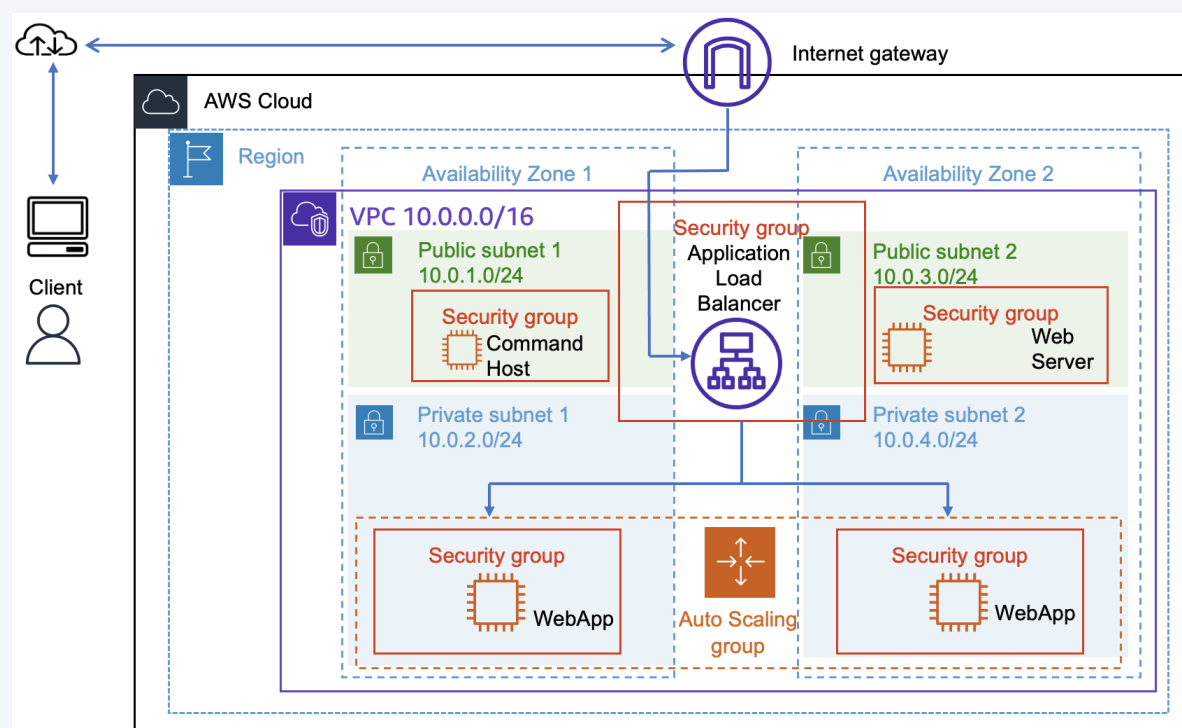## Uso de escalado automático en AWS (Linux)



Amazon EC2 Auto Scaling

# Interactuando con Arquitecturas

Los objetivos son:

- Crear una instancia EC2 mediante un comando de la CLI de AWS.
- Crear una nueva AMI mediante la CLI de AWS.
- Crear una plantilla de lanzamiento de Amazon EC2.
- Cree una configuración de lanzamiento de Amazon EC2 Auto Scaling.
- Configure las políticas de escalado y cree un grupo de Auto Scaling para aumentar y reducir el número de servidores en función de una carga variable.

# Tarea 01

Empezamos creando la AMI de la instancia *Command Host* desde la CLI de AWS



Primero, crearemos una nueva instancia, a continuación se muestra la User Data para esta instancia a crear
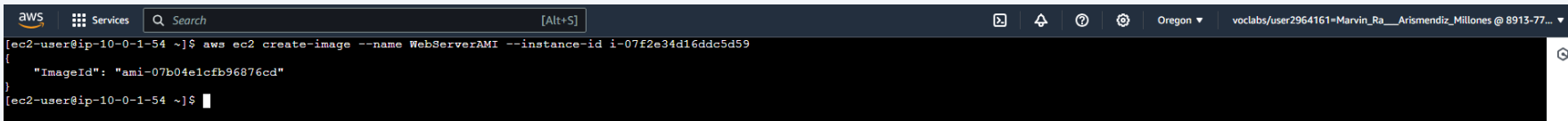


Lanzamos la nueva instancia

# Tarea 01

Ahora, veamos esta nueva instancia:



Luego, procedemos a crear una AMI de esta instancia



Después de ello, procedemos a crear un ELB:

# Tarea 01

## Más de la configuración del ALB (Application Load Balancer)



### Specify group details

Your load balancer routes requests to the targets in a target group and performs health checks on the targets.

**Basic configuration**
Settings in this section can't be changed after the target group is created.

Choose a target type

- **Instances**
  - Supports load balancing to instances within a specific VPC.
  - Facilitates the use of Amazon EC2 Auto Scaling to manage and scale your EC2 capacity.

- **IP addresses**
  - Supports load balancing to VPC and on-premises resources.
  - Facilitates routing to multiple IP addresses and network interfaces on the same instance.
  - Offers flexibility with microservice based architectures, simplifying inter-application communication.
  - Supports IPv6 targets, enabling end-to-end IPv6 communication, and IPv4-to-IPv6 NAT.

- **Lambda function**
  - Facilitates routing to a single Lambda function.
  - Accessible to Application Load Balancers only.

- **Application Load Balancer**
  - Offers the flexibility for a Network Load Balancer to accept and route TCP requests within a specific VPC.
  - Facilitates using static IP addresses and PrivateLink with an Application Load Balancer.

Target group name
webserver-app

### Listeners and routing

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 HTTP:80** — Remove

Protocol: HTTP
Port: 80 (1-65535)
Default action: Forward to webserver-app / Target type: Instance, IPv4 — HTTP

Create target group

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

Add listener tag
You can add up to 50 more tags.

Add listener

**Health checks**

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

Health check protocol
HTTP

Health check path
Use the default path of "/" to perform health checks on the root, or specify a custom path if preferred.
/index.php
Up to 1024 characters allowed.

Advanced health check settings

## Luego, procedemos a crear la plantilla de lanzamiento



### Summary

**Software Image (AMI)**
WebServerAMI
ami-07b04e1cfb96876cd

**Virtual server type (instance type)**
t3.micro

**Firewall (security group)**
HTTPAccess

**Storage (volumes)**
1 volume(s) - 8 GiB

Cancel   **Create launch template**

SWIPE

## Luego creamos el ASG (Auto Scaling Group)

### 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 |
|---|---|---|
| web-app-launch-template 🔗 | Default | web-app-launch-template |
| lt-01bc43cbcc2d4a863 | | |

Instance type
t3.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-07c536590dcad87be (Lab VPC)
10.0.0.0/16

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

us-west-2a | subnet-0cc436d2de1674e74 (Private Subnet 1)  ✕
10.0.2.0/24

us-west-2b | subnet-00dfef2f5052e3e74 (Private Subnet 2)  ✕
10.0.4.0/24

Create a subnet 🔗

### Health checks

Health checks increase availability by replacing unhealthy instances. When you use multiple health checks, all are evaluated, and if at least one fails, instance replacement occurs.

**EC2 health checks**
ⓘ Always enabled

Additional health check types - *optional*   Info
☑ Turn on Elastic Load Balancing health checks  `Recommended`
Elastic Load Balancing monitors whether instances are available to handle requests. When it reports an unhealthy instance, EC2 Auto Scaling can replace it on its next periodic check.

ⓘ EC2 Auto Scaling will start to detect and act on health checks performed by Elastic Load Balancing. To avoid unexpected terminations, first verify the settings of these health checks in the Load Balancer console 🔗   ✕

☐ Turn on VPC Lattice health checks
VPC Lattice can monitor whether instances are available to handle requests. If it considers a target as failed a health check, EC2 Auto Scaling replaces it after its next periodic check.

**Health check grace period**   Info
This time period delays the first health check until your instances finish initializing. It doesn't prevent an instance from terminating when placed into a non-running state.

300   seconds

### 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.

| ○ **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 that you want to attach to your Auto Scaling group.

| ● **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  ▼

webserver-app | HTTP   ✕
Application Load Balancer: WebServerELB

### 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 |
|---|---|
| 2 | 4 |
| 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. |
|---|---|

**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**

50

**Instance warmup**   Info

300   seconds

# Tarea 01

## Y lo probamos, funciona correctamente