# AWS Elastic Load Balancing tips and tricks: from basic to pro

*12 August 2022 - 8 min. read*

| *Advanced Networking* | *Application Load Balancer (ALB)* | *AWS Elastic Load Balancer (ELB)* |

| *Network Load Balancer (NLB)* |

An old Italian commercial claimed: "Two is better than one."

When it comes to **application availability**, having multiple instances available makes a business more resilient: this can help you achieve **fault tolerance** for applications and infrastructures. You know: "Everything can fail, all the time".

**Load balancing** is the fundamental building block that can make us achieve better uptime and application availability: redistributing traffic on different instances in auto scaling and checking their healthiness isn't always as easy as it seems.

In my past career as a system administrator, I always struggled to find a resilient, redundant, and elastic solution. After a lot of work and automation, I could have a decent night of sleep!

Using managed services, as always, can help us reduce the amount of work required to achieve our goals.

When it comes to load-balancing, AWS gives us a lot of flexibility: under the "Elastic Load Balancing (ELB)" umbrella, many options are available.

This article will explain ELB fundamentals and deep dive into not-so-common use cases.

# Basic Key Concepts

Three managed load balancer types are available: **Application**, **Network**, and **Gateway**.

In this article, we'll focus on **Application Load Balancers** and **Network Load Balancers**, while we will not consider the Classic Load balancer, as it is a mix of application and network load balancers with some missing features (check this page for comparison), and Gateway Load Balancers. We will write about them in a dedicated article.

Every type of ELB can span multiple Availability Zones and relies on three main components:

- **listener**
- **target group**
- **health checks**

ELB can be public (internet facing) or private (so that only private routed resources can access them).

A **listener** is a small portion of the ELB configuration that defines the entry point of the traffic. If our application uses the HTTP protocol on port 8080, the listener configuration will route traffic to the same protocol and port.

A **target group** is the set of computational resources that handle application traffic distributed by the load balancer.

A target group can use EC2 instances, ECS containers, IP addresses, lambda functions, or even another application load balancer!

Note that not every ELB type can use all the targets: for example, only an Application Load Balancer can have lambda functions as a target.

A **health check** is a test used by the target group to determine if the target is healthy, so failing compute resources are excluded and will not receive any traffic. For example, if our application receives TCP traffic on port 31337, the health check verifies that the service is listening.

# Application Load Balancer

The application load balancer (ALB) operates at ISO/OSI level 7, thus the name.

It's the most commonly used because it balances HTTP and HTTPS traffic and can **perform redirection, authentication, and SSL offloading** using Amazon Certificate Manager (ACM) certificates.

## Network Load Balancer

The network load balancer (NLB) operates at ISO/OSI level 4 (network or connection level). It can handle any traffic based on TCP or UDP without specialization for a specific application protocol. TLS offloading is supported, offering static IP addresses to use (for both internal and external load balancers).

## Gateway Load Balancer

Gateway load balancers are the last addition to the ELB family. It routes level 3 traffic using GENEVE encapsulation, enabling usage of third-party and custom network and security appliances.

Suppose you want to deploy a custom IDS solution to analyze traffic in transit: Gateway load balancing helps achieve high availability, using multiple instances in different availability zones instead of a single one.

So far so good. And straightforward. Let's now dive deep into **five (+1)** less common **configuration scenarios that can make your day**: how to configure **load balancers like a pro!**

## 1. Use ALB to secure WordPress

Refactoring applications and adapting them to use the Cloud efficiently isn't always an option when you are in a rush.

WordPress is one of the most common examples of an application that isn't easy to migrate; if you have the time, please consider containerization.

When migrating a WordPress installation, you need to keep the same DNS name and the connection on HTTPS. The quick and dirty solution could be to use a public amazon EC2 instance, but exposing a WordPress installation without a WAF can be risky.

Placing an Application Load Balancer in front of the WordPress instance can help **reduce administration overhead** and **improve security**, even if you aren't balancing traffic between different instances:

- The EC2 instance will not be publicly accessible

- You can place a WAF and use the WordPress specialized ruleset

- Amazon Certificate Manager will give us public certificates and renew them

ALB will perform HTTP redirection to the HTTPS listener in this use case, and LetsEncrypt CertBot will issue SSL certificates on the instance.

Associating a Web Application Firewall rule to the load balancer is relatively easy; let's see how to configure our ALB to allow certificate renewal with LetsEncrypt.

After obtaining CertBot (https://certbot.eff.org/instructions), you can issue free certificates for your website, but the verification process requires access to a path on your web server using HTTP

```
(/.well-known-acme-challenge/)
```

To achieve this, we'll need to define a target group that uses port 80 and then insert a rule that will use that target group for the acme-challenge path.

First, go to the management console and create a new target group:

Then select your instance and click "Include as pending below."



After creating the target group, click on "**View/edit rules**" on the HTTP listener to add a new rule:

Remember that this rule should have a higher priority than the default action that redirects HTTP to HTTPS

## 2. Simple and quick maintenance page for your site

It can sound obvious, but sometimes we forget it. Setting up a rule with a fixed response with an HTTP status code of 503 will enable us to serve a maintenance page quickly for a given website:



## 3. Host-based HTTP routing to reduce bot scanning

When using the default action to target our web application, all traffic will be forwarded, even bot scanning for vulnerabilities targeting every public IP address with an open port. This behavior can make our logs "noisy".

To avoid this, change the default action to a fixed response and add a more specific rule for our host:

## 4. Implement custom logic for health checks in ECS tasks balanced using an NLB

For this example, we'll use a Minecraft server, even if it's not a traditional workload. Minecraft runs on port 25565 using its custom protocol.

We can run it in an ECS Fargate Cluster, configuring a service to keep the server running with "auto-healing" and setting the desired count of 1.

With a network load balancer, we can have static IP addresses to connect to, so the healthy container will always receive traffic without worrying when its internal IP changes.

We will configure the target group to check if port 25565 is reachable, but we want to take a step further: we want to be sure that the server will respond to a command.

This bash script can check the number of connected users; unfortunately, we can't configure our target group to execute it.

```
echo -e '\xfe' | nc -w 5 127.0.0.1 25565 | awk -F'\xa7' '$2 {printf
  "users: %s",$2;}'| grep users
```

Lucky for us, ECS Fargate will let us define a custom health check that will run in our task instance.

Simply add to your container definition the health check command so that a failure will stop the container and a new instance will take its place. Remember that the health check script has to return 0 to be considered OK.
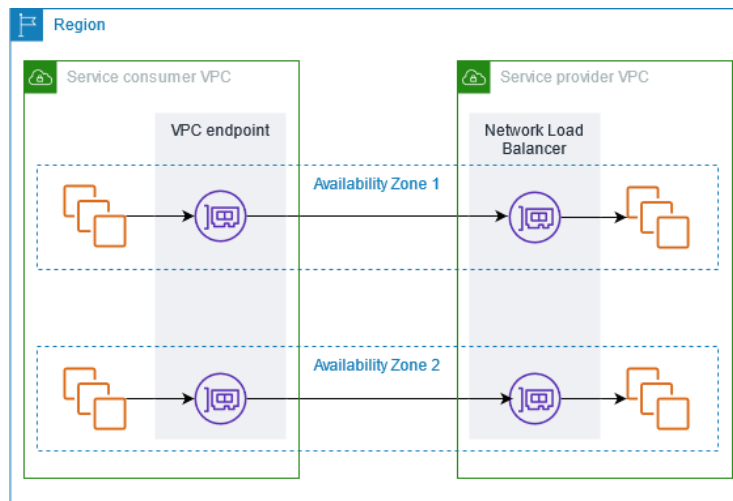
## 5. Share services to other AWS Accounts without crossing the Internet

AWS PrivateLink is a technology that allows you to offer services to customers and partners, taking advantage of the AWS network infrastructure without crossing the internet.

When you share a service with other AWS accounts, you become the service provider; the other AWS accounts become Service consumers.

It works by linking a Network Loadbalancer's network interfaces with other accounts, placing a local interface in the consumer VPC.



To share services, you first need to create an **Endpoint Service**, go to the VPC console, select Endpoint Services, and create a new one by selecting **"Network"** as the load balancer type.

Once the service endpoint is available, please note the Service Name in the console.



**To use the endpoint service as a consumer,** go to the "Endpoints" section in the VPC, select "Create a new endpoint," and then select "**other endpoint services**." Paste the name you previously noted and then click "**Verify Service**." Once the service name is verified, you can select the desired VPC, subnet, and security group to use to deploy the endpoint.

Remember that you need to accept the connection if you enable the "Acceptance Required" checkbox while creating the endpoint service.



After accepting the connection, the endpoint will be ready, and the consumer can access the service without traversing the Internet or using a VPN connection.

## Bonus tip

As the last tip, remember that you can **use an application load balancer as a target for a network load balancer**, so you can enable consumers to use an internal service without having to reconfigure the infrastructure.

## To conclude

The AWS ELB services offer great flexibility in terms of choice and configuration. Mixing and matching managed services can help reduce operational complexity and lower costs. We only scratched the surface, and the possibilities are endless.

Did you find other usage scenarios for elastic load balancing? Let us know in the comments below!

See you again in 14 days on **Proud2beCloud**!

## About Proud2beCloud

Proud2beCloud is a blog by beSharp, an Italian APN Premier Consulting Partner expert in designing, implementing, and managing complex Cloud infrastructures and advanced services on AWS. Before being writers, we are Cloud Experts working daily with AWS services since 2007. We are hungry readers, innovative builders, and gem-seekers. On Proud2beCloud, we regularly share our best AWS pro tips, configuration insights, in-depth news, tips&tricks, how-tos, and many other resources. Take part in the discussion!



### Damiano Giorgi

Ex on-prem systems engineer, lazy and prone to automating boring tasks. In constant search of technological innovations and new exciting things to experience. And that's why I love Cloud Computing! At this moment, the only "hardware" I regularly dedicate myself to is that my bass; if you can't find me in the office or in the band room try at the pub or at some airport, then!