# Auto Scaling

## Introduction to Auto Scaling



*Auto Scaling*

In the previous lesson, we saw how the Load Balancer manages traffic between our existing resources, in other words, instances.

But, what if we want to automate the increasing or decreasing of my existing virtual machines when demand changes over time?

We call **Auto Scaling** that makes this process in AWS for you. Amazon EC2 Auto Scaling is a component that helps you ensure that you have the correct number of Amazon EC2 instances available to handle the load for your application.

Thanks to Auto Scaling, we can add or remove a new virtual machine when it is necessary. So, you can expand your system to respond to that demand according to the size of the request you receive, and automatically collapse when the demand decreases.

With Auto Scaling, you can also define certain rules and create new resources in the system according to these rules.

You can set rules and give direction which depends on some sort of condition. For example, you may say;

- If the virtual machines use more than 90% CPU resources for 5 minutes, add a new virtual machine,

- Keep running until the CPU usage of other machines decreases to 30%,

- If the CPU usage drops below 30%, delete this virtual machine.

## Benefits of Auto Scaling



*Benefits of Auto Scaling*

### Setup Scaling Quickly:

AWS Auto Scaling lets you set target utilization levels for multiple resources in a single, intuitive interface. You can quickly see the average utilization of all of your scalable resources without having to navigate to other consoles.

### Pay Only for What You Need:

AWS Auto Scaling can help you optimize your utilization and cost efficiencies when consuming AWS services so you only pay for the resources you actually need. When demand drops, AWS Auto Scaling will automatically remove any excess resource capacity so you avoid overspending. **AWS Auto Scaling is free to use and allows you to optimize the costs of your AWS environment.**

### Automatically Maintain Performance:

Using AWS Auto Scaling, you maintain optimal application performance and availability, even when workloads are periodic, unpredictable, or continuously changing. AWS Auto Scaling continually monitors your applications to make sure that they are operating at your desired performance levels.

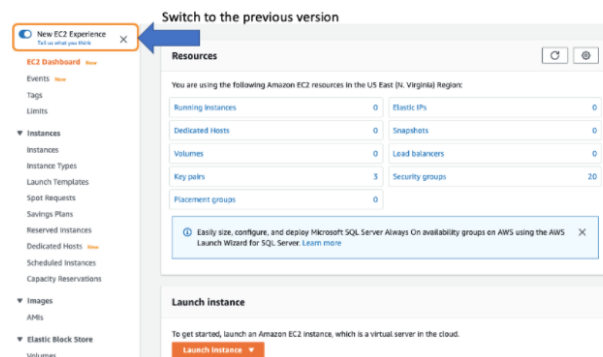### Make Smart Scaling Decisions:

If you try to make scaling instead of AWS Auto Scaling you may not always find an optimum solution. AWS Auto Scaling automatically creates all of the scaling policies and sets targets for you based on your preference. AWS Auto Scaling monitors your application and automatically adds or removes capacity from your resource groups in real-time as demands change.
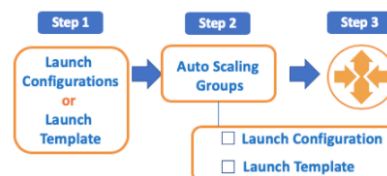
## Auto Scaling Working Process



*Auto Scaling Working Process*

## Auto Scaling Creating Process
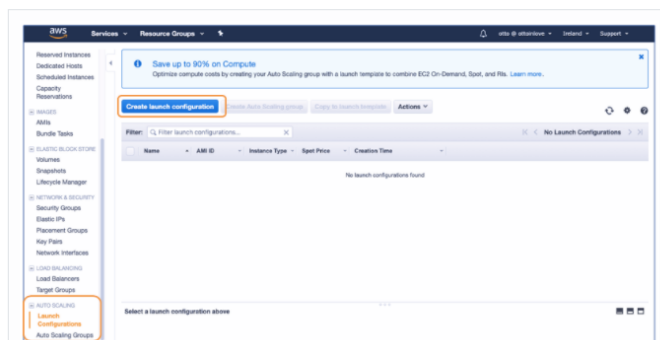


⚠️**Avoid !** :
- Auto Scaling menu is newly changed by AWS . LMS content is more compatible with the previous version. So, it will be better to switch the page version as you see in the picture above.



While creating an Auto Scaling;

- First, we need to **Launch Configuration** or **Launch Templates**. Auto Scaling will create a new instance according to the features determined by Launch Configuration or Launch Templates. In fact, these are templates of instance. **Launch Template** is designed in place of **Launch Configuration**. But, we'll use **Launch Configuration** option for now.

- Then, we create an Auto Scaling Group. While creating an Auto Scaling Group we can choose either **Launch Configuration** or **Launch Templates** that we created before.
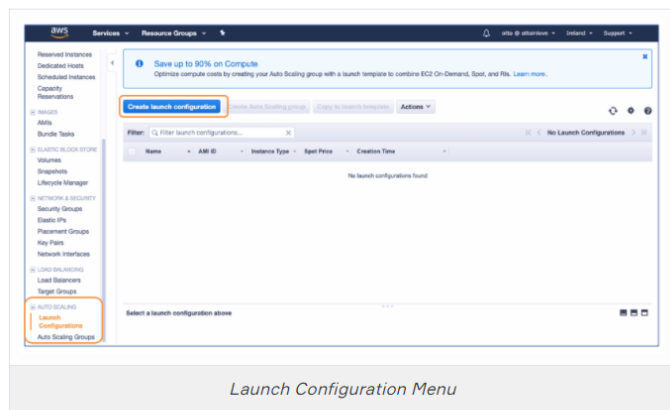
- Finally, we create an Auto Scaling.

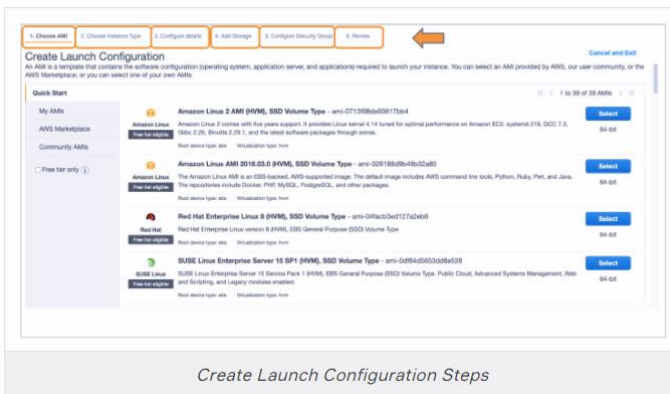## Launch Configurations Menu



*Launch Configuration Menu*

In the EC2 AWS Management console, select Launch Configurations from the left-hand menu and then click **Create Launch Configuration** to start.

## Launch Configurations Menu



*Launch Configuration Menu*

In the EC2 AWS Management console, select Launch Configurations from the left-hand menu and then click **Create Launch Configuration** to start.

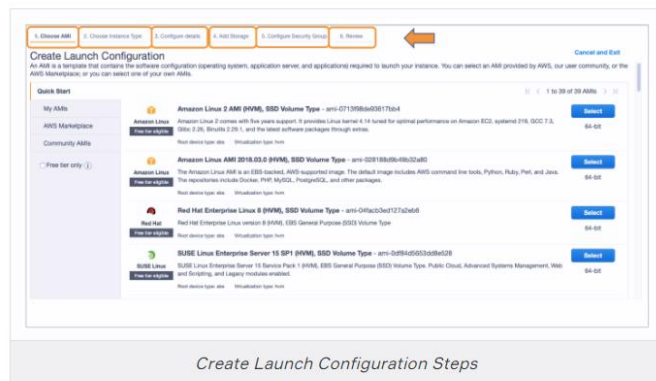## Create Launch Configurations



*Create Launch Configuration Steps*

As you see in the picture above, it is the same menu and the same steps as if you launch an EC2 instance.

We'll determine what we need to choose from disk size, IAM role, key pair to AMI type, we will do everything just like creating a new virtual machine.

But at the end, we will not create a virtual machine. Instead, we'll actually save these things as a template and Auto Scaling will create a new instance according to this configuration when needed.

Let's see.

## Launch Configuration Steps



*Create Launch Configuration Steps*

### Step 1 Choose AMI:.

We select the Linux AMI top of the list.

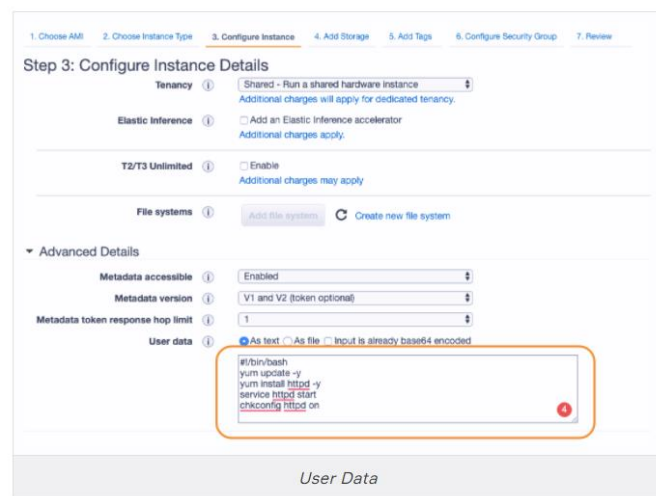### Step 2 Choose Instance Type:

t2.micro is enough for now. We select this one.

### Step 3 Configure Details:

- **Name:** Let's name as First-Launch-Config
- **Purchasing option Request:** No need to Spot Instances.
- **IAM Role:** We choose an existing role.
- **Monitoring:** We don't want to open Detailed Monitoring in, for now, so keep it as is.

- **Advanced Details:** To ensure instance to be installed web server software automatically when created, we enter the following command in the **User Data** section, as you see in the picture below:

```
#!/bin/bash
yum update -y
yum install httpd -y
service httpd start
chkconfig httpd on
```



*User Data*

### Step 4 Choose Instance Type:

In the storage part, the 8 GB gp2 root device is sufficient in our example.

### Step 5 Configure Security Group:

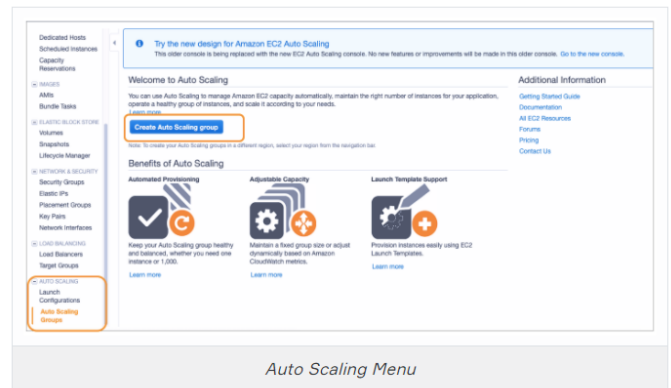In this part, we select the group that we created before.

### Step 6 Review:

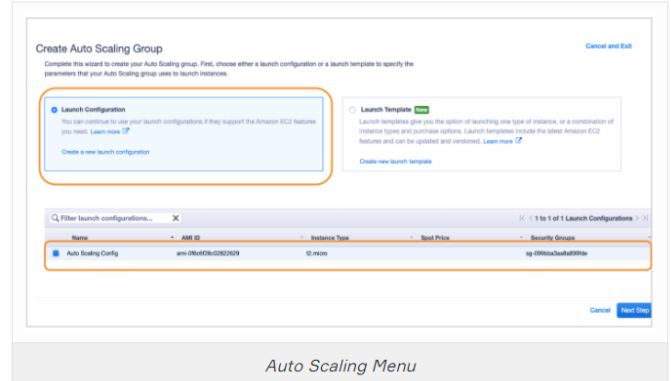Check the details of the configuration here.

- Click the **Create Launch Configuration** tab and pick a key pair. We confirm that we still have the key as we will continue with the existing one.
- We click **Create Launch Configuration** again.

Now, it's ready.

## Creating an Auto Scaling Group



*Auto Scaling Menu*

In the EC2 AWS Management console, select Auto Scaling Group from the left-hand menu at the bottom and then click **Create Auto Scaling Group** to start.



*Auto Scaling Menu*

When the page is opened;

First, we are asked what type of template to use. We have 2 options: Launch Configuration and Launch Template.

Since we created Launch Configuration before, we will continue with **Launch Configuration**.

Then we select the Launch Configuration that we just created, and click **Next** to continue.

## Step 1 - Configure Auto Scaling Group Details



*Configure Auto Scaling Group Details*

**Name:**

The Auto Scaling Group is named as First-Group.

**Group Size:**

This part is critical for the Auto Scaling Group process. Here, we need to decide how many machines will be created by Auto Scaling.

According to our scenario, there should be at least 2 machines at a time so that one can back up the other in case of malfunction. That's why we set the number of Instances that this AutoScaling Group will create as 2.

**Network:**

We select default VPC

**Subnet:**

Let's select the 3 subnets under the default VPC for now.

## Step 1 - Advanced Details



*Advanced Details*

### Advanced Details Section

- **Load Balancing:**

First, AWS asks you whether you want to automatically add the machines you will create with this Auto Scaling to the back of a load balancer you have already created or not.

Check the box that says **Receive traffic from one or more load balancers**. Because we want our Auto Scaling to be coordinated with Load Balancer.

**Classic Load Balancers:**

If you want to check this, you can not handle it. Because we created the application load balancer, not a classic load balancer.

**Target Groups:**

We select our target group named Target-First that is created for the load balancer.

## Step 2 - Configure Scaling Policies



*Configure Scaling Policies*

This is the most important part of creating an Auto Scaling process where we determine the rules.

We have two options here. These are;

- **Keep This Group at Initial Size**

- **Use Scaling Policies to Adjust the Capacity of This Group**

If we choose the first option, Keep This Group at Initial Size, we ensure that Auto Scaling Group will always have 2 virtual machines running under all circumstances. Because we determined Group Size as 2 in "Configure Auto Scaling Group Details" section

So, when this is selected, the Auto Scaling Group will create 2 machines for us.

If we somehow terminate these 2 machines, Auto Scaling will re-create 2 machines.

And, if there is a problem with one of the virtual machines, Auto Scaling will create a new one to delete this machine.

But that's not what we want. In fact,

- We want new virtual machines to be added when we receive too many requests for them or,

- If this intense demand is reduced, we want these added machines to be deleted afterward.

We can do this with the second option. Let's see how.

## Step 2 - Configure Scaling Policies (Using a Scaling Policies)



*Scaling Policies*

This is the most important part where Auto Scaling is planned. Here we asked to prefer scaling values first. We prefer a range of **2** and **4**

This means that;

- Auto Scaling will ensure that there are at least 2 machines at a time.

- But, if the capacity usage increase Auto Scaling will add up to 4 new machines or,

- If the capacity usage decrease, the number of virtual machines will be reduced to 2 again.

**Name:**

We write the name as **Scale Group Size**

## Metric Type:

Now, we will determine the rules according to which rule we will set the number of machines between 2 and 4.

Here we have 4 metric types as you see in options. But the most common ones are **Application Load Balancer Request Count Per Target** and **Average CPU Utilization**.

## Application Load Balancer Request Count Per Target:

This option allows auto-scaling based on the number of user requests selected.

## Average CPU Utilization:

If you select this option, Auto Scaling takes into account CPU usage while adding or removing the instance. Let's continue with this option.

## Target Value:

We select 30 as the target value. It means that,

- If the average CPU of the machines in the environment is below 30%, we want the number of the instance to be stable.

- But, if CPU usage exceeds 30%, this time we order to AWS to add machines here.

Now we have described this in a simple way, but at the bottom of the page, you'll see the **Scale the Auto Scaling Group using step or simple scaling policies** option. If you click this tab you can reach advance settings. Since we prefer the simple scaling, for now, leave it as is.
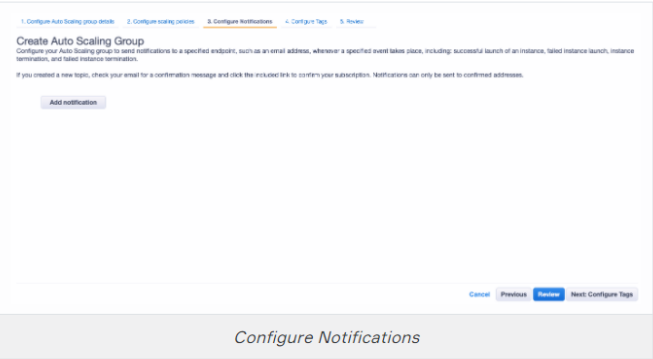
**30% of the Average CPU Usage** option is enough for our scenario.

## Instances Need:

This is the amount of time that your instances need to warm up. Keep it as 300
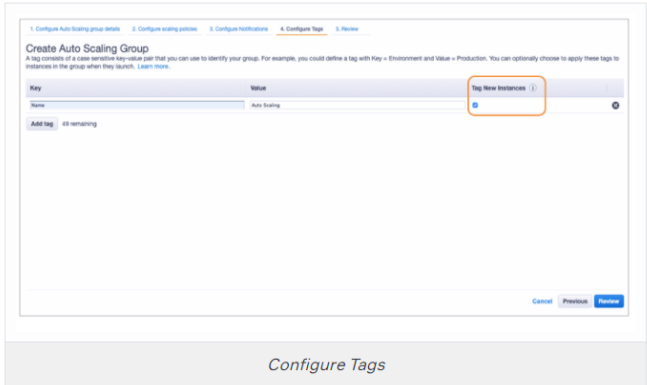
Let's click next and go on.

## Step 3 - Configure Notifications



*Configure Notifications*

If we want AWS to send us a notification when Auto Scaling creates or removes a new machine, we can click **Add Notification**.

Let's skip this step for now and click next.

## Step 4 - Configure Tags
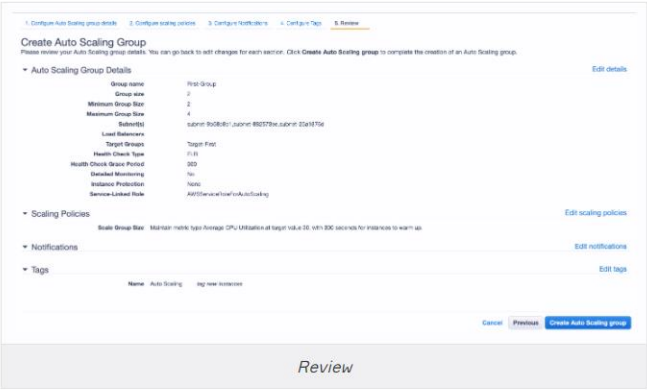


*Configure Tags*

Here, we can determine the tags we want. For example, let's write **Name** for Key and **Auto Scaling** for Value

Then, if we click **Tag New Instances** flag here, it will assign this tag to all instances created with the Auto Scaling.

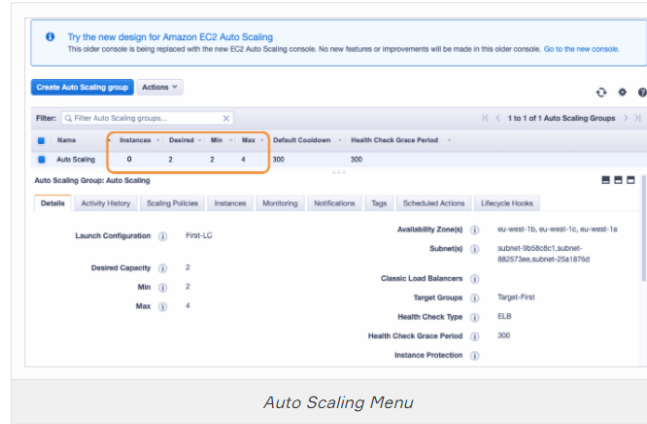We can give the same tag to all newly created instances in this way.

## Step 5 - Review



*Review*

Here we can see our selections. After completing the check, click Create Auto Scaling Group to complete the creation process.

So, we created Auto scaling group.
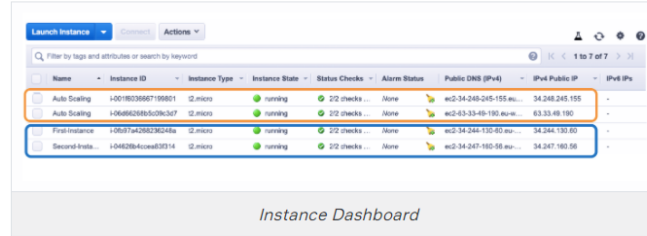
## Checking Auto Scaling Activities - General



*Auto Scaling Menu*

Let's go to the Auto Scaling menu and check the values.

Here we see; **Instance:** 0, **Desired:** 2, **Min:** 2 and **Max:** 4

It means there is no machine appearing in the system now. We desire 2 machines in ideal conditions and at least 2 machines always must be in use. The maximum number of instances that can be in use is 4.

After refreshing, 2 instances will be seen on the dashboard. Let's see the instances on the next page.

## Auto Scaling Group Settings-1



*Instance Dashboard*

Let's go to the Instance Dashboard and examine the number of instances.

If we look at the Instance dashboard shown in the picture above, we see that there are totally four instances in use. The top two are created by Auto Scaling, and their names are **Auto Scaling**. The remaining two instances belong to the Load Balancer.

In the beginning, as you remember, we had no machine appearing in the system, we desired 2 machines in ideal conditions and at least 2 machines would be in use all the time. The maximum number of instances that would be in use was considered as 4. (**Instance:** 0, **Desired:** 2, **Min:** 2 and **Max:** 4)
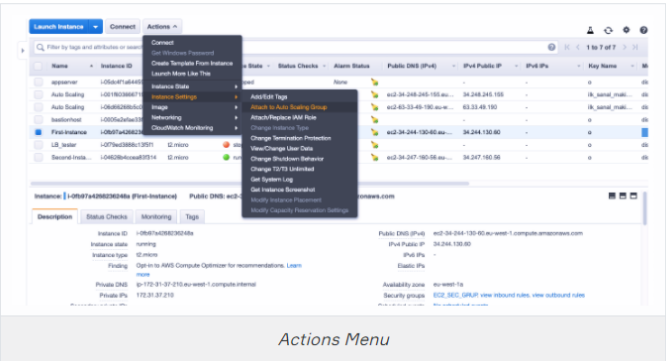
But now, there are 4 servers here. We had 2 servers already in the beginning, but Auto Scaling Group created 2 more servers.

Why Auto Scaling didn't consider the Load balancer's instances?

## Auto Scaling Group Settings - 2

One of the rules is; **The Auto Scaling Group does not consider your current instance capacity.** So, even if you have 200 virtual machines in the environment, it still creates extra machines that you determine when configuring Auto Scaling.
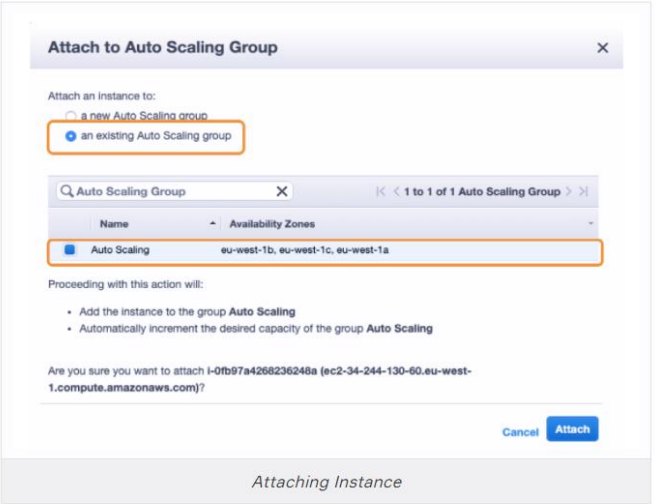
Because the Auto Scaling Group doesn't know that these machines are built to do this job. So, we need to attach the machines to the Auto Scaling Group.



*Actions Menu*

After selecting the Instances menu, we check the instance we want to attach to ASG.

Then we select **Actions> Instance Settings> Attach to Auto Scaling Group** option as you see in the picture above.

After that, we check **an existing Auto Scaling Group** option and select Auto Scaling Group that we have just created. Then, complete the process by clicking **Attach.**
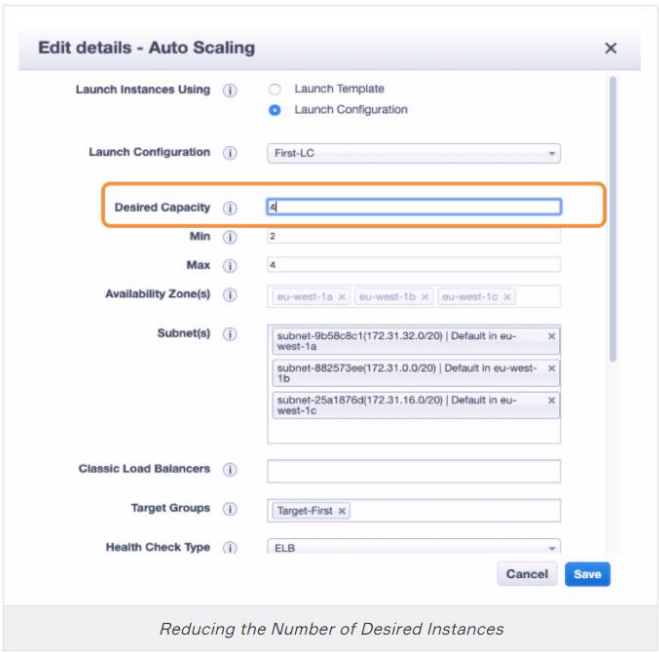


*Attaching Instance*

We repeat the same process for the other virtual machines.

## Auto Scaling Group Settings - 3

Now come back to the Auto Scaling Group menu and control the instances. But this time we see Desired: 4, Min: 2, Max: 4.

Why the desired instance number appears 4, we originally determined it as 2.

Because the other rule for Auto Scaling is; **Manually added servers increase the number of desired servers that you determine before.**



*Reducing the Number of Desired Instances*

To fix this problem, choose your Auto Scaling Group by clicking the checkbox on the left, then select **Edit** from the **Actions** tab for correction as you see in the picture above. Then we reduce the desired capacity value to 2.

Now if you check the instance' status, you'll see that Auto Scaling will terminate two of its virtual machines.

## Testing Auto Scaling

Now let's increase the CPU load of the machines to 30% and check whether the Auto Scaling Group is doing its job depending on this condition.

For this, there is an application called Stress, which is used in Ubuntu machines to increase CPU usage of the machines.

- First connect to ec2 with Terminal
- Let's get root privilege with this command seen below

```
sudo su
```

- Then install stress application with the following commands and increase CPU usage.

```
yum install https://dl.fedoraproject.org/pub/epel/epel-release-latest-7.noarch.
yum install stress
stress --cpu 80 --timeout 2000
```
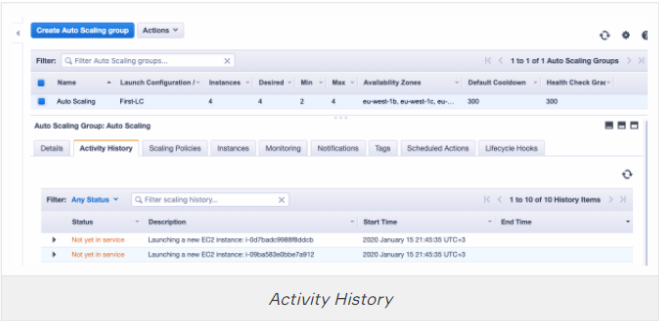
- Thus stress application started to increase CPU for this machine.
- We can do the same process in the other running machine. The CPU increase will gradually increase.

Let's wait for a while then see the result.

## Auto Scaling Test Result

Let's see how Auto Scaling responded to the requests. So click **Auto Scaling Group** on the left-hand menu.

- If you select **Activity History** tab you'll see Auto Scaling Group starts to create new instances as you see in the picture below.



*Activity History*

- You can also see if you click **Instance** tab that 2 machines are in **pending** status as seen in the picture below.



*Instance*

So, as a result, when the CPU usage increased, our Auto Scaling started to create new instances to handle the request.

Congratulations! It works.