# AUTOSCALING

- Autoscaling is a service provided by AWS that automatically increasing the number of EC2 instances during demand spikes to maintain performance and decrease the number of EC2 instances during lulls to reduce costs.

- This autoscaling will increase or decrease the number of instances based on chosen **Cloudwatch** metrics.

- Autoscaling helps you to maintain application availability as per the conditions you define.

For example: If your application's demand increases unexpectedly, autoscaling can automatically scale up ( add instance ) to meet the demand and terminate instances when the demand decreases.  This is known as "elasticity" in the AWS environment.

- Autoscaling is well suited both to applications that have stable demand patterns or that experience hourly, daily, or weekly variability in usage.

-Autoscaling is also used to help ensure that you are running your desired number of amazon EC2 instances.

**Components of Autoscaling**

Autoscaling has 2 main components

1) Launch configuration
2) Auto scaling Group

**Launch Configuration** : The "EC2" template used when the auto scaling group needs to provision an additional instance ( i.e. AMI, instance type, user-data, storage, tags, security groups etc)

**Auto scaling Group :**

All the rules and settings that govern if/when an EC2 instance is automatically provisioned or terminated.

For example

1) number of MIN & MAX allows instances
2) VPC & AZs to launch instances into
3) if provisioned instances should receive traffic from a ELB
4) Scaling policies  (cloudwatch metrics thresholds that trigger scaling)
5) SNS notifications (to keep you informed when scaling occurs)

**NOTE**: To make Highly Available & Fault tolerant architecture, it MUST have and ELB serving traffic to and AutoScaling Group with a MIN of two instances located in separate availability zones.

## Learn to automatically scale up or down your EC2 infrastructure using Auto Scaling Groups

This lab introduces the basics of Auto Scaling in Amazon Web Services. The Amazon Web Services (AWS) Auto Scaling service automatically adds or removes compute resources allocated for your cloud application, in response to changes in demand. For applications configured to run on a cloud infrastructure, scaling is an important part of cost control and resource management.

Scaling is the ability to increase or decrease the compute capacity of your application either by changing the number of servers (horizontal scaling) or by changing the size of the servers (vertical scaling).

Auto Scaling helps you maintain application availability and allows you to scale your Amazon EC2 capacity up or down automatically according to the defined conditions. You can use Auto Scaling to help ensure that you are running your desired number of Amazon EC2 instances. Auto Scaling can also automatically increase the number of Amazon EC2 instances during demand spikes to

maintain performance and decrease capacity during lulls to reduce costs. AS is well suited to applications that have stable demand patterns, or that experience hourly, daily, or weekly variability in usage.

By completing this lab you will learn about:

- Configuring Auto Scaling to automatically launch web server instances

- Building an elastic cluster by integrating Auto Scaling with an Elastic Load Balancer

- Setting CloudWatch alarms to automatically adjust the size of the web farm based on CPU utilization

- Utilizing Auto Scaling to ensure the availability of steady state resources

AGENDA:

## You'll build and learn following these steps:

Log In to the Amazon Web Service Console
*Your first step to start the laboratory experience*

Auto Scaling Overview
*Give an overview about the creation process of an Auto Scaling Group*

Create a load balancer using ELB
*How to create a load balancer using Elastic Load Balancing service.*

Create a Launch Configuration
*How to create an Auto Scaling Launch Configuration*

Create an Auto Scaling Group
*How to create an Auto Scaling Group by using a specific Launch Configuration*

## `Step 1` Log In to the Amazon Web Service Console

This laboratory experience is about Amazon Web Services and you will use the AWS Management Console in order to complete all the lab steps.

# Amazon Web Services

## Compute

**EC2**
Virtual Servers in the Cloud

**EC2 Container Service**
Run and Manage Docker Containers

**Elastic Beanstalk**
Run and Manage Web Apps

**Lambda**
Run Code without Thinking about Servers

## Storage & Content Delivery

**S3**
Scalable Storage in the Cloud

**CloudFront**
Global Content Delivery Network

**Elastic File System**
Fully Managed File System for EC2

**Glacier**
Archive Storage in the Cloud

**Snowball**
Large Scale Data Transport

**Storage Gateway**
Hybrid Storage Integration

## Database

**RDS**
Managed Relational Database Service

**DynamoDB**
Managed NoSQL Database

**ElastiCache**
In-Memory Cache

**Redshift**
Fast, Simple, Cost-Effective Data Warehousing

**DMS**
Managed Database Migration Service

## Networking

**VPC**
Isolated Cloud Resources

**Direct Connect**
Dedicated Network Connection to AWS

**Route 53**
Scalable DNS and Domain Name Registration

## Developer Tools

**CodeCommit**
Store Code in Private Git Repositories

**CodeDeploy**
Automate Code Deployments

**CodePipeline**
Release Software using Continuous Delivery

## Management Tools

**CloudWatch**
Monitor Resources and Applications

**CloudFormation**
Create and Manage Resources with Templates

**CloudTrail**
Track User Activity and API Usage

**Config**
Track Resource Inventory and Changes

**OpsWorks**
Automate Operations with Chef

**Service Catalog**
Create and Use Standardized Products

**Trusted Advisor**
Optimize Performance and Security

## Security & Identity

**Identity & Access Management**
Manage User Access and Encryption Keys

**Directory Service**
Host and Manage Active Directory

**Inspector**
Analyze Application Security

**WAF**
Filter Malicious Web Traffic

**Certificate Manager**
Provision, Manage, and Deploy SSL/TLS Certificates

## Analytics

**EMR**
Managed Hadoop Framework

**Data Pipeline**
Orchestration for Data-Driven Workflows

**Elasticsearch Service**
Run and Scale Elasticsearch Clusters

**Kinesis**
Work with Real-Time Streaming Data

**Machine Learning**
Build Smart Applications Quickly and Easily

## Internet of Things

**AWS IoT**
Connect Devices to the Cloud

## Game Development

**GameLift**
Deploy and Scale Session-based Multiplayer Games

## Mobile Services

**Mobile Hub**
Build, Test, and Monitor Mobile Apps

**Cognito**
User Identity and App Data Synchronization

**Device Farm**
Test Android, iOS, and Web Apps on Real Devices in the Cloud

**Mobile Analytics**
Collect, View and Export App Analytics

**SNS**
Push Notification Service

## Application Services

**API Gateway**
Build, Deploy and Manage APIs

**AppStream**
Low Latency Application Streaming

**CloudSearch**
Managed Search Service

**Elastic Transcoder**
Easy-to-Use Scalable Media Transcoding

**SES**
Email Sending and Receiving Service

**SQS**
Message Queue Service

**SWF**
Workflow Service for Coordinating Application Components

## Enterprise Applications

**WorkSpaces**
Desktops in the Cloud

**WorkDocs**
Secure Enterprise Storage and Sharing Service

**WorkMail**
Secure Email and Calendaring Service

## Resource Groups    Learn more

A resource group is a collection of resources that share one or more tags. Create a group for each project, application, or environment in your account.

**Create a Group**    **Tag Editor**

## Additional Resources

**Getting Started** ⤢
Read our documentation or view our training to learn more about AWS.

**AWS Console Mobile App** ⤢
View your resources on the go with our AWS Console mobile app, available from Amazon Appstore, Google Play, or iTunes.

**AWS Marketplace** ⤢
Find and buy software, launch with 1-Click and pay by the hour.

**AWS re:Invent Announcements** ⤢
Explore the next generation of AWS cloud capabilities. See what's new

## Service Health

✓ All services operating normally.

Updated: Oct 07 2016 11:21:00 GMT-0300
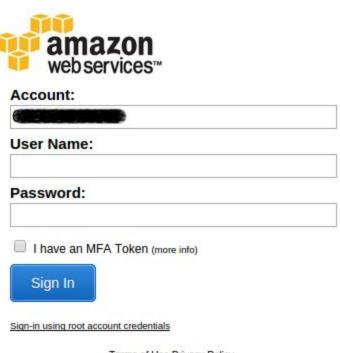
Service Health Dashboard

The AWS Management Console is a web control panel for managing all your AWS resources, from EC2 instances to SNS topics. The console enables cloud management for all aspects of the AWS account, including managing security credentials, or even setting up new IAM Users.

# Log in to the AWS Management Console

In order to start the laboratory experience, open the Amazon Console by clicking this button:

 OPEN AWS CONSOLE

We created a Console User just for you. Log in with the username **student** and the password **Ca1_K0Q2g0ug** .

**amazon**
**web services**™

**Account:**

[redacted]

**User Name:**

**Password:**

☐ I have an MFA Token (more info)

Sign In

Sign-in using root account credentials

Terms of Use Privacy Policy
© 1996-2014, Amazon Web Services, Inc. or its affiliates.

# Select the right AWS Region

Amazon Web Services is available in different regions all over the world, and the console lets you provision resources across multiple regions. You usually choose a region that best suits your business needs to optimize your customer's experience, but you must use the region **US West (Oregon)** for this laboratory.

You can select the **US West (Oregon)** region using the upper right dropdown menu on the AWS Console page.

## **Step 2** Auto Scaling Overview

Before going to the AWS console and creating an Auto Scaling Group, let's take a quick look at the components of an Auto Scaling Group. AWS has done a great job defining them so we'll use the official definition:

**Groups**
Your EC2 instances are organized into groups so that they can be treated as a logical unit for the purposes of scaling and management. When you create a group, you can specify its minimum, maximum, and, desired number of EC2 instances. For more information, see Auto Scaling Groups.
**Launch configurations**
Your group uses a launch configuration as a template for its EC2 instances. When you create a launch configuration, you can specify information such as the AMI ID, instance type, key pair, security groups, and block device mapping for your instances. For more information, see Launch Configurations.

You can read the full documentation here

http://docs.aws.amazon.com/autoscaling/latest/userguide/WhatIsAutoScaling.html

In this lab, we will learn to create an Auto Scaling Group with these components and place it behind an Elastic Load Balancing (ELB). Don't worry if you don't fully understand all the components yet. We will talk in greater detail about each of the components as we create them.

At the end of this lab we'll have an Auto Scaling Group with some web server instances behind an ELB. Although this lab focuses on Auto Scaling, it is important to mention that to have an Auto Scaling Group behind an ELB, it is necessary to create the ELB first. In the next step, we will begin exploring elements in the AWS console by creating an ELB.
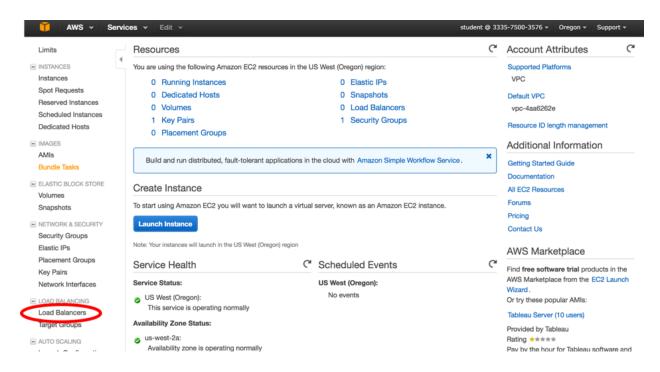
## Step 3 Create a load balancer using ELB

Elastic Load Balancing (ELB) automatically distributes incoming application traffic across multiple Amazon EC2 instances. It enables you to achieve greater fault tolerance in your applications and seamlessly provides the correct amount of load balancing capacity needed in response to incoming application traffic.
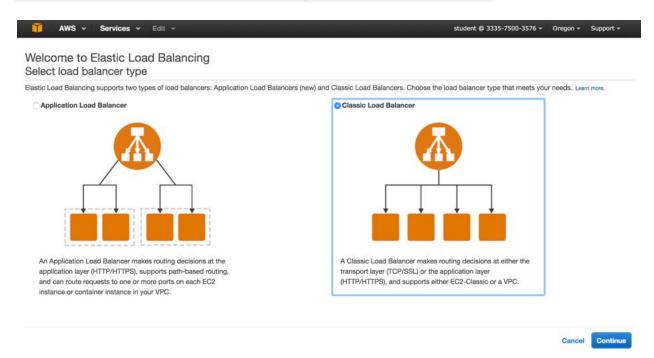
Elastic Load Balancing detects unhealthy instances within a pool and automatically reroutes traffic to healthy instances until the unhealthy instances have been restored to health. Customers can enable Elastic Load Balancing within a single Availability Zone or across multiple zones for greater consistent application performance.

You can create your first ELB by taking the following steps:

1. Select EC2 from the AWS Service List

2. From the EC2 dashboard, click the **Load Balancers** link in the Load Balancing group. The list of all already-created Load Balancers appears--this list will most likely be empty.

3. Click the blue **Create Load Balancer** button

4. Select the **Classic Load Balancer** option and click Continue
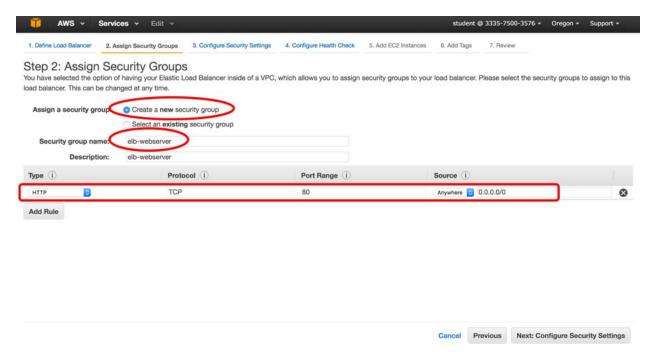


5. On the **Define Load Balancer** step, type a load balancer name (e.g., "web")

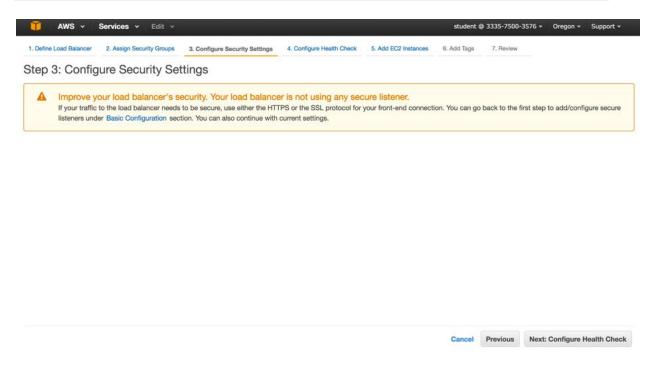6. Select **Enable advanced VPC configuration**

7. Select two subnets, one from the **us-west-2a Availability Zone** and one from the **us-west-2b Availability Zone**



8. Then click the **Next: Assign Security Groups** button

9. In the **Assign Security Groups** section, select *Create a new security group*

10. Type a Security group name (e.g., "elb-webserver") and a description

11. Create a single firewall rule of *type* **HTTP**, *protocol* **TCP**, *port range* **80**, and *source* **Anywhere**

12. Click **Next: Configure Security Settings**.

13. Ignore the warning in the **Configure Security Settings** section. We are only serving the HTTP protocol in this exercise, so these settings are not required



14. Click **Next: Configure Health Check**.

15. In the **Configure Health Check** section, replace the default value of **Ping Path** with a single forward slash ("/")
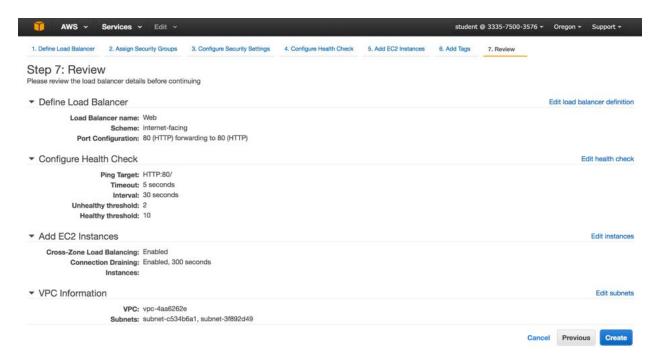


16. Click **Next: Add EC2 Instances**

17. In the **Add EC2 Instances** section, you should see a "No instances available" message. This is because we have yet created and launched our Auto Scaling Group

18. Click **Next: Add Tags to continue**

19. You may leave the fields blank in the **Add Tags** section

20. Click the **Review and Create** button to continue

21. **Review** your settings, then click **Create** when ready

22. Wait for the Load Balancer Creation Status to populate with the message, "Successfully created load balancer." Click **Close**



**Step 4** Create a Launch Configuration

A **Launch Configuration** is a template that the Auto Scaling group uses to launch Amazon EC2 instances. If you've launched an individual EC2 instance before,

you've already walked through the process of defining compute characteristics such as the instance type, security groups, and configuration scripts. A launch configuration allows you to define these same characteristics, which are then applied to any instances launched in the Auto Scaling group.
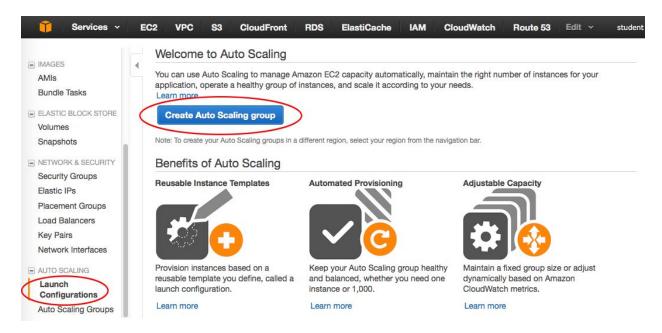
You create the launch configuration by including information such as the Amazon machine image ID to use for launching the EC2 instance, the instance type, key pairs, security groups, and block device mappings, among other configuration settings. When you create your Auto Scaling group, you must associate it with a launch configuration. You can attach only one launch configuration to an Auto Scaling group at a time and it cannot be modified.

Let's start creating our Auto Scaling Group by first defining a **Launch Configuration**.

1.Navigate to the EC2 service from the AWS dashboard:

Compute

EC2
Virtual Servers in the Cloud

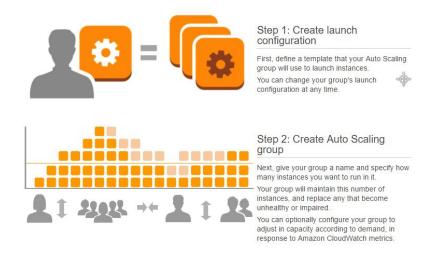2. Open the **Launch Configurations** page and click on the **Create Auto Scaling Group** button.

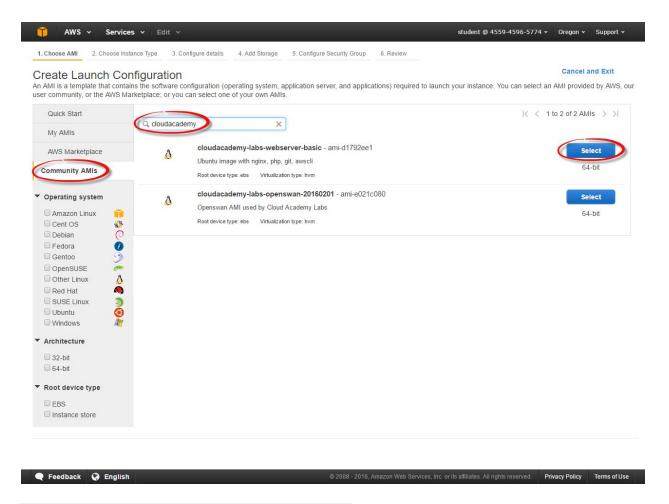This brings you to the Create Auto Scaling group wizard. Click on the **Create Launch configuration** button.



From there the AWS Management Console guides you though each required step and displays a graphical interface that is similar to the Launch Instance Wizard.

The first step is the AMI selection. You have to select the AMI that will be used by all the EC2 instances of the Auto Scaling group. The Cloud Academy DevOps team created a specific AMI for this laboratory. You can find it among the Community AMIs by searching for the word "cloudacademy" in the AMI search box.
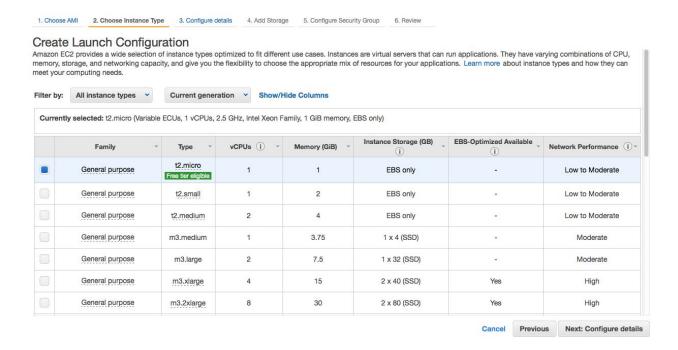
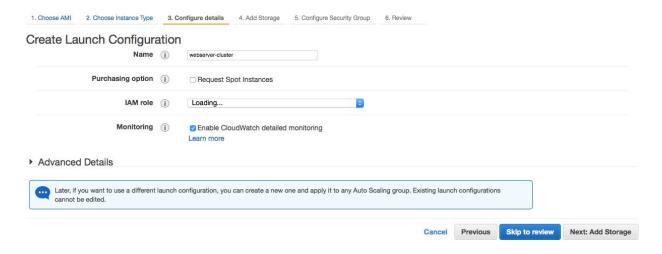3. Select the "**cloudacademy-labs-webserver-basic**" - (ami-d1792dee1) AMI and click Select.

The next step is choosing the instance type

4. Select the t2.micro type and click on the **Next: Configure details** button.

## Create Launch Configuration

Amazon EC2 provides a wide selection of instance types optimized to fit different use cases. Instances are virtual servers that can run applications. They have varying combinations of CPU, memory, storage, and networking capacity, and give you the flexibility to choose the appropriate mix of resources for your applications. Learn more about instance types and how they can meet your computing needs.

Filter by:  All instance types ▾    Current generation ▾    **Show/Hide Columns**

Currently selected: t2.micro (Variable ECUs, 1 vCPUs, 2.5 GHz, Intel Xeon Family, 1 GiB memory, EBS only)

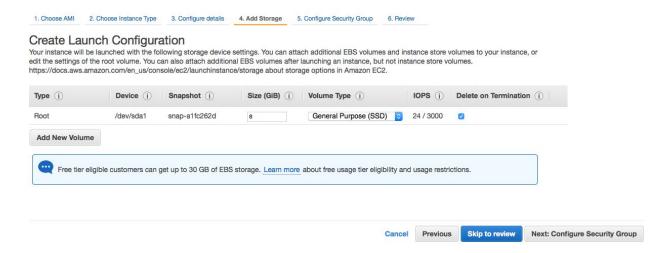| | Family | Type | vCPUs ⓘ | Memory (GiB) | Instance Storage (GB) ⓘ | EBS-Optimized Available ⓘ | Network Performance ⓘ |
|---|---|---|---|---|---|---|---|
| ☑ | General purpose | t2.micro<br>Free tier eligible | 1 | 1 | EBS only | - | Low to Moderate |
| ☐ | General purpose | t2.small | 1 | 2 | EBS only | - | Low to Moderate |
| ☐ | General purpose | t2.medium | 2 | 4 | EBS only | - | Low to Moderate |
| ☐ | General purpose | m3.medium | 1 | 3.75 | 1 x 4 (SSD) | - | Moderate |
| ☐ | General purpose | m3.large | 2 | 7.5 | 1 x 32 (SSD) | - | Moderate |
| ☐ | General purpose | m3.xlarge | 4 | 15 | 2 x 40 (SSD) | Yes | High |
| ☐ | General purpose | m3.2xlarge | 8 | 30 | 2 x 80 (SSD) | Yes | High |

Cancel   Previous   **Next: Configure details**

5. The **Configure details** step asks you to name your launch configuration, enter a friendly name (e.g., webserver-cluster')

6. Enable detailed monitoring

7. Click on **Next: Add Storage**

## Create Launch Configuration

| Name ⓘ | webserver-cluster |
|---|---|
| **Purchasing option** ⓘ | ☐ Request Spot Instances |
| **IAM role** ⓘ | Loading... ▾ |
| **Monitoring** ⓘ | ☑ Enable CloudWatch detailed monitoring<br>Learn more |

▸ Advanced Details

💬 Later, if you want to use a different launch configuration, you can create a new one and apply it to any Auto Scaling group. Existing launch configurations cannot be edited.

Cancel   Previous   **Skip to review**   Next: Add Storage

The **Add Storage** step allows you to add or increment the size of any EBS volume linked to each EC2 instance that will be started by the Auto Scaling group.

8. In order to complete this laboratory exercise, leave the defaults and do not add any EBS volumes. Then click on **Next: Configure Security Group**

**N.B.**: You should use big EBS volumes only if your software requires storage space to process the application data. If you need to store raw or proceseed data, you should use Amazon S3, Redshift, DynamoDB or another storage/database service provided by Amazon.



9. **Create a new Security Group** for your Auto Scaling Group.

10. Choose a name (e.g., Webserver-cluster) and description

11. Add 2 rules:

**1st Rule**

Type=SSH

Protocol=TCP

Port Range=22

Source=My IP

**2nd Rule**

Type=HTTP

Protocol=TCP

Port Range=80

Source=Custom IP (enter 172.31.0.0/16 for the IP)

The default Amazon VPC subnet range is **172.31.0.0/16**. You can use it to allow the HTTP traffic, so the Elastic Load Balancing instance will be able route the HTTP requests to the instances of the Auto Scaling group.



12. Click the blue **Review** button.

13. Once you have reviewed the details for accuracy, click the blue **Create launch configuration** button.



You will be presented with the *Select an existing key pair or create a new key pair* dialogue box. Notice that you will use this Key Pair to access all the instances that are going to be launched by the Auto Scaling service with this Launch Configuration, so secure your Key Pair.

14. Select **Create a new key pair** from the first drop-down menu and type in a Key pair name (e.g., webserver-cluster).

15. Click the **Download Key Pair** button

16. Then click the **Create Launch Configuration** button in this dialogue box



17. Next a screen will appear that tries to configure an Auto Scaling Group, however at this point click on Cancel as we will create this from scratch in the next step
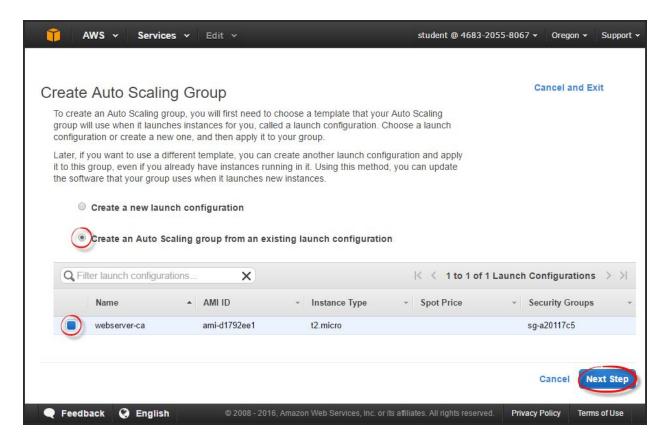
## Step 5 Create an Auto Scaling Group

An Auto Scaling group is a representation of multiple Amazon EC2 instances that share similar characteristics and that are treated as a logical grouping for the purposes of instance scaling and management. For example, if a single application operates across multiple instances, you might want to increase or decrease the number of instances in that group to improve the performance of the application. You can use the Auto Scaling group to automatically scale the number of instances or maintain a fixed number of instances. You create Auto Scaling groups by defining the minimum, maximum, or desired number of running EC2 instances the group must have at any given point of time.

An Auto Scaling group starts by launching the minimum number (or the desired number, if specified) of EC2 instances and then increases or decreases the number of running EC2 instances automatically according to the conditions that you define. Auto Scaling also maintains the current instance levels by conducting periodic health checks on all the instances within the Auto Scaling group. If an EC2 instance within the Auto Scaling group becomes unhealthy, Auto Scaling terminates the unhealthy instance and launches a new one to replace the unhealthy instance. This automatic scaling and maintenance of the instance levels in an Auto Scaling group is the core value of the Auto Scaling service.

1. To create the Auto Scaling group, click on the **Auto Scaling Groups** link in the Auto Scaling menu group and then click the blue **Create Auto Scaling group** button.

2. Select **Create an Auto Scaling group from an existing launch configuration,** select the previously created launch configuration and click Next Step.

3. In the "**Configure Auto Scaling group details**" step, you should use the following settings:

**Group name:** webserver-cluster
**Group size:** 1
**Network:** default
**Subnet:** Select two. The default network in *us-west-2a* and the default network in *us-west-2b*.
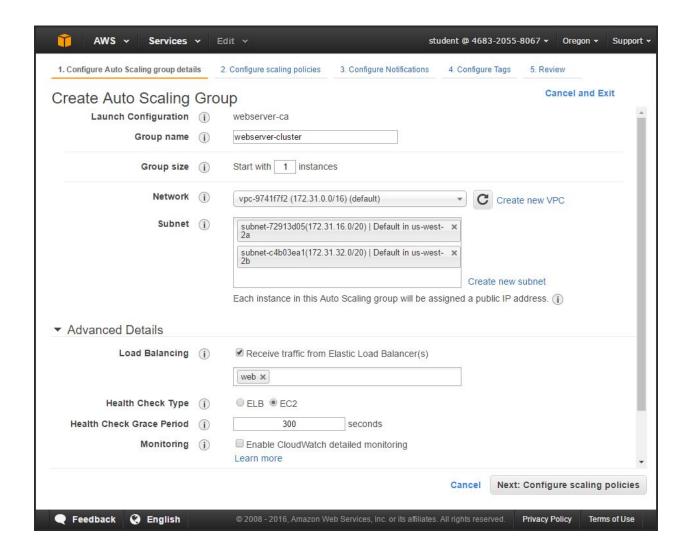
4. Open the Advanced Details, then set as follows:

**Load Balancing:** Check Receive traffic from Elastic Load Balancer(s). Select the "web" ELB you created
**Health Check Type:** ELB
**Monitoring:** Check *Enable CloudWatch detailed monitoring*

5. Once all fields are complete, click **Next: Configure Scaling Policies**.

In this step, you must *Configure scaling policies*, which determine how and when your infrastructure will scale out and scale back.

6. Select the *Use scaling policies to adjust the capacity of this group* button. For this lab you should set your group to scale between **1** and **5** instances.

The Auto Scaling group policies allow you to automatically increase or decrease the group size based upon policies you define. In order to establish an Increase Group size or Decrease Group Size policy, you must create a CloudWatch Alarm and then define which action should be taken if it is triggered.

7. Click **Add new alarm** under the **Increase Group Size section**. A *Create Alarm* dialogue box will pop up.

8. If you want to receive a notification when the alarm is triggered, you need to

set up an **SNS topic.** Check the ***Send a notification to:*** checkbox. Type in a name (e.g., "autoscaling-alarm-up") for the SNS topic and enter at least one email address in the recipients box.

9. Select a metric (e.g., Average, CPU Utilization) and a constraint (e.g., >= 80 percent). Select a count and an interval (e.g., For at least **1** consecutive period of **5 minutes**). Choose a name for the alarm, and then click **Create Alarm**.
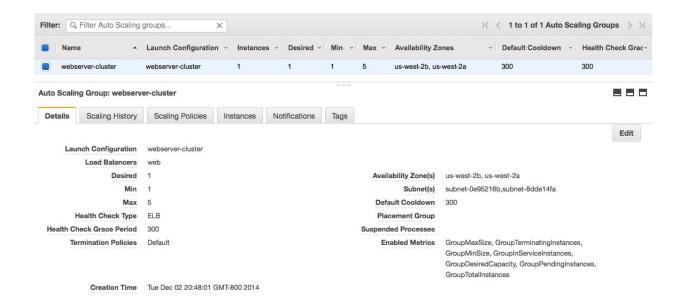


10. Create another alarm with whatever settings you choose for the Decrease Group Size. Click **Next: Configure Notifications**.

11. Configure *Notifications* will notify you whenever an Auto Scaling Group instance is launched or terminated -- with or without success.

12. Click **Add notification**. You can use one of the same SNS topics previously created for the CloudWatch alarms. When you're done, click the blue **Review** button.

13. The **Review** tab allows you to review all the selected options. When you are satisfied, start the creation of your cluster by clicking on **Create Auto Scaling group**.



14. In a few minutes your cluster will be deployed and your EC2 instances will be ready to .

By opening the **Load Balancers** section, selecting your previously created ELB, and then opening the Instances tab, you can see the new Auto Scaling instance(s) automatically added to the ELB configuration.