



[Home](#) » [AWS Cheat Sheets](#) » [AWS Comparison of Services](#) » Step Scaling vs Simple Scaling Policies in Amazon EC2

Step Scaling vs Simple Scaling Policies in Amazon EC2

Amazon's EC2 Auto Scaling provides an effective way to ensure that your infrastructure is able to dynamically respond to changing user demands. For example, to accommodate a sudden traffic increase on your web application, you can set your Auto Scaling group to automatically add more instances. And when traffic is low, have it automatically reduce the number of instances. This is a cost-effective solution since it only provisions EC2 instances when you need them. EC2 Auto Scaling provides you with several dynamic scaling policies to control the scale-in and scale-out events.

In this article, we'll discuss the differences between a simple scaling policy and step scaling policy. And we'll show you how to create an Auto Scaling group with step scaling policy applied.

Simple Scaling

Simple scaling relies on a metric as a basis for scaling. For example, you can set a CloudWatch alarm to have a CPU Utilization threshold of 80%, and then set the scaling policy to add 20% more capacity to your Auto Scaling group by launching new instances. Accordingly, you can also set a CloudWatch alarm to have a CPU utilization threshold of 30%. When the threshold is met, the Auto Scaling group will remove 20% of its capacity by terminating EC2 instances.

When EC2 Auto Scaling was first introduced, this was the only scaling policy supported. It does not provide any fine-grained control to scaling in and scaling out.

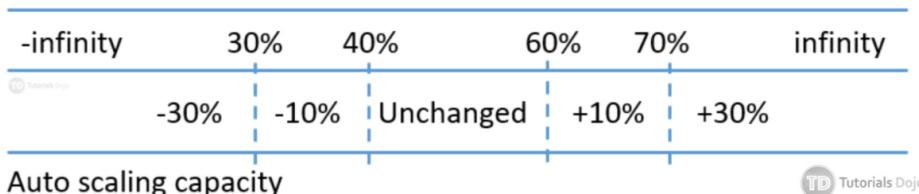
Step Scaling

Step Scaling further improves the features of simple scaling. Step scaling applies "step adjustments" which means you can set multiple actions to vary the scaling depending on the size of the alarm breach.

When a scaling event happens on simple scaling, the policy must wait for the health checks to complete and the cooldown to expire before responding to an additional alarm. This causes a delay in increasing capacity especially when there is a sudden surge of traffic on your application. With step scaling, the policy can continue to respond to additional alarms even in the middle of the scaling event.

Here is an example that shows how step scaling works:

Metric value: CPU Utilization



In this example, the Auto Scaling group maintains its size when the CPU utilization is between 40% and 60%. When the CPU utilization is greater than or equal to 60% but less than 70%, the Auto Scaling group increases its capacity by an additional 10%. When the utilization is greater than 70%, another step in scaling is done and the capacity is increased by an additional 30%. On the other hand, when the overall CPU utilization is less than or

NEW Course – AWS Certified Data Analytics Specialty Practice Exams



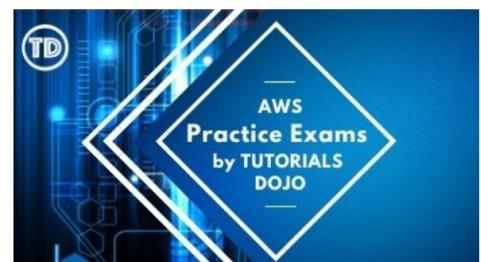
Pass your AWS and Azure Certifications with the Tutorials Dojo Portal



Our Bestselling AWS Certified Solutions Architect Associate Practice Exams



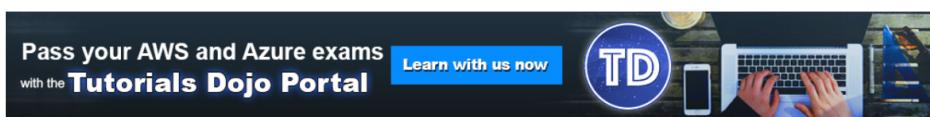
Enroll Now – Our AWS Practice Exams with 95% Passing Rate



Enroll Now – Our Azure Certification Exam

equal to 40% but greater than 50%, the Auto Scaling group decreases the capacity by 10%. And if utilization further dips below 30%, the Auto Scaling group removes 30% of the current capacity.

This effectively provides multiple steps in scaling policies that can be used to fine-tune your Auto Scaling group response to dynamically changing workload.



Creating a Step Scaling Policy for an Auto Scaling Group

Based on the step scaling policy described above, the following guide will walk you through the process of applying this policy when creating your Auto Scaling group.

1. First, create your Launch Configuration for your EC2 instances. Check [this guide](#) if you haven't created one yet.

2. Go to **EC2 > Auto Scaling Groups > Create Auto Scaling group**

3. Select your **Launch Configuration** and click **Next Step**.

4. Configure details for your Auto Scaling group.

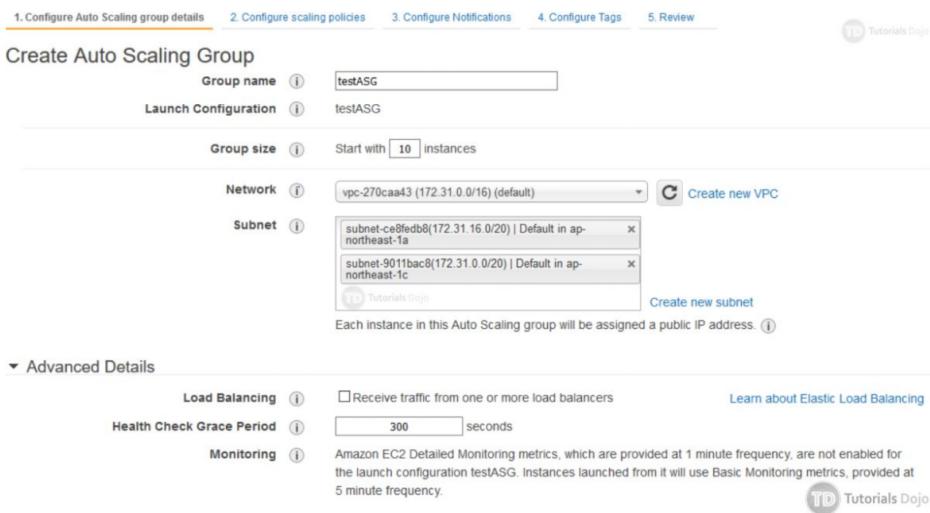
a. **Group name** – descriptive name for this ASG.

b. **Group size** – the initial size of your ASG. Let's set this to 10 for this example.

c. **Network** – the VPC to use for your ASG.

d. **Subnet** – the subnets in the VPC on where to place the EC2 instances. It's recommended to select subnet in multiple availability zones to improve the fault tolerance of your ASG.

e. **Advanced Details** – in this section, you can check the **Load Balancing** option to select which load balancer to use for your ASG. (We won't configure a load balancer for this example). You can also set the **Health Check Grace Period** in this section. This is the length of time that Auto Scaling waits before checking the instance's health status. We'll leave the default to 300 seconds but you can adjust this if you know your EC2 instances need more or less than 5 minutes before they become healthy.



5. Click **Next: Configure scaling policies** to proceed.

6. Here, we'll configure the step scaling policy. Select the "**Use scaling policies to adjust the capacity of this group**" option and this will show an additional section for defining scaling policy. For this example, let's set 5 and 15 as the minimum and maximum size for this Auto Scaling group.

Keep this group at its initial size

Use scaling policies to adjust the capacity of this group

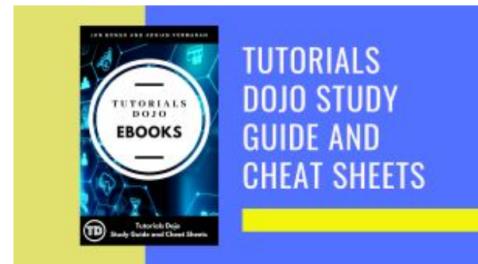
Scale between and instances. These will be the minimum and maximum size of your group.

Tutorials Dojo

Reviewers



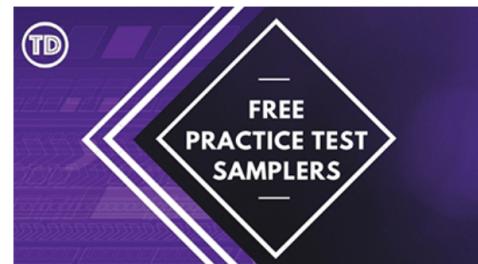
Tutorials Dojo Study Guide and Cheat Sheets eBooks



FREE Intro to Cloud Computing for Beginners



FREE AWS Practice Test Samplers



Browse Other Courses



7. On the Scale Group Size section, you will be able to set the scaling policy for the group. But this is only for simple scaling so you have to click the “**Scale the Auto Scaling group using step or simple scaling policies**” link to show more advanced options for step scaling. You should see the **Increase Group Size** and **Decrease Group Size** section after clicking it.

Increase Group Size

- Name: Increase Group Size
- Execute policy when: No alarm selected
- Take the action: Add 0 capacity units
- Instances need: 300 seconds to warm up after each step

Decrease Group Size

- Name: Decrease Group Size
- Execute policy when: No alarm selected
- Take the action: Remove 0 capacity units

Create a simple scaling policy

8. Now, we can set the step scaling policy for scaling out.

- Set a name for your “**Increase Group Size**” policy. Click “**Add a new alarm**” to add a CloudWatch rule on when to execute the policy.
- On the **Create Alarm** box, you can set an SNS notification. (We won’t add it for this example).
- Create a rule for whenever the **Average CPU Utilization** is greater than or equal to 60 percent for at least 1 consecutive period of 5 minutes. Set a name for your alarm. Click **Create Alarm**.

Create Alarm

You can use CloudWatch alarms to be notified automatically whenever metric data reaches a level you define.

Whenever: Average of CPU Utilization Is \geq 60 Percent For at least 1 consecutive period(s) of 5 Minutes

Name of alarm: awsec2-testASG-CPU-Utilization

CPU Utilization Percent

Time	Value
5/27 08:00	60
5/27 10:00	60
5/27 12:00	60

- For the “**Take the action**” setting, we’ll **Add 10 percent** of the group when CPU Utilization is greater than or equal to **60 and less than 70 percent**.
- Click “Add Step” to add another action, we’ll **Add 30 percent** of the group when CPU Utilization is **greater than or equal to 70 percent**.

Increase Group Size

Name: Increase Group Size

Execute policy when: awsec2-testASG-CPU-Utilization

Recent Posts

- > Amazon Aurora Machine Learning
- > Amazon Comprehend
- > AWS DeepLens

breaches the alarm threshold: CPUUtilization >= 60 for 300 seconds
for the metric dimensions AutoScalingGroupName = testASG

Take the action:

- Add 10 percent of group when 60 <= CPUUtilization < 70
- Add 30 percent of group when 70 <= CPUUtilization < +infinity X

[Add step](#) i

Add instances in increments of at least 1 instance(s)

Instances need: 300 seconds to warm up after each step

 Tutorials Dojo

f. Set 1 for “**Add instances in increments of at least**”. This will ensure that at least 1 instance is added

when the threshold is reached.

g. Set instances need **300 seconds to warm up** after each step.

Instance warmup – this specifies the timeout before the instance’s own metric can be added to the group. Until the warmup time expires, the instance metric (CPU utilization in this case) is not counted toward the aggregated metric of the whole Auto Scaling group.

While scaling in, instances that are terminating are considered as part of the current capacity of the group. Therefore, it won’t remove more instances from the Auto Scaling group than necessary.

9. Next, we can set the step scaling policy for the scaling in.

a. Set a name for your “**Decrease Group Size**” policy. Click “**Add a new alarm**” to add a CloudWatch rule on when to execute the policy.

b. On the **Create Alarm** box, you can set an SNS notification. (We won’t add it for this example).

c. Create a rule for whenever the **Average CPU Utilization** is less than or equal to 40 percent for at least 1 consecutive period of 5 minutes. Set a name for your alarm. Click **Create Alarm**.

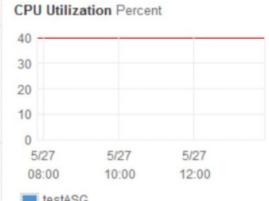
Create Alarm

You can use CloudWatch alarms to be notified automatically whenever metric data reaches a level you define.
To edit an alarm, first choose whom to notify and then define when the notification should be sent.

Send a notification to: No SNS topics found... X

 Tutorials Dojo Whenever: Average of CPU Utilization Is: <= 40 Percent For at least: 1 consecutive period(s) of 5 Minutes Name of alarm: awsec2-testASG-High-CPU-Utilization

CPU Utilization Percent



Cancel Create Alarm

 Tutorials Dojo

d. For the “**Take the action**” setting, we’ll **remove 10 percent** of the group when CPU Utilization is less than or equal to **40 and greater than 30**.

e. Click “Add Step” to add another action, we’ll **remove 30 percent** of the group when CPU Utilization is **less than or equal to 30 percent**.

Decrease Group Size

Name: Decrease Group Size

Execute policy when: awsec2-testASG-High-CPU-Utilization [Edit](#) [Remove](#)
breaches the alarm threshold: CPUUtilization <= 40 for 300 seconds
for the metric dimensions AutoScalingGroupName = testASG

Take the action:

- Remove 10 percent of group when 40 >= CPUUtilization > 30
- Remove 30 percent of group when 30 >= CPUUtilization > -infinity X

[Add step](#) i

Remove instances in increments of at least 1 instance(s)

 Tutorials Dojo

f. Set 1 for “**Remove instances in increments of at least**”. This will ensure that at least 1 instance is removed when the threshold is reached.

10. Click **Next: Configure Notifications** to proceed. On this part, you can click “**Add notification**” so that you can receive an email whenever a specific event occurs. Here’s an example:

Send a notification to: testASG

With these recipients: alerts@tutorialsdojo.com

Whenever instances: launch
 terminate
 fail to launch
 fail to terminate

11. Click **Next: Configure Tags**. Create tags for instances in your Auto Scaling group.

12. Click **Review** to get to the review page.

13. After reviewing the details, click **Create Auto Scaling group**.



Your Auto Scaling group with step scaling policies should now be created. Remember, the initial desired size is 10, with a minimum of 5 and a maximum of 15.

The scale-out rule will have a step scaling policy, a 10% increase if CPU utilization is 60 – 70%, and will add 30% more instance if utilization is more than 70%.

The scale-in rule will have a step scaling policy, a 10% decrease if CPU utilization is 30 – 40%, and will remove 30% more instances if the utilization is less than 30%.

Sources:

<https://docs.aws.amazon.com/autoscaling/ec2/userguide/as-scaling-simple-step.html>

<https://docs.aws.amazon.com/autoscaling/ec2/userguide/Cooldown.html>

<https://docs.aws.amazon.com/autoscaling/ec2/userguide/GettingStartedTutorial.html>

<https://docs.aws.amazon.com/autoscaling/ec2/userguide/create-launch-config.html>

AWS, Azure, and GCP Certifications are consistently among the top-paying IT certifications in the world, considering that most companies have now shifted to the cloud. Earn over \$150,000 per year with an **AWS, Azure, or GCP certification!**

Follow us on [LinkedIn](#), [Facebook](#), or join our [Slack study group](#). More importantly, answer as many [practice exams](#) as you can to help increase your chances of passing your certification exams on your first try!

Did you find our content helpful?

Subscribe to our Newsletter

Sign up to our newsletter and get a **FREE** copy of our Introduction to Cloud Computing eBook.

Name:

Email:

[Submit](#)

We respect your [email privacy](#).

Related Posts



[Amazon Aurora Machine Learning](#)

November 30th, 2020 | 0 Comments



[Amazon Comprehend](#)

November 25th, 2020 | 0 Comments



[AWS DeepLens](#)

November 25th, 2020 | 0 Comments



[Amazon SageMaker](#)

November 25th, 2020 | 0 Comments



[Amazon Lex](#)

November 25th, 2020 | 0 Comments



Tutorials Dojo

Proudly Made in the Philippines

Founded in Manila, Philippines, Tutorials Dojo is your one-stop learning portal for technology-related topics, empowering you to upgrade your skills and your career.

[portal.tutorialsdojo.com](#)

[support@tutorialsdojo.com](#)

Tutorials Dojo Courses

- AWS eBooks
- AWS Certified Cloud Practitioner Courses
- AWS Associate Level Courses
- AWS Professional Level Courses
- AWS Specialty Courses
- Free AWS Practice Test Samplers
- Azure Reviewers



Join us on Slack!

Meet other IT professionals in our [Slack Community](#). Communicate your IT certification exam-related questions (AWS, Azure, GCP) with other members and our technical team.

Follow us on:



© 2020 Tutorials Dojo. All rights reserved.

