

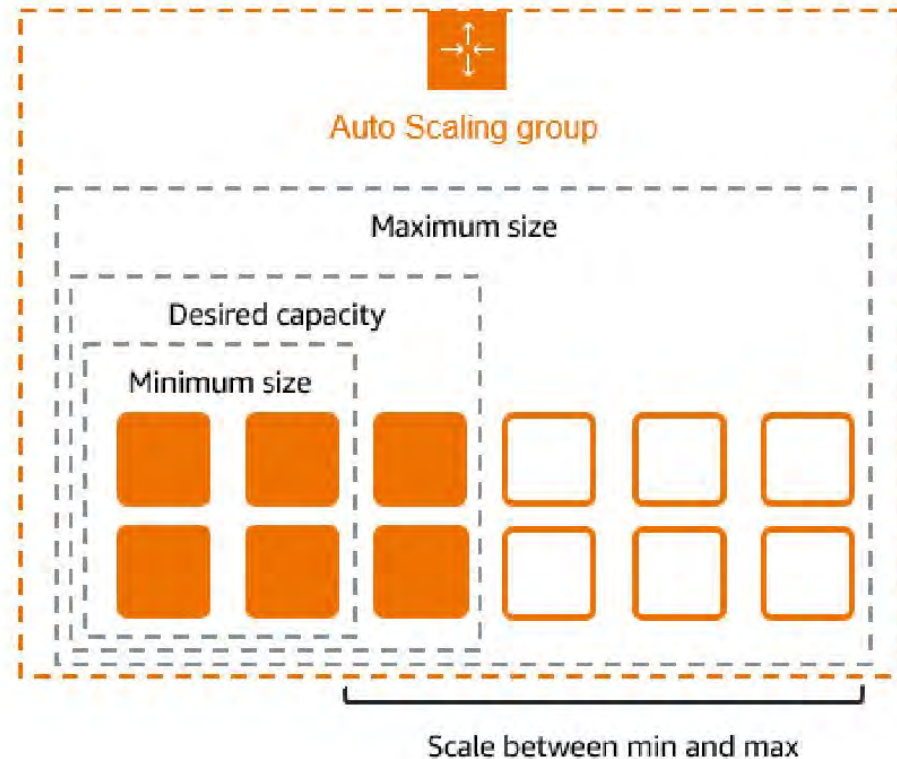
iam**neo**

AWS AutoScaling

What is autoscaling and why is it important?

- Autoscaling is a feature of cloud computing that allows you to automatically adjust the number of computing resources in your system based on demand. This means that if your application suddenly becomes very popular and starts receiving a lot of traffic, autoscaling will automatically add more resources to handle the increased load.
- Autoscaling is important because it helps you ensure that your application is always available and responsive, even during times of high traffic.
- It also helps you save costs by only using the number of resources that you need, rather than paying for resources that are sitting idle.
- With autoscaling, you can achieve better performance, higher availability, and lower costs.

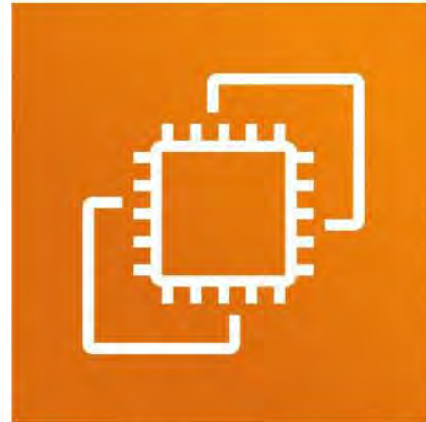
What is autoscaling and why is it important?



What is autoscaling and why is it important?



Scaling based on demand

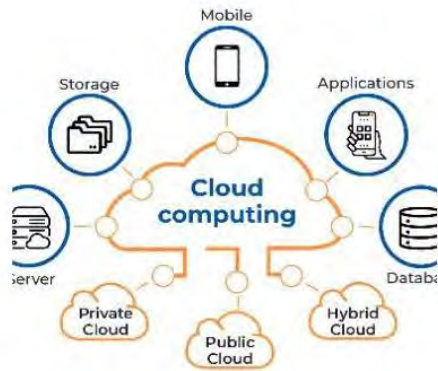


Multiple instance types supported



Integration with CloudWatch

Why Autoscaling Matters?



Scalability and Flexibility



Cost Efficiency

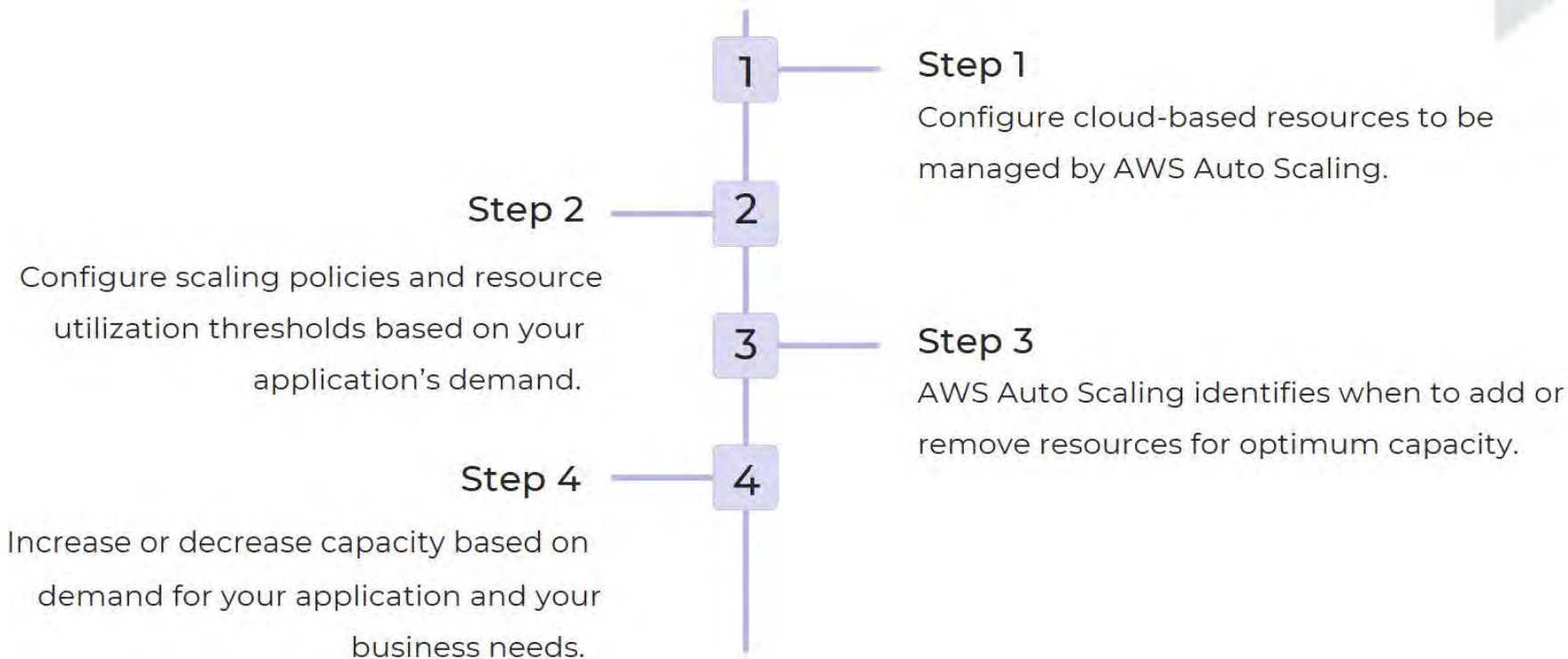


Performance and Reliability

How does AWS Auto Scaling work?



Implementation of AWS Auto Scaling



Understanding the components of Auto Scaling

Auto Scaling is an AWS service that allows you to automatically adjust the capacity of your EC2 instances based on demand. There are three main components to Auto Scaling:

- **Auto Scaling groups:** A collection of EC2 instances that are created and managed together. You can specify the minimum and maximum number of instances in the group, and Auto Scaling will automatically adjust the number of instances based on demand.
- **Launch configurations:** A template that defines the settings for new instances that are launched by the Auto Scaling group. This includes the AMI, instance type, and security groups.
- **Scaling policies:** Rules that determine when and how to adjust the capacity of the Auto Scaling group. For example, you might create a scaling policy that adds 2 instances when CPU usage exceeds 80% for 5 minutes.

By using these components together, you can ensure that your application always has the right amount of capacity to handle traffic. You can also save money by only paying for the instances that you need.

AutoScaling Groups

Defining Group Size Limits

Creating an AutoScaling Group

Configuring AutoScaling Triggers

Creating Multi-Zone Deployments

Understanding Autoscaling Groups

- An autoscaling group is a collection of Amazon EC2 instances that are designed to work together to handle incoming traffic.

When you create an autoscaling group, you specify the minimum and maximum number of

- instances that should be running at any given time.

If the traffic to your application increases, the autoscaling group will automatically add more

- instances to handle the load. If the traffic decreases, the autoscaling group will remove some of the instances to save costs. This allows your application to handle variable levels of traffic without any manual intervention.

Creating Auto Scaling groups and defining launch configurations

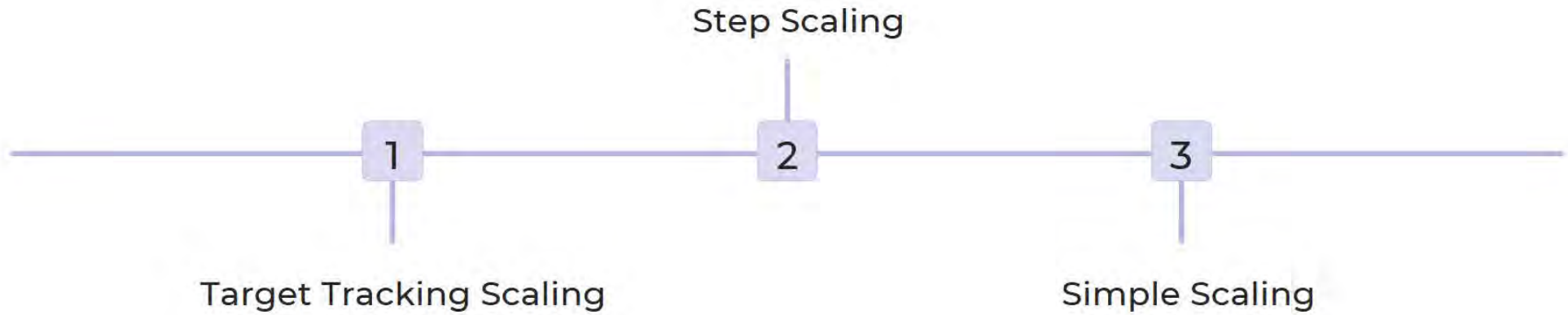
- Auto Scaling groups are a core component of AWS Auto Scaling. An Auto Scaling group contains a collection of Amazon EC2 instances that are created from a common Amazon Machine Image (AMI).
- The group automatically scales the number of instances up or down in response to changes in demand for the application. To create an Auto Scaling group, you'll need to:
 1. Create an Amazon Machine Image (AMI) that contains the software and configuration for your application.
 2. Create a launch configuration that describes the settings for the instances that will be launched by the Auto Scaling group. This includes the AMI ID, instance type, and security groups.
 3. Create an Auto Scaling group and specify the minimum and maximum number of instances that the group should maintain. You'll also need to specify the launch configuration that you created in step 2.

Configuring Autoscaling Group Capacity

- When creating an autoscaling group, you need to specify the minimum, maximum, and desired capacity. The minimum capacity is the smallest number of instances that can be running at any time. The maximum capacity is the largest number of instances that can be running at any time. The desired capacity is the number of instances that should be running at any given time.
- When the autoscaling group launches, it will start with the desired capacity. If the traffic to your application increases and exceeds the desired capacity, the autoscaling group will automatically add more instances up to the maximum capacity. If the traffic decreases and goes below the desired capacity, the autoscaling group will remove some of the instances down to the minimum capacity.
- It's important to choose appropriate values for these parameters based on the expected traffic to your application. If the minimum capacity is too high, you'll be paying for instances that you don't need. If the maximum capacity is too low, your application won't be able to handle spikes in traffic.



AutoScaling Policies



Customizing AutoScaling with Lifecycle Hooks



Extending AWS Services



Terminate Protection



Amazon **EBS**

EBS Volume Creation

Key Features and Benefits

Cost Optimization💰

Quick and Accurate
Scaling🚀

Improved Reliability🛡️

Use Cases



E-commerce applications



Mobile apps



Cloud infrastructure

Challenges and Limitations

1 Application design

2 Scaling limitations

3 Costs

Best Practices for AWS Auto Scaling

Algorithm selection

Choose the algorithm that best meets the needs of your application.

Application architecture

Consider implementing per-application scaling for optimal resource allocation.

Monitoring and testing

Continuously monitor and test to ensure AWS Auto Scaling is working optimally with your applications.

Scheduling

Use AWS Auto Scaling scheduled actions to proactively manage capacity.

AutoScaling Case Studies: Examples of Successful Implementation

Netflix

Netflix relies heavily on AutoScaling to stream content for millions of users worldwide. By dynamically allocating streaming resources as needed, Netflix ensures that users can watch their favorite shows and movies whenever they want.

Nasa JPL

NASA's Jet Propulsion Laboratory (JPL) uses AutoScaling to process data and simulate workloads for various space missions. By scaling up or down based on demand, JPL ensures that its compute resources are optimized and cost-effective.

Kajabi

Kajabi, an e-learning platform, uses predictive scaling to optimize compute resources during their peak hours. By predicting workload patterns and scaling in advance, they ensure students have smooth interaction with courses delivered on the platform.