# AUTO SCALING

Auto scaling is a cloud computing feature that automatically adjusts the number of compute resources in a server farm or application deployment based on predefined conditions such as traffic demand, resource utilization, or custom metrics. It is a key component of cloud elasticity, allowing applications to dynamically scale up or down to handle fluctuations in workload without manual intervention.

Here's how auto scaling typically works:

Monitoring: Auto scaling begins by continuously monitoring various metrics, such as CPU utilization, network traffic, or queue length. These metrics provide insights into the current workload and resource usage.

Scaling Policies: Based on the monitored metrics, auto scaling policies define rules and thresholds for scaling actions. These policies specify conditions under which the system should scale up or down.

Scaling Actions: When the monitored metrics exceed or fall below the thresholds defined in the scaling policies, auto scaling triggers scaling actions. Scaling actions can involve adding more instances or resources when demand increases (scaling out) or removing instances or resources when demand decreases (scaling in).

Dynamic Provisioning: Auto scaling dynamically provisions or de-provisions compute resources, such as virtual machines, containers, or serverless functions, to match the current workload demand. This ensures that applications can efficiently utilize resources without over-provisioning or under-provisioning.

Integration with Load Balancing: Auto scaling often works in conjunction with load balancing to distribute incoming traffic across multiple instances or resources. As instances are added or removed by auto scaling, the load balancer automatically adjusts its configuration to ensure even distribution of traffic.

Benefits of auto scaling include:

High Availability: Auto scaling helps maintain high availability by automatically adding resources to handle increased demand or replacing unhealthy instances.

Cost Efficiency: By scaling resources based on demand, auto scaling helps optimize resource usage and reduce operational costs. Resources are only provisioned when needed, avoiding over-provisioning.

Performance Optimization: Auto scaling ensures that applications can dynamically adjust to changes in workload, maintaining optimal performance and responsiveness.

Simplified Management: With auto scaling, administrators can set up scaling policies and let the system manage resource provisioning and de-provisioning automatically, reducing the need for manual intervention.

Auto scaling is a fundamental feature of cloud computing platforms like Amazon Web Services (AWS), Microsoft Azure, Google Cloud Platform (GCP), and others, enabling highly scalable and resilient applications and services.

## Get started with auto scaling using the console

Sure, here's a general overview of how you can get started with auto scaling using the AWS Management Console:

Sign in to the AWS Management Console:

Go to the AWS Management Console (https://aws.amazon.com/console/), and sign in to your AWS account.

Navigate to the Auto Scaling Service:

Once logged in, navigate to the "Services" menu at the top of the console. Under the "Compute" section, select "Auto Scaling."

Create a Launch Configuration:

In the Auto Scaling dashboard, click on "Create launch configuration." A launch configuration defines the instance type, AMI, key pair, security groups, and other configurations for the instances that auto scaling launches.

Select an Amazon Machine Image (AMI):

Choose an Amazon Machine Image (AMI) for your instances. An AMI is a pre-configured template for the root volume of the instance.

Choose an Instance Type:

Select the instance type for your auto-scaled instances. Instance types define the hardware specifications (CPU, memory, storage, etc.) of the instances.

Configure Security Groups and Key Pairs:

Configure security groups and key pairs for your instances to control inbound and outbound traffic and enable SSH access if needed.

Configure Auto Scaling Groups:

After creating a launch configuration, go back to the Auto Scaling dashboard and click on "Create Auto Scaling group." An Auto Scaling group defines the scaling policies and other configurations for your instances.

Set Auto Scaling Group Details:

Provide a name and other details for your Auto Scaling group, such as the desired capacity, minimum and maximum instances, subnets, and VPC settings.

Configure Scaling Policies:

Define scaling policies for your Auto Scaling group. You can create scaling policies based on various metrics such as CPU utilization, network traffic, or custom CloudWatch metrics.

Review and Create:

Review your configurations and click "Create Auto Scaling group" to create your auto scaling group.

Monitor Auto Scaling Activity:

Once your Auto Scaling group is created, you can monitor its activity, including scaling events, instance launches, terminations, and health checks, from the Auto Scaling dashboard.

This is a high-level overview of how to get started with auto scaling using the AWS Management Console. Keep in mind that specific configurations and options may vary depending on your requirements and AWS region. Always refer to the AWS documentation for detailed guidance and best practices.

## Maintain a fixed number of running EC2 instances

To maintain a fixed number of running EC2 instances using Auto Scaling in the AWS Management Console, you can create an Auto Scaling group with the desired number of instances and configure the scaling policies accordingly. Here's a step-by-step guide:

Sign in to the AWS Management Console:

Go to the AWS Management Console (https://aws.amazon.com/console/), and sign in to your AWS account.

Navigate to the Auto Scaling Service:

Once logged in, navigate to the "Services" menu at the top of the console. Under the "Compute" section, select "Auto Scaling."

Create a Launch Configuration:

In the Auto Scaling dashboard, click on "Create launch configuration." Follow the wizard to configure the launch configuration, including selecting an AMI, instance type, security groups, and other settings.

Create an Auto Scaling Group:

After creating the launch configuration, go back to the Auto Scaling dashboard and click on "Create Auto Scaling group." Follow the wizard to configure the Auto Scaling group:

Choose the launch configuration you created earlier:

Set the desired capacity to the fixed number of EC2 instances you want to maintain.

Configure other settings such as VPC, subnets, and health checks as needed.

Configure Scaling Policies:

Since you want to maintain a fixed number of instances, you don't need to configure scaling policies for scaling out or scaling in. Leave the scaling policies section empty.

Review and Create:

Review your configurations and click "Create Auto Scaling group" to create your group.

Monitoring:

Once your Auto Scaling group is created, you can monitor its activity from the Auto Scaling dashboard. It will ensure that the specified number of instances is always running, and if any instance fails, Auto Scaling will replace it automatically to maintain the desired capacity.

## Dynamic scaling

Dynamic scaling with Auto Scaling allows you to automatically adjust the number of EC2 instances in response to changes in demand or based on predefined conditions. Here's how to set up dynamic scaling using the AWS Management Console:

Sign in to the AWS Management Console:

Log in to your AWS account through the AWS Management Console.

Auto Scaling Service:

Go to the AWS Management Console, select "Services" from the top menu, then choose "Auto Scaling" under the "Compute" section.

Create a Launch Configuration:

If you haven't already created a launch configuration, start by creating one. A launch configuration specifies the instance type, AMI, key pair, security groups, and other settings for the EC2 instances.

Create an Auto Scaling Group:

Click on "Create Auto Scaling group." Follow the wizard to create an Auto Scaling group:

Choose the launch configuration you created earlier.

Set the initial number of instances and define the minimum and maximum number of instances you want your group to scale between.

Configure the network settings, subnets, and health checks for your Auto Scaling group.

Configure Scaling Policies:

In the Auto Scaling group creation wizard, set up scaling policies to define when and how the group should scale:

Scale Out Policy: Define conditions for scaling out, such as CPU utilization reaching a certain threshold or an increase in the number of incoming requests. Specify the number of instances to add when scaling out.

Scale In Policy: Similarly, define conditions for scaling in, such as CPU utilization dropping below a certain threshold or a decrease in the number of incoming requests. Specify the number of instances to remove when scaling in.

Review and Create:

Review your configurations and click "Create Auto Scaling group" to create your group.

Monitoring and Adjustments:

Once your Auto Scaling group is created, monitor its activity and adjust scaling policies as needed:

Use CloudWatch metrics to monitor the performance and health of your EC2 instances.

Review scaling events and alarms in the Auto Scaling console.

Adjust scaling policies based on performance patterns and workload changes to optimize resource usage and maintain performance.

With dynamic scaling configured, your Auto Scaling group will automatically adjust the number of EC2 instances in response to changes in demand, ensuring that your applications can handle varying workloads efficiently while maintaining availability and cost-effectiveness.

## The lifecycle of auto scaling

The lifecycle of auto scaling refers to the stages an instance goes through from its creation to termination within an Auto Scaling group. These stages include:

Launching: When Auto Scaling determines that additional instances are needed, it launches new EC2 instances based on the configured launch configuration. The new instances are initialized and configured according to the launch configuration's specifications.

In-Service: Once the newly launched instance passes its health checks and is ready to handle traffic, it transitions to the "In-Service" state. In this state, the instance is actively serving requests and is part of the active fleet of instances handling the workload.

Terminating: When Auto Scaling determines that it needs to scale in and remove instances, it selects instances for termination based on its termination policy. The instance enters the "Terminating" state, and any in-flight requests are allowed to complete. After that, the instance is deregistered from the load balancer (if applicable) and removed from service.

Terminated: Once the termination process is complete, the instance enters the "Terminated" state. In this state, the instance no longer exists, and its resources are released. The instance is removed from the Auto Scaling group, and its status is no longer monitored.

Throughout this lifecycle, Auto Scaling continuously monitors the health and performance of instances, adjusting the number of instances in the group as needed to maintain the desired capacity and meet the defined scaling policies. It ensures that the Auto Scaling group can dynamically scale in response to changes in demand while maintaining high availability and reliability for the application or service.

## Policies of auto scaling

Auto Scaling policies define the conditions and rules that Auto Scaling follows to automatically adjust the number of EC2 instances in an Auto Scaling group based on changes in demand, workload, or other metrics. There are two main types of Auto Scaling policies:

Scaling Policies**:**

Scaling policies define the conditions under which Auto Scaling should scale the number of instances in an Auto Scaling group. There are two types of scaling policies:

Scale Out Policy (Scaling Out): This policy defines conditions that trigger the addition of instances to the Auto Scaling group, typically in response to increased demand or workload. For example, you might create a scale-out policy that adds instances when CPU utilization exceeds a certain threshold or when the number of incoming requests surpasses a predefined limit.

Scale In Policy (Scaling In): This policy defines conditions that trigger the removal of instances from the Auto Scaling group, typically in response to decreased demand or workload. For example, you might create a scale-in policy that removes instances when CPU utilization drops below a certain threshold or when the number of incoming requests decreases.

Scheduled Actions**:**

Scheduled actions allow you to define specific times or dates when Auto Scaling should perform scaling actions, regardless of current workload or metrics. Scheduled actions are useful for predictable changes in demand, such as increased traffic during peak hours or reduced demand during off-peak periods. For example, you could schedule Auto Scaling to add more instances before a scheduled event or scale down instances during periods of low usage to reduce costs.

Auto Scaling policies can be configured using the AWS Management Console, AWS CLI, or AWS SDKs. You can attach multiple scaling policies to an Auto Scaling group and combine them with scheduled actions to create complex scaling strategies that meet the needs of your application or service.

It's important to carefully define and fine-tune Auto Scaling policies to ensure that your Auto Scaling group can efficiently and effectively handle changes in demand while maintaining performance, availability, and cost-effectiveness. Regular monitoring and adjustment of policies based on performance metrics and workload patterns are essential for optimizing Auto Scaling behavior.