```
Amazon EC2 Auto Scaling helps you ensure that you have the correct number of Amazon EC2 instances available to handle the load for your application.
                                                                 You can specify the minimum number of instances and Amazon EC2 Auto Scaling ensures that your group never goes below this size.
                                                                 You can specify the maximum number of instances and Amazon EC2 Auto Scaling ensures that your group never goes above this size.
                                                                  If you specify the desired capacity Amazon EC2 Auto Scaling ensures that your group has this many instances.
                                                                                                       Auto Scaling can detect when an instance is unhealthy, terminate it, and launch an instance to replace it.
                                                                    Better fault tolerance
                                                                                                        You can also configure Auto Scaling to use multiple Availability Zones. If one Availability Zone becomes unavailable, Auto Scaling can launch instances in another one to compensate.
                                                                                            Auto Scaling can help you ensure that your application always has the right amount of capacity to handle the current traffic demands.
                                                                     Better availability
                                                                                                                         Vertical Scaling
                                                                                                                                                                                                                                                                           Scale Out
                                                                                                                                                                           Add more nodes to a system, such as adding a new computer to a distributed software application.
                                                                                      Scaling Terminology
                                                                                                                         Horizontal Scaling
                                                                                                                                                                                                                                           Scale In
                                                                                                                                                                          Remove nodes from a system
                                                                                                                         Auto Scaling Provides Horizontal Scaling
                                                                                                                                                                                                                                                                                                                  Auto Scaling group
                                                                                                                                                                                                                                                                                                    Your EC2 instances are organized into groups so that they can be treated as a logical unit for the purposes of scaling and management. In Auto scaling Group You Can Specify
                                                                                     EC2 Autoscaling Groups
                                                                                                                                                                                                                                                                                                                                    Scale out as needed
                                                                                                                                                                                                                                                                                                 Minimum size
                                                                                                                                                                                                                                                                                                     Desired capacity
                                             Auto Scaling components
                                                                                                                                                                                                                                                                                                                         Maximum size
                                                                                                                                                                              A launch configuration is an instance configuration template that an Auto Scaling group uses to launch EC2 instances. When you create a launch configuration, you specify information for the instances. Include the ID of the Amazon Machine Image (AMI), the instance type, a key pair, one or more security groups, and a block device mapping.
                                                                                                                     Launch configurations OR Launch templates
                                                                                                                                                                                           allows you to have multiple versions of a template With versioning,
                                                                                    Launch Templates
                                                                                                                                                                                                                                                                                              current generation of EBS volume types (gp3 and io2)
                                                                                                                     Launch Templates Advantage Over Launch Configuration
                                                                                                                                                                                           launch templates to ensure that you're accessing the latest features and improvements Like
                                                                                                                                                                                                                                                                                             EBS volume tagging
                                                                                                                                                                                           you cannot create an Auto Scaling group that launches both Spot and On-Demand Instances or that specifies multiple instance types or multiple launch templates. You must use a launch template to configure these features.
                                                                                                         Manual scaling is the most basic way to scale your resources.
                                                                                                          You Need to specify the change in the maximum, minimum, or desired capacity of your Auto Scaling group.
                                                                           Manual Scaling
                                                                                                                                 To increase resources for an infrequent event, such as the release of a new game version that will be available for download and require a user registration.
                                                                                                            Means that scaling actions are performed automatically as a function of time and date
                                                                                                             This is useful when you know exactly when to increase or decrease the number of instances in your group,
                                                                           Scheduled Scaling
                                                                                                                                    periodic events such as end-of-month, end-of-quarter, or end-of-year processing, and also other predictable, recurring events.
                                                                                                           Dynamic scaling lets you define parameters that control the Auto Scaling process in a scaling policy.
                                                                                                           For example, you can create a policy that calls for enlarging your fleet of EC2 instances whenever the average CPU utilization rate stays above ninety percent for fifteen minutes.
                                                                           Dynamic Scaling
                                                                                                                           It is used to Handle Uncertain Load
                                                                                                           SINGLE THRESHOLD
                                                                         SIMPLE SCALING
                                                                                                           THRESHOLD - add 1 instance If CPU Utilization is between 50% and 60%
                                                                                                         MULTIPLE thresholds
                                                                                                         Threshold A - add 1 instance when CPU Utilization is between 50% and 60%
                                                                         STEP SCALING
                                            Scaling Policy
                                                                                                          Threshold B - add 2 instances when CPU Utilization is between 60% and 70%
                                                                                                          Threshold C - add 3 instances when CPU Utilization is between 70% and 90%
                                                                                                                       It's automatic
                                                                         TARGET TRACKING SCALING
                                                                                                                       Set the value and that's Done
                                                                                                                       Self Optimization
                                                                                                         Predictive Scaling is an advanced feature that enables your cloud infrastructure, specifically EC2 instances, to automatically scale in or out based on predicted traffic
                                                                                                         By using machine learning models to analyze historical traffic data (usually three weeks of data is a minimum requirement), Predictive Scaling accurately forecasts future demand and proactively manages the compute capacity of your Auto Scaling group.
                                                                              Introduction
                                                                                                          Predictive Scaling can be used in conjunction with dynamic scaling and scheduled scaling. While dynamic scaling responds to changes in real-time, predictive scaling prepares for expected changes, offering a more proactive approach.
                                                                                                                                                          Predictive Scaling will automatically adjust the number of instances in your Auto Scaling group based on its predictions of future traffic.
                                                                                                          Scale based on forecast
                                                                                                                                                           Predictive Scaling will still make predictions about future traffic, but it will not take any automatic actions to scale your resources.
                                                                                                                                                                            The chosen metric is what the system will analyze historically to predict future needs.
                                                                                                                                                                             These could be CPU utilization, network I/O, request count per target (for load balancers), or any other relevant metric that reflects the load or performance of your instances.
Auto Scaling
                                                                                                                                                                                            If you choose CPU utilization, the algorithm will look at past CPU usage trends to forecast future CPU requirements.
                                                                                                          Metrics and target utilization
                                            Predictive Scaling
                                                                                                                                                                                     Target Utilization refers to the desired level of resource usage that you want to maintain.
                                                                                                                                                                                     It's a threshold that guides the scaling process.
                                                                                                                                                      Target utilization
                                                                             Configuration
                                                                                                                                                                                                     If the metric is CPU utilization, setting a target utilization of 50% means you're indicating that the ideal operational state is when each server is using around 50% of its available CPU power.
                                                                                                                                              This setting configures the amount of time by which the instance launch time can be advanced
                                                                                                          Pre-launch instances
                                                                                                                                                              The forecast says to add capacity at 10:00 AM, and you choose to pre-launch instances by 5 minutes. In that case, the instances will be launched at 9:55 AM.
                                                                                                                                                                             This is an extra percentage or number of instances that you choose to add on top of the forecasted capacity.
                                                                                                          Buffer maximum capacity above the forecasted capacity
                                                                                                                                                                                                       Suppose the Predictive Scaling algorithm forecasts that you will need 10 instances to meet the demand at a certain time.
                                                                                                                                                                                                      If you set a buffer maximum capacity of 20%, it will actually provision 12 instances (10 + 20% of 10), giving you 2 extra instances as a buffer.
                                                                                                                                                                                                                                                   Updating instances to newer versions of the AMI.
                                                                                      The Instance Maintenance Policy in AWS Auto Scaling allows you to define how your Auto Scaling group replaces instances during various operations like
                                                                                                                                                                                                                                                    Replacing unhealthy instances identified by health checks.
                                                                                                                                                                                                                                                   Maintaining the desired capacity across Availability Zones.
                                            Instance Maintenance Policy
                                                                                                                   Unhealthy instances are terminated immediately, followed by new launches
                                                                                                                       New instances are launched first before unhealthy ones are terminated, ensuring minimal downtime
                                                                                       Launch before terminating
                                                                                                              Fine-tune scaling behavior for unique needs with advanced control over replacements and health checks.
                                                                                        Custom behavior
                                                                                                            The termination policy of an Auto Scaling group determines which instances the Auto Scaling group will terminate first when it needs to scale In.
                                                                                Introduction
                                                                                                           AWS provides several built-in termination policies, and you can also set up custom termination policies to suit your specific needs.
                                                                                                                                                                                                                                                     The default termination policy aims to maintain an even distribution of instances across Availability Zones. In our scenario, Ap-South-1A has more instances than the others. So, the policy targets this AZ for scale-in to balance the distribution.
                                                                                                                                                                                                    Balance Across Availability Zones
                                                                                                                                                                                                                                                     Target for Consideration: Instances A1, A2, A3 in Ap-South-1A
                                                                                                                                                                                                                                               All three instances in Ap-South-1a are checked for scale-in protection
                                                                                                                                                                                                                                               If one of them is marked as protected from scale-in actions, it's excluded from termination.
                                                                                                                                                                                                     Instance Protection Status
                                                                                                                                                                                                                                               Assuming none of these instances have scale-in protection enabled:
                                                                                                                                                                                                                                               Still in Consideration
                                                                                                                                                                                                                                                                          Instances A1 and A2 and A3
                                                                                                                                                                                                                                                              Suppose Instance A1 and A2 were launched using the same launch configuration, and Instance A3 was launched with a newer version because you updated the launch template after A1 and A2 were already running.
                                                                                                                                                                                                     Oldest Launch Configuration or Template
                                                                                                                                                                                                                                                                                                     Instances A1 and A2, because they have the oldest launch configuration.
                                                                                                                                                                                                                                                               Target for Further Consideration
                                                                                                                                                                                                                                                    The default termination policy considers which instance is closest to the next billing hour to optimize cost savings.
                                                                                                              Default Termination Policy
                                                                                                                                                                                                                                                    If Instance A1 is 50 minutes into its current hour, and Instance A2 is only 10 minutes in
                                                                                                                                                                                                     Closest to the Next Billing Hour
                                                                                                                                                                                                                                                     Terminating A1 would result in a smaller portion of unused time that you've already paid for.
                                                                                                                                                                                                                                                     Instance Selected for Termination
                                                                                                                                                                 Terminate instances in the Auto Scaling group to align the remaining instances to the allocation strategy for the type of instance that is terminating (either a Spot Instance or an On-Demand Instance).
                                                                                                                                                                                                        Terminate instances that have the oldest launch template or OldestLaunchConfiguration
                                                                                                                                 OldestLaunchTemplate Or OldestLaunchConfiguration
                                                                                                                                                                                                     This policy is useful when you're updating a group and phasing out the instances from a previous configuration.
                                                                                                                                                                              Terminate instances that are closest to the next billing hour.
                                                                                                                                                                              This policy helps you maximize the use of your instances that have an hourly charge.
                                                                                                                                  ClosestToNextInstanceHour
                                                                                                                                                                              Only instances that use Amazon Linux, Windows, or Ubuntu are billed in one-second increments.
                                                                               Other Built-in Termination Policies
                                                                                                                                                                   Terminate the newest instance in the group.
                                            Termination policy
                                                                                                                                  NewestInstance
                                                                                                                                                                  This policy is useful when you're testing a new launch configuration but don't want to keep it in production.
                                                                                                                                                                  Terminate the oldest instance in the group.
                                                                                                                                 OldestInstance
                                                                                                                                                                 This option is useful when you're upgrading the instances in the Auto Scaling group to a new EC2 instance type.
                                                                                                                                                                 You can gradually replace instances of the old type with instances of the new type.
                                                                                                                           A custom termination policy provides better control over which instances are terminated, and when.
                                                                                                                           A Lambda function can be used in conjunction with custom termination policies to provide advanced, programmable logic for determining which instances to terminate during a scale-in event.
                                                                                                                                                                                                                                                                                                                  Sometimes, the usual ways AWS decides which EC2 Instance to shut down don't match exactly what you need.
                                                                                                                                                                                                                                                                Standard Rules Might Not Fit All Needs
                                                                                                                                                                                              Better Control Over Which Instances to Terminate
                                                                                                                                                                                                                                                                                                                      Picking specific servers based on special labels (tags) you've given them,
                                                                                                                                                                                                                                                               Custom Checks for Better Choices
                                                                                                                                                                                                                                                                                                                     Looking at very specific metric of how well your application is running on each instances to decide which ones to shut down.
                                                                                                                                                                                                                                                 Regular rules for terminating EC2 instances might stop them suddenly while they're still working.
                                                                                                                                                                                              Graceful Shutdown of Applications
                                                                                                                                                                                                                                                 Your Lambda function can initiate a graceful shutdown sequence within the instance before termination
                                                                                                                                                                                                                                                 This prevents data loss, service disruptions, and improves overall application resilience.
                                                                                Custom Termination Policies
                                                                                                                                                                                                                                                                                                                                                                     Backing up critical data to other instances or external storage.
                                                                                                                           5 Solid Reasons to Use a Custom Termination Policy
                                                                                                                                                                                                                                                         You might need to perform specific actions before terminating an instance, such as
                                                                                                                                                                                                                                                                                                                                                                     Draining connections or gracefully terminating running processes
                                                                                                                                                                                              Pre-Termination Actions and Data Backups
                                                                                                                                                                                                                                                                                                                                                                     Unregistering the instance from any service discovery mechanism.
                                                                                                                                                                                                                                                         Your Lambda function can execute these pre-termination steps, ensuring data safety and a smooth transition even when scaling down.
                                                                                                                                                                                                                                                             Certain termination decisions might require access to information beyond what Auto Scaling exposes directly.
                                                                                                                                                                                              Integration with External Systems or Databases
                                                                                                                                                                                                                                                              Your Lambda function can interact with external systems or databases (e.g., for license management, custom health checks) to make informed termination decisions based on broader context.
                                                                                                                                                                                                                                                           As your app gets more complex or your needs change, you might need flexible rules for turning off EC2 instances.
                                                                                                                                                                                    5 Flexible and Dynamic Termination Logic
                                                                                                                                                                                                                                                           By customizing the logic within your Lambda function, you can easily update the termination criteria as your needs change, without code modifications to the Auto Scaling group itself.
```