Amazon EC2 instance types. ( Requirements for compute memory or storage capabilities.

**1. General purpose instances.**

provide a balance of <u>compute</u>, <u>memory</u> and <u>networking</u> resources.
uses them for a variety of workloads.

* application servers.
* gaming server
* Backend servers for enterprises
* small and medium data bases.

**2. compute optimized instance.**

* are ideal for compute-bond applications that benefit from <u>high performance processors</u>.

* use for work lode os web, application and gaming servers.

**3. Memory optimized instance.**

are designed to deliver <u>fast performance</u> for workloads that process <u>large datasets in memory</u>.

**4. Accelerated computing instance.**

* use <u>hardware accelerators</u>, or <u>coprocessors</u>, to perform some funtions more efficiently than is possible in software running on CPUs.

* Accelerated computing instances ore ideal for workloads such as graphicsapplication, game streaming and application streaming.

**5. storage optimized instances.**

designed for workloads that require <u>high, sequential read and write</u> access to <u>large datasets</u> on <u>local storage</u>.

workloods includes distributed file systems, data ware housing applications, and high-frequency online transaction processing syst.
( OLTP)

(IOPS) input/output operations per second is a metric that measure performance of a storage device.

## 1. On-Demand.

* ideal for short-term, irregular workloads that cannot be interrupted.
* No-upfront costs or mini contracts applys.
* The instances run continuously until you stop them and you pay for only the compute time you use.

## 2. Saving plans.

* enables you to reduce your compute costs by committing to a consistent amount of compute usage for a 1-year or 3-year ter.
* Saving of upto 66% or 72% over on-Demand costs.

## 3. Reserved Instances.

* are billing discount applied to use of on-Demand Instance in your accout.
* you can purchase standard Reserved and Convertible Reserved. Instence for a 1-year or 3-year terms and scheduled reserved instances for a 1-year term.
* you realise greater cost saving with the 3-year option.

## 4. Spot Instances:

* ideal for workloads with flexible start and end time, or that can withstand interruptions.
* Spot Instances use unused Amazon $EC_2$ computing capacity and offer you cost savings od up to 90% off of on-Demand prices.
* If Amazon $EC_2$ capacity is available, spot Instances lanched. or not.

## 5. Dedicated Hosts:

* are physical servers with Amazon $EC_2$ instance capacity that is fully dedicated to your use.
* Dedicated Hosts are more expensiver.
* use your existing per-socked - per-core, per-VM software licest.

# Scalability:-

* involves beginning with only the resources you need and designing your architecture to automatically respond to changing demand by scaling out or in.

* As a result, you pay for only the resources you use.

* The Aws services that provides this functionality for Amazon EC2 instances is Amazon EC2 Auto Scaling.

## Amazon EC2 Auto Scaling.

* try to access a website that wouldn't load and frequently timedout.

* Amazon EC2 Auto Scaling enables you to automatically add or remove Amazon EC2 instance in responce to changing application demand

→ Dynamic Scaling responds to changing demand

→ predictive Scaling automtically schedule the right number of amazon EC2 instances based on predicted demand

To scale faster, you can dynamic Scaling and prediction Scaling together.

* Mini. No. of Amazon EC2 instance at one.

* If you do not specify the desired number Amazon EC2 instance in an Auto Scaling group, the desired capacity defaults to your mini. capacity.

* Max. capacity.

# Elastic Load Balancing

* is the Aws services that automatically distributes incoming application traffic across multiple rescources such as Amazon EC2 instance.

* A load Balancer acts as a Single point of contact for all incoming web traffic to your Auto Scaling group.

* these request route to the load balancer frist, distribues the workload across the multiple instance so that no single instance has to carry bulk of it.

* Although Elastic Load Balancy and Amazon EC2 Auto Scaly are separate Services, they work together.