

---

---

---

---

---



# Mod1: Intro to AWS

What is a Client-Server model?



Learning Objectives:

Benefits of AWS

On-demand delivery vs Cloud deployments

Pay-as-you-go pricing model

What is Cloud Computing?

- On demand delivery
- Over the internet
- Pay-as-you-go pricing

3 different deployment models

- Cloud-Based Deployment

- run all parts of application in cloud
- Migrate existing applications → Cloud
- Design + Build new applications in the Cloud
- Can build low level or high-level infrastructure

- On-premises Deployment

- Deploy resources using virtualization + resource management tools
- increase resource utilization w/ application management & virtualization technologies
- aka "private cloud" deployment

- Hybrid Deployment

- Connect cloud-based resources ↔ On-premises infrastructure
- integrate cloud-based resources w/ legacy IT applications

# Mod2: Compute in the Cloud

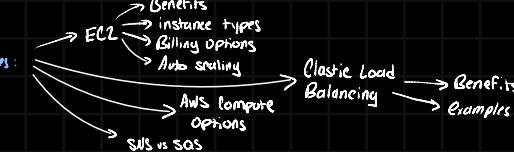
Amazon Elastic Compute Cloud

- EC2 → Secure, resizable compute capacity in cloud as EC2 instances
- Avoid overhead of local servers
  - hardware costs, wait time, installation, configurations
- Virtualization technology w/ EC2
  - isolate VMs
  - Share resources ↗ Multitenancy
  - Keep data secure
  - launch instances in minutes
  - Pay only for what you use
- Vertical Scaling
  - Allocate more memory & CPU when needed
- Control network requests for server

Amazon EC2 Instance Types

- Optimized for different tasks
  - ↳ Consider computing, memory, & storage
  - ↳ Each type grouped under instance family
- General Purpose Instances
  - provide balance of computing, memory, storage & network resources
  - Applications
    - ↳ application servers
    - ↳ gaming servers
    - ↳ backend servers
    - ↳ databases
  - Use if resource needs for compute, memory, & network are similar
- Compute Optimized Instances
  - Ideal for compute-bound applications / high-performance processes
  - can use for batch processing (many transactions)
  - high-performance web servers, compute-intensive app servers, dedicated gaming servers

Learning Objectives:



① Launch

- Select config (Op Sys, Application, Instance type, Hardware type)
- Specify security settings (to control network traffic)

② Connect

- Connect to instance

- multiple ways
- exchange data
- users connect by logging in

③ Use

- run commands
- add software
- add/move files
- more...

## - Memory Optimized instances

→ Fast performance memory intensive tasks → processing large datasets

→ processing large amounts of data

## - Accelerated Computing Instances

→ good for floating point number calculations  
graphics processing  
data pattern matching

→ use hardware accelerators or co-processors to perform functions efficiently

→ More efficient than software running on CPU

→ Hardware accelerator: component that expedites data-processing

## - Storage Optimized instances

→ Good for work needing high performance w/ locally stored data  
high sequential access speeds

→ Input/Output operations per second (IOPS):

↪ metric for performance of storage device

## Amazon EC2 Pricing

### - Multiple Purchase Options

#### On-demand

- Good for short term, irregular workloads
- No upfront costs or min costs
- Only pay for what you use

#### Reserved Instances

- Predictable usage (month-to-month)
- 1 or 3 yr terms w/ 3 payment options
  - ① Upfront
  - ② Partial upfront
  - ③ No upfront

#### Spot Instances

- Flexible start & end times of interruptions
- Spare computing EC2
  - ↪ Amazon can always reclaim EC2, so work can be interrupted

#### Savings Plan

- lower prices w/ commitment to consistent usage of 1-3 yr term
- Usage beyond commitment is charged at on-demand rate

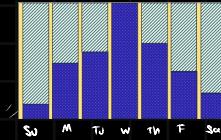
#### Dedicated Hosts

- Physical Servers fully dedicated to your use

## Scaling Amazon EC2



→ Can't handle peak hours  
or  
You waste resources at non-peak hours



Su M Tu W Th F Sa

- Scalability: beginning w/ only resources you need & designing architecture respond to demand by scaling in or out

↪ Only pay for usage

↪ Amazon provides this service

- Amazon EC2 Auto Scaling automatically adds & removes EC2 instances in response to demand

- Dynamic Scaling: responds to changing demand

→ Can use together for more speed!!!

- Predictive Scaling: Sets schedule based on predicted demand

- ↑↑↑↑ Application Availability

- Scaling UP

- More power to machines VS  
that are running

- Scaling Out

- Creating more instances

→ right amt of power for each process

### - Auto Scaling Group

→ Minimum Capacity: launched at creation

→ Desired Capacity: optional, default capacity

→ Maximum Capacity: max instances that can be used

## Elastic Load Balancing

- AWS Service that distributes incoming application traffic across multiple resources

### - Load Balancer

- Single Point of Contact for all web traffic → Auto Scaling Group

- As EC2 instances are added/removed

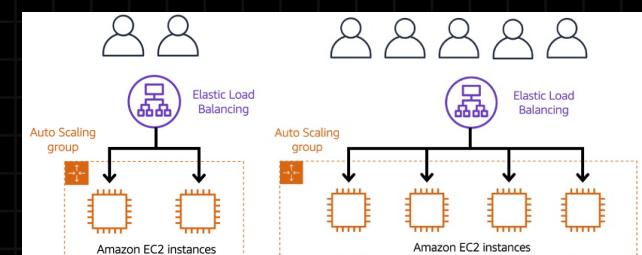
↳ requests route to load balancer

↳ requests spread across multiple resources that will handle them

- No single instance carries bulk of workload

- Auto Scaling & Load Balancing are SEPARATE services

↪ Work together → high performance availability



Ex: Coffee shop employee directs customers to which register they go to (ELB)

More registers are opened as traffic is increased and more customers (Auto Scaling)

Each register is (EC2 instance)

# Messaging & Queuing

## - Monolithic Applications & Microservices

- application components communicate w/ each other
- Monolithic Application - App w/ tightly knit components

↳ if one component fails, entire app may fail

- Microservices Approach - Components in app are loosely coupled

↳ Maintains app availability when component fails

↳ Uses services & components

## - Amazon Simple Queue Service (SQS)

- message queue service

- send, store, receive messages btw software components



Simple Queue Service (SQS)

Simple Notification Service (SNS)