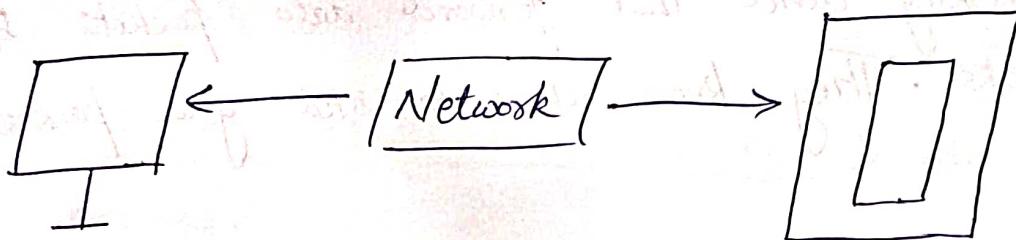


CLOUD COMPUTING

How Website Work?



Client → Network ← Server

We have a server hosted somewhere, and we as a web browser want to get access to that server to visualize a website.
 (Use a network)

The browser/client/us will find the network & will use network to route the packets, the data to the server & then the server will respond to us and we can view a website.

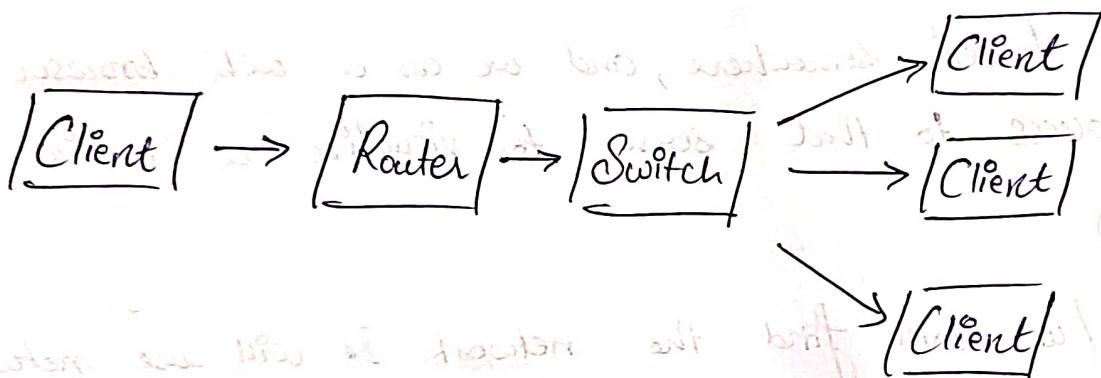
- Client and servers both have IP addresses to locate each other

What is a server composed of?

- Compute : CPU (Computations / Calculations)
- Memory : RAM (Very fast, helps store info & quickly get)
- Storage : Data (Storing data somewhere)
- Database & Store data in a structured way.
- Network : Routers, switch, DNS Server

Important Terms

- ① Network : cables, routers & servers connected with each other
- ② Router : A networking device that forwards data packets b/w computer networks. They know where to send your packets on the internet
- ③ Switch : Takes a packet & send it to correct server/client on your network



Problems With Traditional IT Approach

- ① Pay for rent for data center
- ② Pay for power supply, cooling, maintenance
- ③ Adding & replacing hardware takes time
- ④ Limited scaling
- ⑤ Hire 24/7 team to monitor infrastructure
- ⑥ Disasters (earthquake, fire, -)

What is Cloud Computing?

- ① Cloud Computing is the on-demand delivery of computer power, database storage, applications & other IT Resources.
- ② Through a service (cloud) platform with pay-as-you-go pricing.
- ③ You can provision exactly the right type & size of computing resources you need.
- ④ You can access as many resources as you need, almost instantly.
- ⑤ Simple way to access servers, storage, databases & a set of application services.

→ Amazon Web Service (AWS) owns & maintains the network connected hardwares required for these applications - services, while you provision & use what you need via a web application.

Cloud Services

- ① GMAIL : Email Cloud Service
- ② DROPBOX : Cloud Storage Service
- ③ NETFLIX : Built on Aws (Video on demand)

THE Deployment Models of the Cloud

Private Cloud

- Cloud services used by a single organization, not exposed to the public
- Complete Control
- Security for sensitive applications
- Meet specific business needs.

Public Cloud

- Cloud services owned & operated by a third party cloud service provider delivered over the Internet
Eg → Microsoft Azure, Google Cloud, AWS etc

Hybrid Cloud

- Keep some servers on premises & extend some capabilities to the Cloud
- Hybrid of our own infrastructure & the cloud (AWS)
- Control over sensitive assets in your private infrastructure
- Flexibility & cost-effectiveness of the public cloud

Characteristics of Cloud Computing

- ① On-demand self service → Users can provision resources & use them without human interaction from the service provider
- ② Broad Network Access → Resources available over the network & can be accessed by diverse client platforms
- ③ Multi-tenancy & resource pooling →
 - ① Multiple customers can share the same infrastructure & applications with security & privacy
 - ② Multiple customers are serviced from the same physical resources
- ④ Rapid elasticity & scalability →
 - ① Automatically & quickly acquire and dispose resources when needed
 - ② Quickly and easily scale based on demand
- ⑤ Measured Service → Usage is measured, users pay correctly for what they have used.

Advantages of Cloud Computing

- ① Trade Capital Expense (CAPEX) for operational expense (OPEX)
 - ① Pay on demand → Don't own hardware
 - ② Reduced total cost of ownership (TCO) & Operational expense (OPEX)
- ② Benefit from massive economies of scale
 - ① Prices are reduced as AWS is more efficient due to large scale

- ③ Stop guessing capacity
 - Scale based on actual measured usage.
- ④ Increase speed & agility
- ⑤ Stop spending money running & maintaining data centers

Problems Solved By Cloud

- ① Flexibility
- ② Cost-effectiveness
- ③ Scalability
- ④ Elasticity
- ⑤ High availability & fault tolerance
- ⑥ Agility

Types Of Cloud Computing

① Infrastructure as a Service (IAAS)

- Provides Building Blocks for Cloud IT
- Provides networking, computers, data storage space
- Highest level of flexibility
- Easy parallel with traditional on-premises IT

Eg: EC2 GCP, Azure, Linode

② Platform as a Service (PaaS)

- Removes the need for your organization to manage underlying infrastructure
- Focus on deployment & management of your applications.
Eg - Elastic Beanstalk, Heroku, Google App Engine, Azure

③ Software as a Service (SaaS)

- Completed product that is run and managed by the service provider
- Eg - Google Apps (Gmail), Dropbox, Zoom

AWS Regions

- Aws has regions all around the world.
- Names can be us-east-1, eu-west-3, etc.
- A region is a cluster of data centers
- Most Aws services are region-scoped

How to choose an aws region?

- Compliance with govt.
- Proximity to customers
- Available services within a region
- Pricing

AWS Availability Zones

- Each region has many availability zones (usually 3, min is 3, max is 6). Eg - ap-south-1a, ap-south-1b, ap-south-1c
- Each availability zone (AZ) is one or more discrete data centers with redundant power, networking & connectivity
- They are separated from each other, so that they are isolated from disasters
- They are connected with high bandwidth, ultra low latency networking.

AWS Points of Presence

- ① Amazon has 400+ Point of Presence (400+ Edge locations and 10+ Regional Cache) in 90+ cities across 40+ countries.
- ② Content is delivered to end user with low latency.

Tour of AWS Console

- ① AWS has global services
 - ① Identity and Access Management (IAM)
 - ② Route 53 (DNS)
 - ③ CloudFront (Content Delivery Network)
 - ④ WAF (Web Application Firewall)
- ② Most AWS Services are region-scoped.
 - ① Amazon EC2 (IaaS)
 - ② Elastic Beanstalk (PaaS)
 - ③ Lambda (FaaS)
 - ↳ function as a service
 - ④ Rekognition (SaaS)

IAM : Users & Groups

- ① IAM → Identity and Access management, Global Service
- ② Root account created by default, shouldn't be used or shared
- ③ Users are people within your organization, and can be grouped
- ④ Groups ^{only} contain users, not other groups
- ⑤ Users don't have to belong to a group, and users can belong to multiple groups.

IAM : Permissions

- ⑥ Users or groups can be assigned JSON documents called policies
- ⑦ These policies define the permissions of the users
- ⑧ In AWS, you apply the least privilege principle = Don't give more permissions than a user need

IAM : Password Policy

IAM > Account Settings > Edit Password Policy

- ⑨ Strong password = higher security for your account
- ⑩ In AWS, you can setup a password policy:
 - Set a minimum password length
 - Require specific character types
 - Including uppercase letters
 - Lowercase letters
 - numbers
 - non-alphanumeric characters (?) --)

- Allow all IAM users to change their own password
- Require users to change their passwords after sometime (password expiration)
- Prevent Password re-use (don't use previously used pass)

MULTI-FACTOR AUTHENTICATION (MFA)

- Login > Security Credentials > Assign MFA > Install app on your other device (MFA can also be removed)
- ① Users have access to your account & can possibly change configuration or delete resources in your AWS account.
- ② We want to protect our root accounts & IAM user
- ③

MFA = Password you know + Security device you own
- ④ Main benefit of MFA is if a password is stolen or hacked, the account is not compromised

MFA device options in AWS

- ⑤ Virtual MFA device → Support for multiple token on a single device. Eg: Google Authenticator (phone only)
- ⑥ Autify (multi-device)
- ⑦ Universal 2nd factor (U2F) Security Key → Support for multiple root & IAM users using a single security key
Eg: Yubikey by yubico (3rd party device)

- Allow all IAM users to change their own password
- Require users to change their passwords after sometime (password expiration)
- Prevent password re-use (don't use previously used pass.)

MULTI-FACTOR AUTHENTICATION (MFA)

- Login > Security Credentials > Assign MFA > Install app on your other device (MFA can also be removed)
- ① Users have access to your account & can possibly change configuration or delete resources in your AWS account.
- ② We want to protect our root accounts & IAM user
- ③

MFA = Password you know + Security device
 (you own)

- ④ Main benefit of MFA
 if a password is stolen or hacked, the account is not compromised

MFA device options in AWS

- ⑤ Virtual MFA device → Support for multiple tokens on a single device. Eg → Google Authenticator (phone only)

Physical MFA Authy (multi-device)

- ⑥ Universal 2nd factor (U2F) Security key → Support for multiple root & IAM users using a single security key
 Eg → Yubikey by yubico (3rd party device)

- ① Hardware key fob MFA Device → provided by Gemalto (3rd party)
- ② Hardware key fob MFA Device for AWS GovCloud → provided by SurePass ID (3rd party)

How can users access AWS?

- ① To access AWS, you have 3 options

→ AWS Management Console

Protected by password + MFA (optional)

→ AWS Command Line Interface (CLI)

Protected by access keys

→ AWS Software Development Kit (SDK)

for code: protected by access keys

- ① Access keys are generated through AWS console

- ② Users manage their own access keys

- ③ Access keys are secret, just like a password

- ④ Access key ID ≈ username] (Just like this!)

Access key ≈ password

↳ Secret

AWS CLI

- ⑤ A tool that enables you to interact with AWS service using commands in your command line shell

- ⑥ Direct access to the public APIs of AWS services

- ⑦ Alternative for using AWS management console

Aws SDK

- ① Aws Software Development Kits
- ② Language Specific APIs (Set of Libraries)
- ③ Enables you to access & manage AWS services programmatically.
- ④ Embedded within your application
- ⑤ Supports
 - SDKs (Javascript, Python, PHP, .Net, Ruby, Java, GO, Node.js, C++)
 - Mobile SDKs (iOS, Android, --)
 - IoT Device SDKs (Embedded C, Arduino --)
- ⑥ Example :- Aws CLI is build on Aws SDK for python

⇒ Instead of using the Aws management console, we can also use the Aws CloudShell (`[>-]`)

Commands (CloudShell & CLI)

- ① aws
- ② aws --version
- ③ aws iam list-users
- ④ ls
- ⑤ echo "test" > demo.txt Add "text" as content in demo.txt
- ⑥ cat demo.txt
- ⑦ pwd (present working directory)
(Download / Upload files using Actions button)

IAM Roles for Services

- ① Some AWS services will need to perform actions on your behalf
- ② To do so, we will assign permissions to AWS services with IAM Roles
- ③ Common Roles:
 - EC2 Instance Roles
 - Lambda Function Roles
 - Roles for CloudFormation

IAM Security Tools

IAM > Users > User > Access Advisor

① IAM Credentials Report (Account-level)

A report that lists all your account's users and the status of their various credentials.

② IAM Access Advisor (User-level)

- Access advisor shows the service permissions granted to a user and when those services were last accessed
- You can use this information to revise your policies.

SHARED RESPONSIBILITIES MODEL FOR IAM

AWS

YOU

- ① Infrastructure (global network security)
- ② Configuration & vulnerability analysis
- ③ Compliance validation
- ④ Users, groups, roles, policies management & monitoring
- ⑤ Enable MFA on all device accounts
- ⑥ Rotate all your keys often
- ⑦ Use IAM tools to apply appropriate permission
- ⑧ Analyze access patterns & review permissions

IAM Section (Summary)

- ① Users & mapped to a physical user, has a password for aws console
- ② Groups & contains users only, groups cannot contain another groups
- ③ Policies & JSON document that defines permission for users or groups
- ④ Roles & for EC2 Instances or AWS services
- ⑤ Security & MFA + password policy
- ⑥ AWS CLI & manage your aws service using command line
- ⑦ AWS SDK & " " " " " programming language
- ⑧ Access Keys & access aws using CLI or SDK
- ⑨ Audit & IAM Credentials Reports & IAM Access advisor

IAM Credentials Report lists all your Account's users and the status of their various credentials. The other IAM Security tool is IAM Access Advisor. It shows the service permission granted to a user & when those services were last accessed.

EC2 - Elastic Compute Cloud (IAAS)

AWS Budget Setup → AWS budgets lets you set custom cost budgets that alert you when your budget thresholds are exceeded

Amazon EC2

- ① EC2 is one of the ^{most} popular of AWS' offering
- ② EC2 = Elastic Compute Cloud = IAAS
- ③ It mainly consists in capability of
 - Renting virtual machines (EC2)
 - Storing data on virtual drives (EBS)
 - Distributing load across machines (ELB)
 - Scaling the services using an auto scaling group (ASG)
- ④ Knowing EC2 is fundamental to understand how the cloud works

EC2 sizing & configuration options

- ① Operating System (OS) : Linux, Windows or Mac OS
- ② How much compute power & cores (CPU)
- ③ How much RAM
- ④ How much storage space
 - Network attached (EBS & EFS)
 - Hardware (EC2 Instance store)
- ⑤ Network speed : speed of the card, public IP addresses
- ⑥ Firewalls & Security group
- ⑦ Bootstrap script (config. at first launch) : EC2 user data

EC2 - Elastic Compute Cloud (IAAS)

AWS Budget Setup → AWS budgets lets you set custom cost & usage budgets that alert you when your budget thresholds are exceeded

Amazon EC2

- ① EC2 is one of the ^{most} popular of AWS' offering
- ② EC2 = Elastic Compute Cloud = IAAS
- ③ It mainly consists in capability of
 - Renting virtual machines (EC2)
 - Storing data on virtual drives (EBS)
 - Distributing load across machines (ELB)
 - Scaling the services using an auto scaling group (ASG)
- ④ Knowing EC2 is fundamental to understand how the Cloud works

EC2 sizing & configuration options

- ① Operating System (OS) : Linux, Windows or Mac OS
- ② How much compute power & cores (CPU)
- ③ How much RAM
- ④ How much storage space
 - Network attached (EBS & EFS)
 - Hardware (EC2 Instance store)
- ⑤ Network speed & speed of the card, public IP addresses
- ⑥ Firewalls & Security group
- ⑦ Bootstrap script (Config, at first launch) : EC2 user data

EC2 User Data

(EAT) (hands) (informed) (standard - 3)

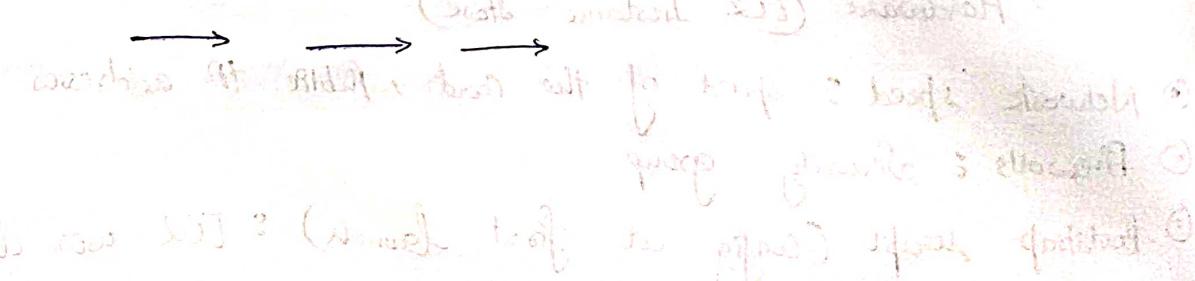
- ① It is possible to bootstrap our instances using an EC2 User data script
 - ② Bootstrapping means launching commands when a machine starts
 - ③ The script is only run once at the instance first start
 - ④ EC2 user data is used to automate boot tasks such as
 - Installing updates
 - " " softwares
 - Downloading common files from the Internet
 - Anything you can think of
 - ⑤ The EC2 user data script runs with root user
- Eg - t2.micro
t2.xlarge
- (a) from which class we have to choose with file -

EC2 Instance Types

AWS has the following naming convention

- m - instance class
- 5 - generation
- &large - size within instance class

EC2 instance types are -



① General Purpose → ○ Great for a diversity of workloads such as web servers or code repositories

○ Balance between

→ Compute

→ Memory

→ Networking

○ In the course, we use ~~Euler~~ ~~TensorFlow~~, which is a general purpose ~~EC2~~ instance

② Compute Optimized → ○ Great for compute-intensive tasks that require high performance processors

→ Batch processing workloads

→ Media transcoding

→ High performance web servers

→ High performance computing (HPC)

→ Scientific modeling & machine learning

→ Dedicated gaming servers

③ Memory Optimized → Fast performance for workloads that process large data sets in memory.

○ Use cases

→ High performance relational/non-relational databases

→ Distributed web scale cache stores

→ In-memory databases optimized for BI (Business Intelligence)

→ Applications performing real time processing of big unstructured data

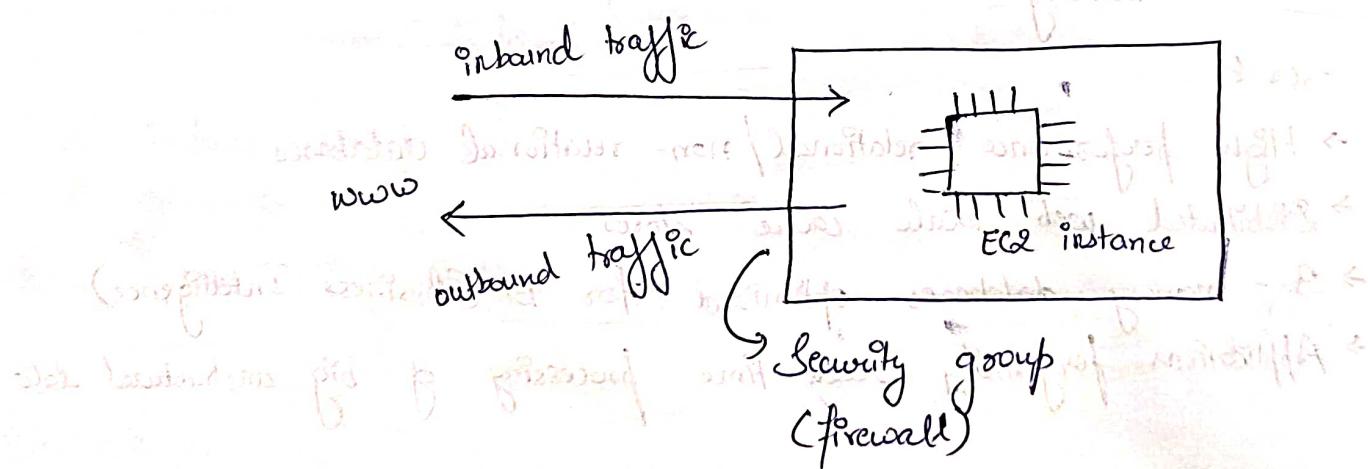
④ Storage Optimized → Great for storage-intensive tasks that require high, sequential read & write access to large data sets on local storage.

⑤ Use cases →

- High frequency online transactions processing (OLTP) systems
- Relational & NoSQL databases
- Cache for in-memory databases (Eg., Redis)
- Data warehousing applications
- Distributed file system

SECURITY GROUPS

- ⑥ Fundamental of network security in AWS
- ⑦ They control how traffic is allowed into or out of our EC2 instances
- ⑧ Security groups only contain allow rules
- ⑨ Security groups rules can reference by IP or by security group

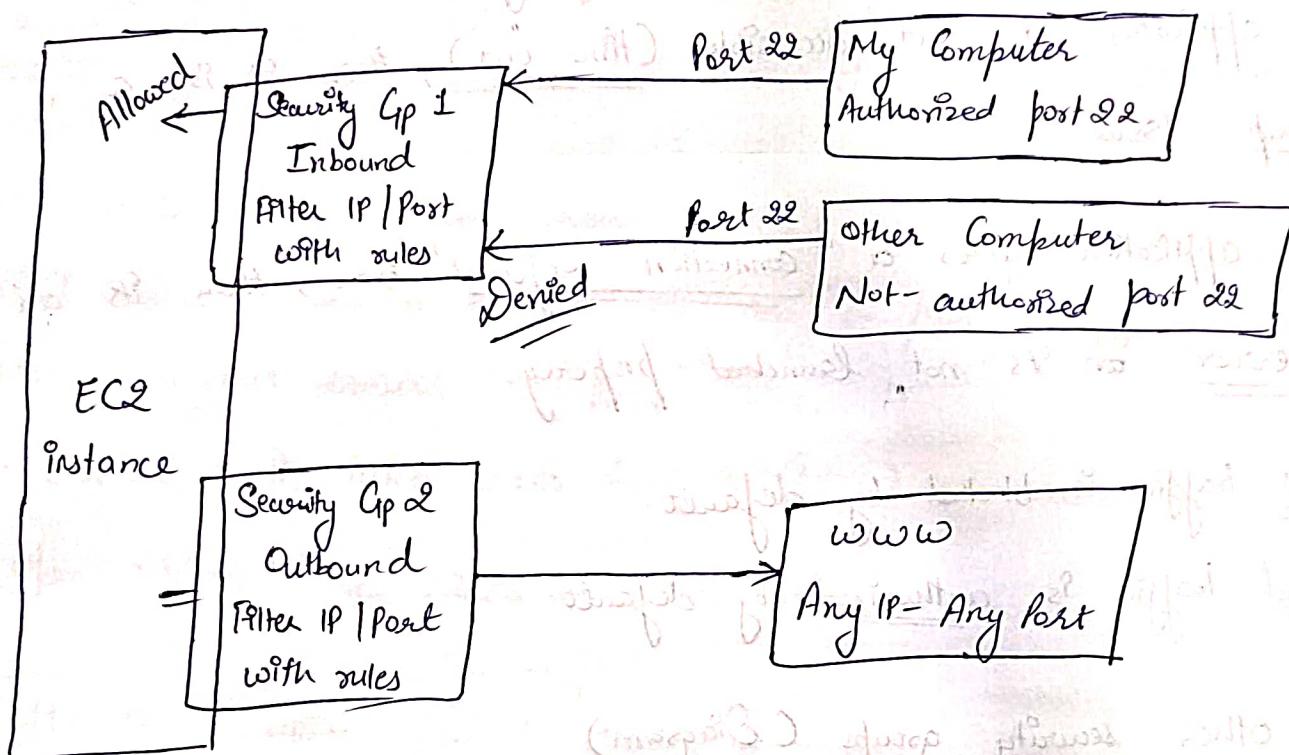


⑩ Security group acts as a firewall on EC2 instances.

① They regulate:

- access to ports
- Authorized IP ranges - IPv4 or IPv6
- Control of inbound network
- " outbound network

Diagram



Our Computer is authorized on port. (say 22) to our EC2 instance

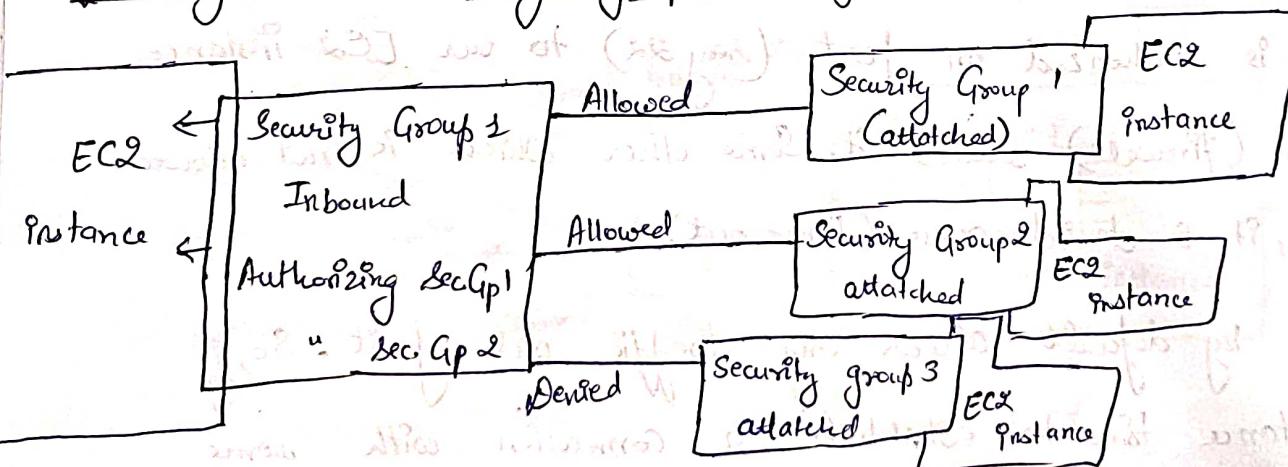
Security group (firewall) allows it. Some other which is not allowed at port 22, it is denied access (time-out)

Security group by default allows any traffic out of it. So, if our EC2 instance tries to establish a connection with some website, our security group allows it

Security Groups - Good to Know

- ① Can be attached to multiple instances
- ② Locked down to a region/VPC combination (present in one VPC, not in other)
- ③ Does 'live' outside the EC2 - if the traffic is blocked the EC2 instance won't see it
- ④ It's good to maintain one separate security group for SSH access
- ⑤ If your application is not accessible (time-out), then it's a security group issue
- ⑥ If your application gives a "connection refused" error, then it's an application error or it's not launched properly.
- ⑦ All inbound traffic is blocked by default.
- ⑧ All outbound traffic is authorized by default

Referencing other security groups (Diagram)



Ports

- ① 22 = SSH (Secure Shell) - log into a linux instance
- ② 21 = FTP (File Transfer Protocol) - upload files into a file share
- ③ 22 = SFTP (Secure FTP) - upload files using SSH
- ④ 80 = HTTP - access unsecured websites
- ⑤ 443 = HTTPS - access secured websites
- ⑥ 3389 = RDP (Remote Desktop Protocol) - log into a window instance

⇒ We can't see the text when we click on open address under Public IPv4 address (timeout), but if we copy & paste that respective IP address in our browser, it works fine.

Ctrl + C → It takes uses port 80 (HTTP) not secured

Open address → uses HTTPS (secured) It doesn't work

SSH

→ It allows you to control a remote machine, all using the Command Line.

Shared Responsibility Model for EC2

AWS

YOU

- ① Infrastructure (global network security)
- ② Isolation on physical hosts
- ③ Replacing faculty hardware
- ④ Compliance validation
- ④ Security group rules
- ④ Operating system patches & updates
- ④ Software & utilities installed on EC2 instance.
- ④ IAM Roles assigned to EC2 & IAM user access management
- ④ Data security on your instance

EC2 Instance Pricing Options

- ① On-Demand Instances → Short workload, predictable pricing, pay as you go
- ② Reserved (1 & 3 years)
 - Reserved Instances: long workload
 - Convertible Reserved Instance: long workload with flexible instances
- ③ Savings Plan (1 & 3 years) → Commitment to an amount of usage, long workload
- ④ Spot Instances → Short workloads, cheap, can loose instance (less reliable)
- ⑤ Dedicated Hosts → book an entire physical server, control instance placement
- ⑥ Dedicated Instance → No other customer will share your hardware
- ⑦ Capacity Reservation → Reserve capacity in a specific Availability zone for any duration

On-Demand

- ① Pay for what you use
 - Linux or Windows: Billing per second, after first min
 - All other OS: Billing per hour
- ② Has the highest cost but no ~~upfront~~ payment.
- ③ No long term commitment
- ④ Recommended for short workload, when you can't predict how the application will behave.

Reserved

- ① Up to 72% discount compared to On-Demand
- ② You reserve a specific instance attribute
 - ↳ (Instance type, Region, Tenancy, OS)
- ③ Reservation Period is 1 year (tds) or 3 yrs (+tds)
- ④ Recommended for steady-state usage applications
- ⑤ You can buy & sell in reserved instance marketplace

Convertible Reserved Instance

- ① Up to 66% discount
- ② Can change the EC2 instance type, instance family, OS, scope & tenancy

Saving Plan

- ① Get a discount based on long term usage (Up to 72% - same as Reserved)
- ② Commit to a certain type of usage
- ③ Usage beyond EC2 saving plan is billed at the on demand price
 - Instance size (e.g. m5.2xlarge)
 - OS (Linux / Windows)
 - Tenancy (Host, dedicated, default)

Spot Instances

- ① Get a discount of 90% compared to on-Demand
- ② Instances that you can "lose" at any point of time if your max price is less than the current spot price
- ③ Most cost-efficient instances in AWS

- ① Useful for workloads that are resilient to failure
 - Batch Jobs
 - Data Analysis
 - Image processing
 - Any distributed workloads
 - Workloads with a ~~fixed~~ flexible start & end time

- ② Not suitable for critical jobs or databases.

Dedicated Hosts

- ① A physical server with EC2 instance capacity fully reserved (dedicated) to your use
- ② Allows you to address compliance requirements & use your existing server-based software licenses (per-socket, per-core, per-VIN software licenses)
- ③ Purchasing options
 - On-demand - pay fee sec. for active dedicated hosts
 - Reserved (1 or 3 yr)
- ④ Most expensive option
- ⑤ Useful for software that have complicated licensing model (BYOL - Bring your own license)
- ⑥ Or for companies that have strong or compliance needs.

Dedicated Instances

- ① Instances run on hardware that's dedicated to you
- ② May share hardware with other instances in same account
- ③ No control over instance placement
- ④ Can move hardware after stop/start

Capacity Reservations

- ① Reserve on demand instances capacity in a specific AZ for any duration
- ② You always have access to EC2 capacity when you need it.
- ③ No firm commitments (create / cancel anytime), no billing discounts
- ④ Combine with Regional Reserved Instances & saving plan to benefit from billing discounts.
- ⑤ You are charged at on-demand rate whether you run instances or not
- ⑥ Suitable for short-term, uninterrupted workloads that needs to be in a specific AZ.

EC2 - Summary

- ① EC2 instance : AMI (OS) + Instance size (CPU + RAM) + Storage + Security groups + EC2 user data
- ② Security groups & Firewall attached to the EC2 instance
- ③ EC2 user data & Script launched at first start of an instance
- ④ SSH & start a terminal into our EC2 instances (port 22)
- ⑤ EC2 Instance Role & Link to IAM Roles.
- ⑥ Purchasing Options : On demand, spot, Reserved (Standard + Convertible), Dedicated host, Capacity Reservations, Dedicated Instances

Types of Storage you can attach to EC2 instance

- ① EBS (Elastic Block Store)
- ② Instance Store
- ③ EFS (Elastic File System)

- ① When you start instance, it changes its Public IP, to have fix public IP, you need elastic IP
- ② You can have 5 elastic IP in your acc.

EBS Volume

- ① An EBS (Elastic Block Store) volume is a network drive you can attach to your instances while they run
- ② It allows your instances to persist data, even after their termination
- ③ They can only be mounted to one instance at a time (at the CCP level)
 - | Except I/O1/Io2 family | → Certified Cloud Practitioner
- ④ They are bounded to specific A2 instances.
- ⑤ 30GB of free EBS storage of type General Purpose (SSD)
- ⑥ It's a network drive (ie not a physical drive)
 - It uses the network to communicate the instance, which means there might be a bit of delay
 - It can be detached from an EC2 instance & attached to another one quickly
- ⑦ We can have one EBS attached to one EC2 instance
- ⑧ We can have multiple EBS attached to one EC2 instance
- ⑨ No single EBS can be attached to 2 instances
- ⑩ An EBS can be left idle (un-attached)

EBS - Delete On Termination Attribute

- ⑪ Controls the EBS behavior when an EC2 instance terminates
 - By default, the root EBS volume is deleted (attribute enabled)
 - By default, any other attached EBS volume is not deleted (attribute disabled)
- ⑫ This can be controlled by AWS Console/AWS CLI

① Use Case 1

Preserve root volume when instance is terminated

⇒ EC2 instance > Storage

EBS > Volume > Create a volume

2 GiB

gp2

AZ → has to be same as of instance
AZ → us-east-1a

actions > attach volume > select instance
(EBS)

&

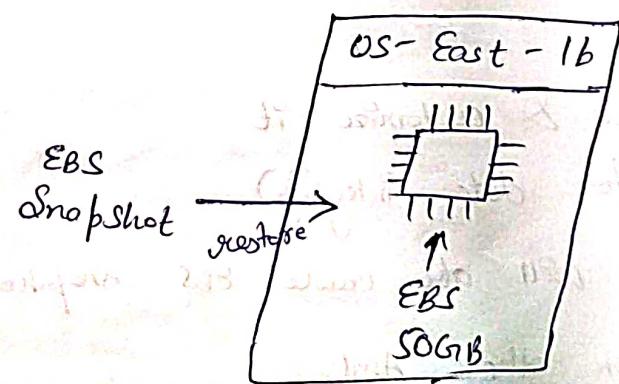
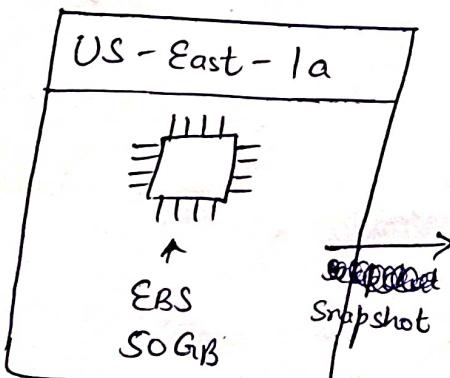
EBS is attached to our instance
[Create multiple EBS and attach them to one instance]

On default EBS, delete on terminates → Yes

Terminate instance → that EBS will be gone

EBS Snapshot

- ② Make a backup (snapshot) of your EBS volume at a point in time
- ③ Not necessary to detach volume to do snapshot but recommended
- ④ Can copy snapshots across AZ or regions



Features

① EBS Snapshot Archive

- Move a snapshot to an "archive tier" that is 75% cheaper
- Takes within 24 to 72 hours for restoring the archive

② Recycle bin for EBS snapshots

- Setup rules to retain deleted snapshots so you can recover them after an accidental deletion
- Specify retention (from 1 day to 1 yr)

AMI Overview

① AMI - Amazon Machine Image

② AMI are a customization of EC2 instances

- You add your own software, configuration, operating system, monitoring
- Faster boot / configuration time because all your software is pre-packaged

③ AMI are built for a specific region (and can be copied across regions)

④ You can launch EC2 instance from

- A Public AMI or AWS provided

- Your own AMI if you make & maintain them yourself

- the AWS Marketplace AMI (someone else made & potentially sells)

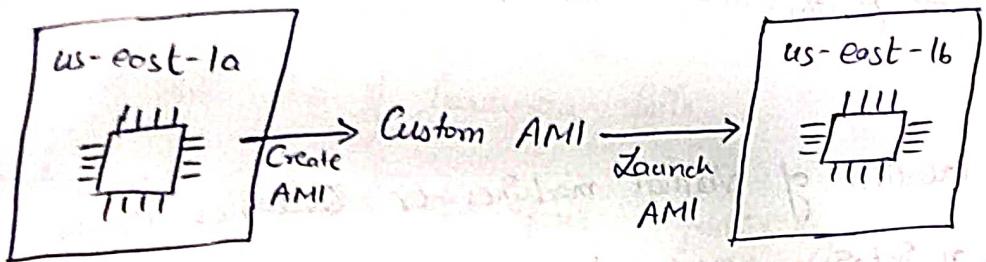
AMI Process

① Start an EC2 instance & customize it

② Stop the instance (for data integrity)

③ Build an AMI → this will also create EBS snapshot

④ Launch instances from other AMIs.



Create AMI Yourself

Create an instance, place the below code in user data (Advanced Settings)

```
#!/bin/bash
# Use this for your user data (script top to bottom)
# Install httpd (Linux or version)
```

```
yum update -y
yum install -y httpd
systemctl start httpd
systemctl enable httpd
```

(Installs Apache Server)

Instance will run through EC2 script. Wait for some time to copy paste Public IPv4 because we need to give it sometime even if it says running for EC2 user data script to run for first time.

AMI - Give the state for our EC2 instance

Right Click the instance > Image & Template > Create Image
name: Demo Image
Create Image.

AMIs > Demo Image is being created

Launch Instances from AMI

name: From AMI

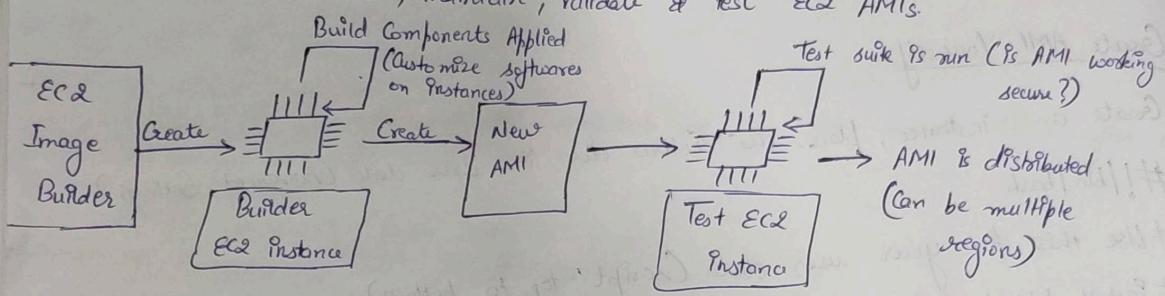
AMI > My AMIs (Owned by me)

Search the Public IPv4 & its loaded very soon. (not installing httpd again)

EC2 Image Builder

→ Automation framework

- ① Used to automate the creation of virtual machines or container images
- ⇒ Automate the creation, maintain, validate & test EC2 AMIs.



- ① Free Service (only pay for the underlying resources)
 - Like if we create an instance using this Image builder, we are only going to pay for our instances)

EC2 Instance Store

- ① "EBS Volumes" are network drives with "good" but limited performance
 - Sometimes we need even higher performance & that is going to be a hardware disk attached onto your EC2 instance
- ② If you need a high performance hardware disk use EC2 instance store
 - ⇒ EC2 instance is a virtual machine but it is obviously attached to a real hardware server
- ③ Better I/O performance
- ④ If you stop or terminate EC2 instance, that has an instance store, then the storage will be lost
 - ∴ It is called ephemeral storage
- ⑤ Good for buffer/cache/scratch data/temporary content.

- ① For long term
- ② Risk of data loss
- ③ Backups & Replication

If you see very few EC2 instances, consider

EFS (Elastic File System)

- ① Managed NFS
- ② EFS
- ③ FSS

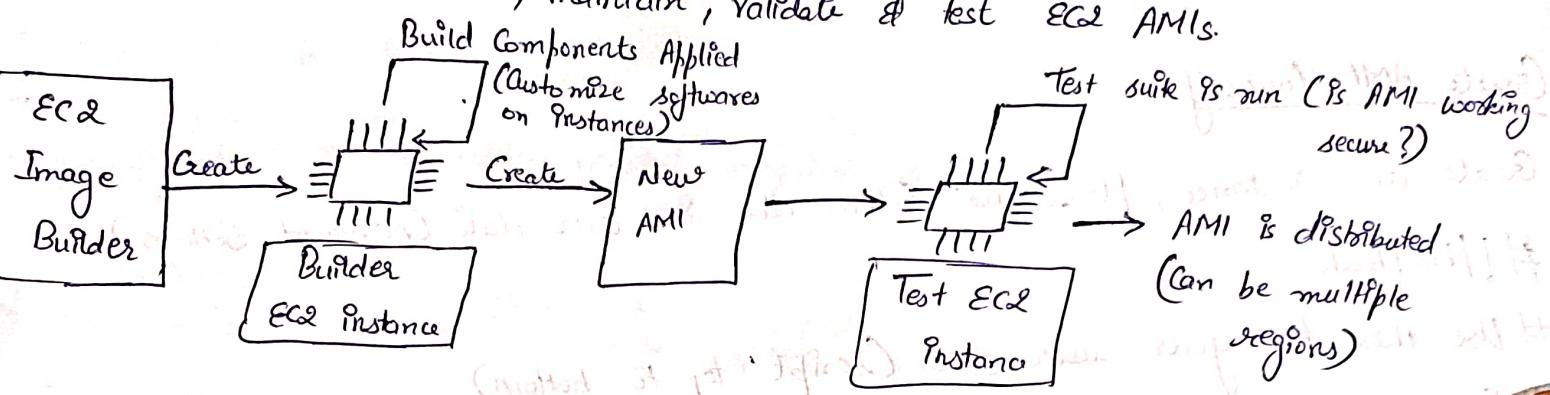
EBS vs EFS

With EBS,
EBS volume
But if we
we could
availability
But this

EC2 Image Builder

→ Automation framework

- ① Used to automate the creation of virtual machines or container images
- ⇒ Automate the creation, maintain, validate & test EC2 AMIs.



- ① Free service (only pay for the underlying resources)
 - Like if we create an instance using this image builder, we are only going to pay for our instances)

EC2 Instance Store

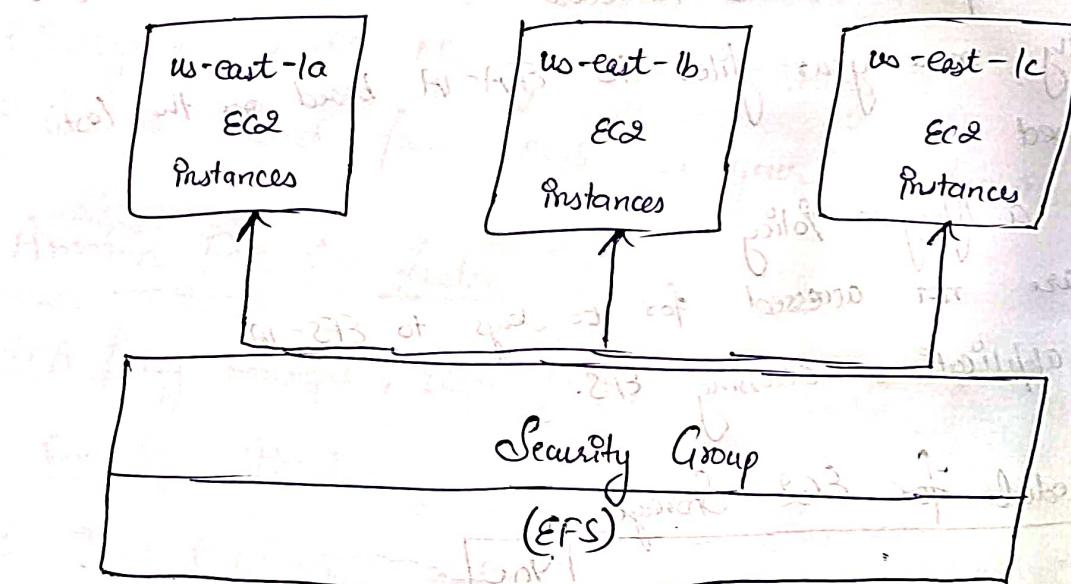
- ① "EBS Volumes" are network drives with "good" but limited performance
 - Sometimes we need even higher performance & that is going to be a hardware disk attached onto your EC2 instance
- ② If you need a high performance hardware disk use EC2 Instance Store
 - ⇒ EC2 instance is a virtual machine but it is obviously attached to a real hardware server
- ③ Better I/O performance
- ④ If you stop or terminate EC2 instance, that has an instance store, then the storage will be lost
 - ↳ It is called ephemeral storage
- ⑤ Good for buffer/cache/scratch data/temporary content.

- ① For long term memory storage, EBS is a better choice
- ② Risk of data loss if hardware fails
- ③ Backups & Replication are your responsibility

If you see very high performance hardware attached volumes for your EC2 instances, consider that it's local EC2 instance store.

EFS (Elastic File System)

- ④ Managed NFS (Network File System) that can be mounted on 100s of EC2 instances
- ⑤ EFS works with Linux EC2 instances in multi-AZ
- ⑥ Highly available, scalable, expensive (3x cost), pay per use, no capacity planning



EBS vs EFS

With EBS,

EBS volume can be attached to one specific instance in one specific AZ. But if we wanted to move over EBS volume from one AZ to other, we could create a snapshot & then restore that snapshot onto a new availability zone. But this is a copy, this is not an in-sync replica.

EFS is a network file system. This means that whatever is on the EFS drive is shared by everything that is mounted to it. Say, we have many instances in AZ 1, on one or many instances in AZ 2.

At the same time, all these instances can mount the same EFS drive, using a mount target & they will all see the same files. So, this makes it a shared file system.

EFS Infrequent Access (EFA-IA)

- ① Storage class that is cost optimized for files that are not accessed everyday
- ② Up to 98% lower cost compared to EFS standard
- ③ EFS will automatically move your files to EFA-IA based on the last time they were accessed
- ④ Enable EFS-IA with a lifecycle policy
- ⑤ Eg → move files that are not accessed for 60 days to EFS-IA
- ⑥ Transparent to the applications accessing EFS.

Shared Responsibility model for EC2 storage



- ① Infrastructure
- ② Replication for data for EBS volumes & EFS drives
- ③ Replacing faulty hardware
- ④ Ensuring (data privacy) i.e. employees can access your data
- ⑤ Setting up backup / snapshot procedures
- ⑥ Setting up data encryption
- ⑦ Responsibility of any data on the drives
- ⑧ Understanding the risk of using EC2 instance store

- ① Launch 3rd party high performance file systems on AWS
- ② Fully managed service
- ③ FSx for Lustre
- ④ FSx for Windows File Server
- ⑤ FSx for NetApp ONTAP

FSx for Windows File Server

- ① A fully managed, highly reliable & scalable windows native shared file system
- ② Built on Windows File Server
- ③ Supports SMB protocol & Windows NTFS.
- ④ Integrated with Microsoft Active Directory
- ⑤ Can be accessed from AWS or your on-premises infrastructure

Amazon FSx for Lustre

- ① A fully managed, high performance, scalable file storage for high performance computing (HPC).
- ② Linux + Cluster → Lustre
- ③ Machine learning, Analytics, Video processing, Financial Modeling
- ④ Scales upto 100 Gb/s, millions of IOPS, sub-ms latencies

EC2 Instance Storage - Summary

- ① EBS Volume
 - network drives attached to one EC2 instance at a time
 - mapped to an availability zone
 - Can use EBS snapshots for backups / transferring EBS volumes across AZ

- ① AMI : Create ready to use EC2 instances with customizations
- ② EC2 Image Builder : automatically build, test & distribute AMI's
- ③ EC2 Instance Store
 - High performance hardware disk attached to base of our EC2 instance
 - lost if our instance is stopped / terminated
- ④ EFS : network file system, can be attached to lots of instances in a region
- ⑤ EFS - IA : cost optimized storage class for infrequent, accessed files
- ⑥ FSx for windows : File system for windows servers
- ⑦ FSx for Lustre : High performance computing Linux file system

Scalability & High Availability

- ① Scalability means that an application/system can handle greater loads by adapting
- ② 2 kinds of scalability
 - Vertical Scalability
 - Horizontal
- ③ Scalability is linked but diff to High Availability

⇒ Vertical Scalability (Scale up / down)

- ④ It means \uparrow ing the size of instance
- ⑤ Eg, your application runs on tel. micro, scaling that application vertically means running it on a tel. large
- ⑥ It is very common for non distributed systems such as database
- ⑦ There is usually a limit to how much you can vertically scale

⇒ Horizontal Scalability (Scale out / in)

- ⑧ It means \uparrow ing the no. of instances / systems for your application
- ⑨ Horizontal scaling implies distributed systems
- ⑩ This is very common for web applications / modern applications
- ⑪ It's easy to horizontally scale.

High Availability

- ① It usually goes hand in hand with horizontal scaling
- ② It means running your application in at least 2 availability zones
- ③ The goal is to survive a data center loss (Disaster)

High Availability & Scalability

- ① Vertical Scaling → ↑ Instance size (= scale up/down)
 - From: t2.nano - 0.5GB of RAM, 1 vCPU
 - To: u-12tb.t1.micro - 12.3TB of RAM, 4vCPUs
- ② Horizontal Scaling → ↑ no. of instances (= scale out/in)
 - Auto Scaling group
 - Load Balancer
- ③ High Availability → Run instances for the same application across multiple AZs.

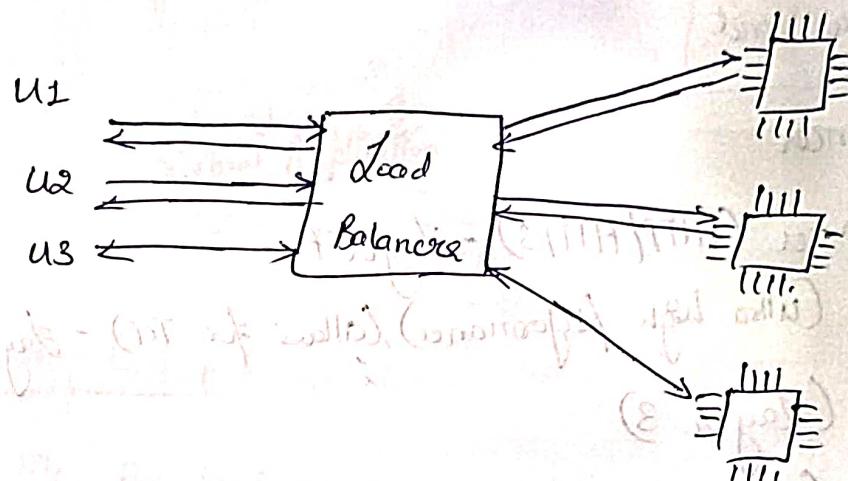
Scalability vs Elasticity

- ① Scalability → Ability to accommodate a larger load by making the hardware stronger (scale up) or by adding nodes (scale out)
- ② Elasticity? Once a system is scaled, elasticity means that there will be some "autoscaling" so that system can scale based on load. This is "cloud-friendly": pay per use, match demand, optimize cost

- ① **Agility** - New IT resources are a click away, which means that you reduce the time to make those resources available to your developers from weeks to just minutes.

What is Load Balancing?

- ② Load Balancers are servers that forward Internet traffic to multiple servers (EC2 Instances) downstream.



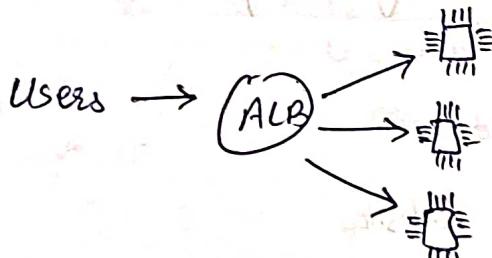
Why we use a Load Balancer?

- ③ Spread load across multiple downstream instances
- ④ Expose a single point of access (PNS) to your application
- ⑤ Seamlessly handle failures of downstream instances
- ⑥ Do regular health checks to your instances
- ⑦ Provide SSL termination (HTTPS) to your websites
- ⑧ High availability across zones.

Why Use an Elastic Load Balancer?

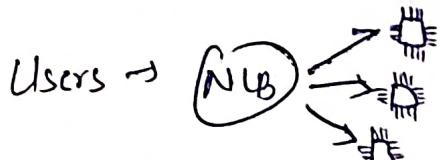
- ① An ELB is a managed load balancer
 - AWS guarantees that it will be working
 - AWS takes care of upgrades, maintenance, high availability
 - AWS provides only a few configuration knobs
- ② It costs less to setup your own load balancer but it will be a lot more effort on your end
- ③ 4 kinds of Load Balancers
 - Application load Balancer (HTTP/HTTPS) - Layer 7
 - Network " " " (Ultra high performance), (allows for TCP) - Layer 4
 - Gateway " " " (Layer 3)
 - Classic " " " (retired in 2023) - Layer 4 & 7

- Application Load Balancer →
- ① HTTP/ HTTPS/ gRPC protocols (Layer 7)
 - ② HTTP Routing features
 - ③ Static DNS (URL)



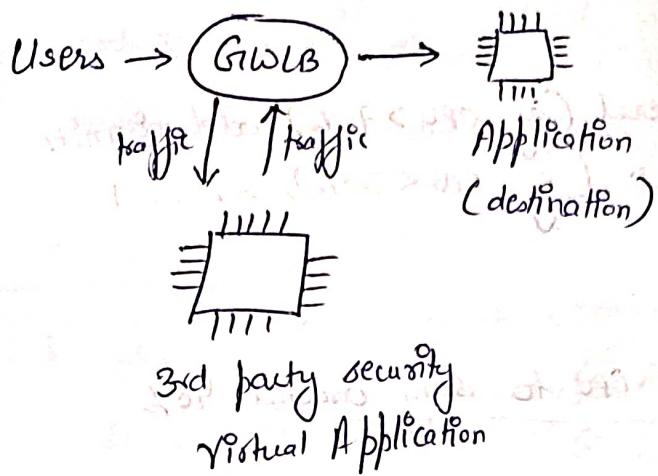
- Network Load Balancer →
- ① TCP/ UDP protocols (Layer 4)

- ② High performance : millions of requests per second
- ③ Static IP through Elastic IP.



Gateway Load Balancer → GENEVE protocol on IP packets (Layer 3)

- Route traffic to firewalls that you manage on EC2 instances
- Intrusion detection

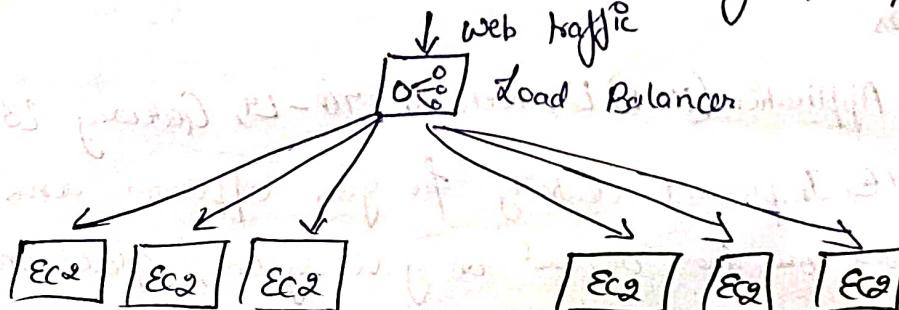


Health Checks

They enable load balancer to know if instances are available to accept requests or not (200)

Auto-Scaling Group?

- In real life, the load on your websites and application can change
- In the cloud, you can create and get rid of servers very quickly
 - Scale out (add EC2 instances) to match rising load
 - Scale in (remove EC2 instances) to match falling load
 - Ensure we have minimum & max. no. of machines running
 - Automatically register new instances to load balancer
 - Replace unhealthy instances
- Cost savings = only run at optimal capacity. (Principle of cloud)



Auto Scaling Groups - Scaling Strategies

- ① Manual Scaling → Update size of ASG manually
- ② Dynamic " → Respond to changing demand
 - Simple/Step scaling
 - ① When a CloudWatch metric is triggered (e.g. CPU > 70%), add 2 units
 - ② When another metric (e.g. CPU < 30%) is triggered, remove 1
 - Target Tracking Scaling
 - ① E.g. I want the avg. ASG CPU to stay around 40%
- Schedule Scaling
 - ① Anticipate a scaling based on known patterns
 - ② E.g. Increase the min capacity to 10 at 5pm on Fridays
- ③ Predictive Scaling
 - Use Machine Learning to predict future traffic ahead of time
 - Automatically provisions to the right no. of EC2 Instances in advance
 - Useful when your load has predictable time based patterns

ELB and ASG = Summary

- ④ High Availability vs Scalability (Vertical & horizontal) vs Elasticity vs Agility
In the Cloud
- ⑤ Elastic Load Balancer
 - Distributed traffic across backend EC2 Instances, can be Multi-AZ
 - Supports health checks
 - 4 types: Classic (old), Application (HTTP) L7, Network, TCP-L4, Gateway L3
- ⑥ ASG (Auto Scaling groups) → ⑦ Implement elasticity for your application across multiple AZ
 - ⑧ Scale EC2 Instances based on the demand on your system, replace unhealthy
 - ⑨ Integrated with ELB

Amazon S3

- ① Amazon S3 is one of the main building blocks of AWS
- ② It's advertised as "infinite ~~storage~~ scaling" storage
- ③ Many websites use Amazon S3 as backbone
- ④ Many AWS services use Amazon S3 as an integration as well

Amazon S3 Use Cases

- ① Backup & storage
- ② Disaster Recovery
- ③ Archive
- ④ Hybrid Cloud Storage
- ⑤ Application Hosting
- ⑥ Media Hosting
- ⑦ Data lakes & big data analytics
- ⑧ Software delivery
- ⑨ Static website

Amazon S3 - Buckets

- ① Amazon S3 allows people to store objects (files) in "buckets" (directories)
- ② Buckets must have a global unique name (across all regions, all accounts)
- ③ Buckets are defined at the region level
- ④ S3 looks like a global service but buckets are created in a region
- ⑤ Naming Convention
 - No upper case, no underscore
 - 3-63 characters long
 - Not an IP
 - Must start with lowercase letter or number
 - Must NOT start with the prefix `x-`
 - Must NOT end with the suffix `.s3.amazonaws.com`

Amazon S3 - Objects

- ① Objects (files) have a key
 - The key is the full path
 - `s3://my-bucket/[my-file.txt]`
 - `s3://my-bucket/[my-folder]/[another-folder]/[my-file.txt]`
- ② The key is composed of prefix + object name
 - `s3://my-bucket/[my-folder]/[another-folder]/[my-file.txt]`
- ③ There is no concept of "directories" within buckets
- ④ Just keys with very long names that contain slashes (/)
- ⑤ Object values are the content of the body
 - Max object size is 5TB (5000 GB)
 - If uploading more than 5GB, must use "multipart upload"
- ⑥ Metadata (list of text key / value pairs - system or user data)
- ⑦ Tags (unique key / value pair - up to 10) - useful for security / lifecycle
- ⑧ Version ID (if versioning is enabled)

Amazon S3 - Security

- ① User-Based
 - IAM-policies & which API calls should be allowed for a specific user from IAM
- ② Resource Based
 - Bucket policies & Bucket wide rules from S3 console - allows cross account
 - Object Access Control List (ACL) → finer grain (can be disabled)

→ Bucket Access Control List (ACL) - less common (can be disabled)

① NOTE: an IAM principal can access an S3 object if

- The user IAM permission Allow it or the resource policy Allows;
- AND there is no explicit deny

② Encryption: encrypt objects in S3 using encryption keys

③ JSON-Based policy

→ Resources: buckets and objects

→ Effect: Allow / Deny

→ Actions: set of API to Allow or deny

→ Principal: The account or user to apply the policy to

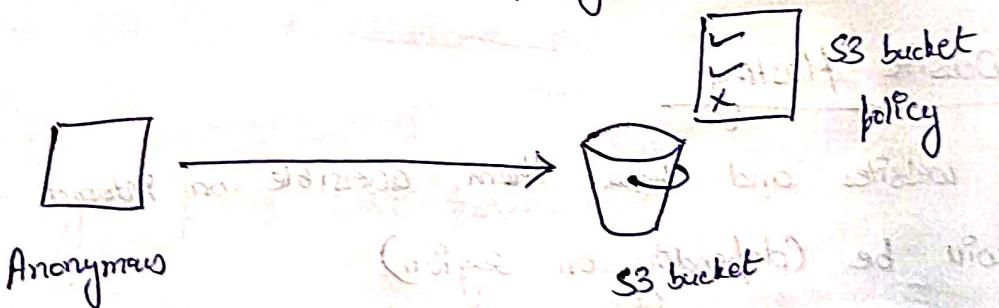
④ Use S3 bucket for policy to

→ Grant public access to the bucket

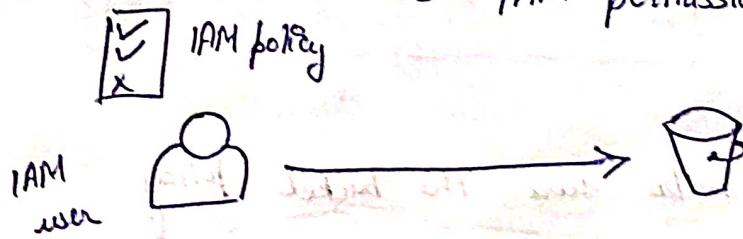
→ Force objects to be encrypted at upload

→ Grant access to another account (Cross account)

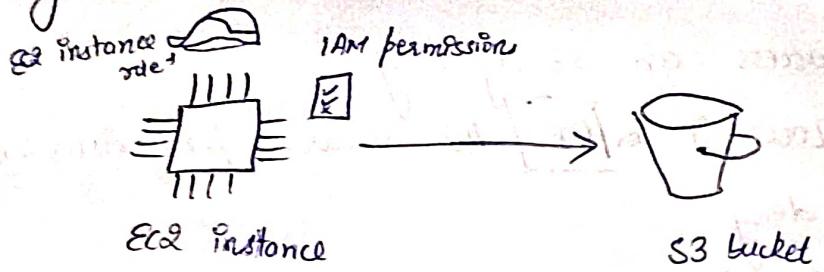
Eg Public Access - Use Bucket policy



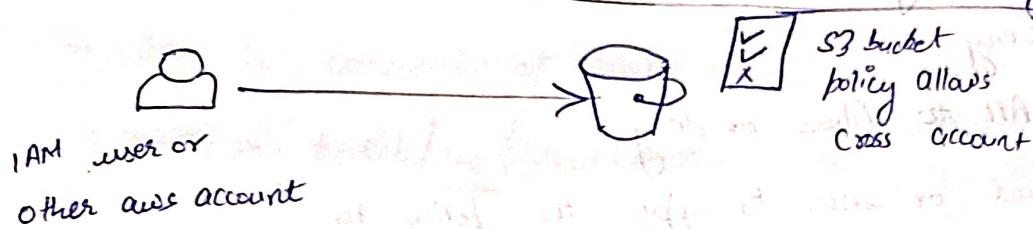
Eg → User access to S3 - IAM permissions



Eg → EC2 Instance Access - Use IAM Roles



Advanced : Cross - Account Access - Use Bucket Policy



Bucket settings for Block Public access

- ① These settings were created to prevent company data leaks
- ② If you know your bucket should never be public, leave these on
- ③ Can be set at account level.

Amazon S3 - Static Website Hosting

- ① S3 can host static websites and make them accessible on Internet
- ② The website URL will be (depending on region)
 - ③ `http://bucket-name.s3-website-us-west-2.amazonaws.com`
 - OR
 - ④ `http://bucket-name.us-west-2.amazonaws.com`
- ⑤ If you get a 403 forbidden error, make sure the bucket policy allows public reads

Amazon S3 - Versioning

22. 8. 2023 (Wednesday)

- ① You can version your files in Amazon S3.
 - ② It is enabled at bucket level.
 - ③ Same key overwrite will change the "version": 1, 2, 3, ...
 - ④ It is best practice to version your buckets.
 - Protect against unintended delete (ability to restore at version 2 or 3)
 - Easy roll back to previous version
 - ⑤ Notes:
 - Any file that is not versioned prior to enabling versioning will have version "null".
 - Suspending versioning does not delete previous versions.
- Amazon S3 - Replication (CRR & SRR)
- Cross Region Some Region replication

- ① Must enable versioning in source & destination.
- ② Buckets can be in different AWS accounts.
- ③ Copying is asynchronous.
- ④ Must give proper IAM permissions to S3.
- ⑤ Use cases:

- CRR → Compliance, lower latency access, replication across accounts.
- SRR → log aggregation, live replication between production & test accounts.

Creating Bucket in S3

- ① Create Bucket
 - ② Upload an image
 - ③ When we click open, our image is displayed on the Internet
 - ④ But when we use 'Object URL', it does not work (Access denied)
- (S3-pre-signed URL)

↳ A signature that verifies I am the one making the request & has my credentials

Deleting a version → Permanently delete (can't be undone)

" an object → delete

(Delete marker) (Checkpoint to rollback to in future)

versioning

S3 Storage Classes

- ① Amazon S3 Standard - General Purpose
- ② " " " - Infrequent Access (IA)
- ③ " " " one zone - Infrequent Access
- ④ " " " instance retrieval
- ⑤ " " " glacier flexible retrieval
- ⑥ " " " glacier Deep Archive
- ⑦ " " " Intelligent Tearing.
- ⑧ Can move b/w classes manually or using S3 lifecycle configurations.

S3 Durability & Availability

(S3) 200A Singapore - India 32

Durability → High durability (99.99999999%) of object across all multiple AZ

① If you store 10 million objects with Amazon S3, you can on average expect to incur a loss of a single object once every 19,000 years.

② Durability represents the concept of how many objects will be lost by S3

③ Same for all storage classes.

④ Measures how readily a service is

⑤ Varies depending on storage class.

⑥ Eg - S3 Standard has 99.99% availability = not available for 53 minutes a year.

⑦ 99.99% availability

⑧ Used for frequently accessed data

⑨ Low latency & high throughput

⑩ Sustain & concurrent facility failures

Big data, Analytics, mobile & gaming applications, content distribution

S3 Standard - Infrequent Access (IA)

infrequent & infrequent

- ① for data that is less frequently accessed, but requires rapid access when needed.
- ② lower cost than S3 standard
- ③ 99.99% availability
- ④ Use Cases → Disaster recovery, backups

S3 One-Zone - Infrequent Access

- ① High durability (99.99999%) in a single AZ, data lost when AZ is destroyed
- ② 99.5% availability
- ③ Use Cases → Storing secondary backup copies of on-premises data or data you can re-create

S3 Glacier Storage Classes

- ① low-cost object storage meant for archiving / backup
- ② Pricing = price for storage + object retrieval cost
- ③ Amazon S3 Glacier Instant Retrieval
 - Millisecond retrieval, great for data accessed once a quarter
 - Minimum storage duration of 90 days.
(I want data instantly)
- ④ Amazon S3 Glacier Flexible Retrieval (Formerly S3 glacier)
 - Expedited (1 to 5 mins), Standard (3 to 5 hrs), Bulk (5 to 12 hrs) — free
 - Minimum storage duration is 90 days

- ① Amazon S3 Glacier Deep Archive - for long term storage
(I can wait to receive data for 12hrs)
 - Standard (12hrs), Bulk (48 hrs)
 - Minimum storage duration of 180 days
 - Least expensive

S3 Intelligent-Tiering

- ① Small monthly monitoring & auto-tiering fee
- ② Moves objects automatically between Access Tiers based on usage
- ③ There are no retrieval charges in S3 Intelligent-Tiering
- ④ Frequent Access Tier : (Automatic) : default tier
- ⑤ Infrequent Access Tier (Automatic) : objects not accessed for 30 days
- ⑥ Archive Instant Access Tier (Automatic) : objects not accessed for 90 days
- ⑦ Archive Access Tier (Optional) : configurable from 90 days to 1000 days
- ⑧ Deep Archive, to Access Tier (Optional) : Configurable from 180 days to 1000 days

Management > Lifecycle Rules > Set automated storage classes (ace to certain

S3 Encryption

- ① Server Side Encryption (default)

By default, applied, whenever we create a bucket, it will be encrypted

A user uploads an object in S3, that object when it arrives in the bucket is encrypted by S3 for securing purposes

Server encrypts the file after receiving it.

- ① Client Side Encryption
- When a user takes a file & encrypts it before uploading to the server
- S3 uses AWS KMS (Key Management Service) to store encrypted keys
- S3 uses AWS CloudWatch Metrics to monitor access patterns
- IAM Access Analyzer for S3

- ① Ensures that only intended people have access to your S3 buckets
- ② Eg → Publicly accessible buckets, bucket shared with other AWS account ...
- ③ Evaluates S3 Bucket policies, S3 ACLs, S3 Access Point Policies
- ④ Powered by IAM Access Analyzer

Shared Responsibility Model for S3



- ① Infrastructure (global security, durability, availability, sustainability, sustain concurrent loss of data in two facilities)
- ② S3 versioning
- ③ S3 bucket policies
- ④ S3 replication setup
- ⑤ Logging & monitoring
- ⑥ S3 storage classes
- ⑦ Data encryption at rest & in transit
- ⑧ Configuration & vulnerability analysis
- ⑨ Compliance Validation

AWS Snow Family

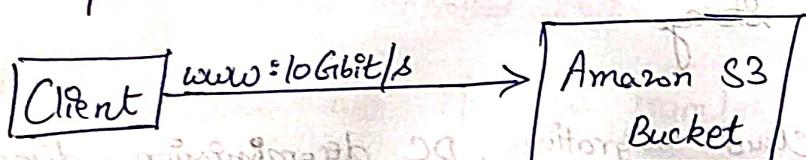
- ① Highly secure, portable devices to collect & process data at the edge, and migrate data into & out of AWS
- ② Data migration: Snow cone, Snowball edge, Snow mobile
- ③ Edge computing: Snow cone, Snowball (Edge not yet supported)

Data Migration with AWS Snow Family

- Challenges &
- ① Limited Connectivity
 - ② Limited Bandwidth
 - ③ High network cost
 - ④ Shared Bandwidth (can't maximise the line)
 - ⑤ Connection Stability

AWS Snow family → offline devices to perform data migrations. If it takes more than a week to transfer over the network, use snowball devices

- ① Direct upload to S3



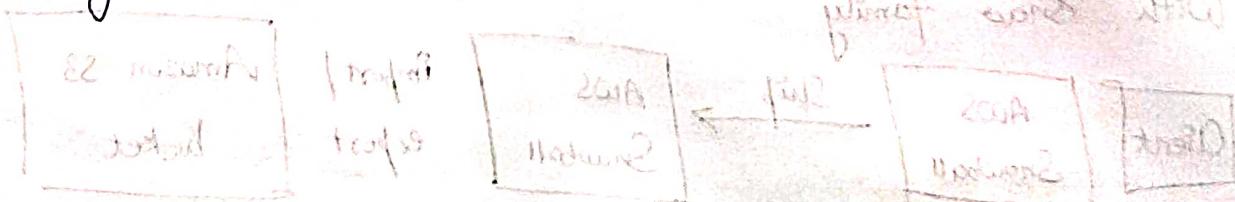
- ② With Snow family



Client request Snowball which is delivered by AWS. We load the data directly onto the devices locally, and then we ship back the device to AWS and then they will plug that device into their infrastructure to import/export file from/to Amazon S3 bucket.

SnowBall Edge (for data transfer)

- ① Physical data transfer solution: move TBs or PBs of data in or out of AWS
- ② Alternating to moving data over the network (& paying network fees)
- ③ Pay per data transfer job
- ④ Provide block storage and Amazon-S3-compatible object storage
- ⑤ Snowball Edge Storage Optimized
 - 80 TB of HDD capacity for block volume & S3-compatible object storage
- ⑥ Snowball Edge Compute Optimized
 - 42 TB of HDD or 28 TB NVMe capacity for block volume & S3 compatible object storage
- ⑦ Use Cases → Large scale cloud migrations, DC decommission, disaster management



AWS Snowcone & Snowcone SSD

- ① Small, portable computing, anywhere, rugged & secure, withstands harsh environments
- ① Light (4.5 pounds, 2.1 kgs)
- ① Device used for edge computing, storage & data transfer
- ① Snowcone - STB of HDD Storage
- ① Snowcone SSD - 14 TB of SSD storage
- ① Use snowcone where snowball does not fit (space-constrained environment)
- ① Must provide your own battery/cables.
- ① Can be sent back to AWS offline, or connect it to Internet & use AWS DataSync to send data

AWS SnowMobile

- ① Transfer exabytes of data ($1EB = 1,000 PB = 1,000,000 TBs$)
- ① Each snowmobile has 100 PB of capacity (use multiple in parallel)
- ① High security: Temp. controlled, GPS, 24/7 video surveillance
- ① Better than snowball if you transfer more than 10PB.

	Snowcone & Snowcone SSD	Snowball Edge Storage Optimized	Snow mobile
Storage Capacity	8TB of HDD 14TB of SSD	80 TB usable	>100 PB
Migration size	Up to 24 TB, online & offline	Up to Petabytes, offline	Up to exabytes, offline
Data Sync agent	Pre-Installed		

Snowball family : Usage Process

- ① Request Snowball device from the AWS console for delivery
- ② Install the Snowball Client / AWS Command Line Interface (CLI) on your servers
- ③ Connect the Snowball to your servers & copy files using the client
- ④ Ship back the device when you are done (goes to the right AWS facility)
- ⑤ Data will be loaded into an S3 Bucket
- ⑥ Snowball is completely wiped

What is Edge Computing?

- ① Process data while it's being created on an edge location.
Eg → A truck on the road, a ship in the sea etc
- ② These locations may have
 - limited / no Internet access
 - limited / no easy access for computing power
- ③ We setup a Snowball Edge / Snowcone device to do edge computing
- ④ Use cases of edge computing
 - Pre process data
 - Machine learning at the edge
 - Transcoding media streams
- ⑤ Eventually, (if need be) we can ship back the device to AWS (for transferring data for example)

Snow Family - Edge Computing

period 3 notes

- ① Snowcone & Snowcone SSD (Smaller)
 - 2 CPUs, 4Gb of memory, wired or wireless access
 - USB-C power using a standard cord or the optional battery
- ② Snowball Edge - Compute Optimized
 - 104 vCPUs, 416 GiB of RAM
 - Optional GPU (useful for video processing or machine learning)
 - 28 TB NVMe or 42 TB HDD usable storage
 - Storage clustering available (upto 16 nodes)
- ③ Snowball Edge - Storage Optimized
 - Up to 40 vCPUs, 80 GiB of RAM, 80 TB storage
 - All can run EC2 instances & AWS Lambda functions (using AWS IoT Greengrass)
 - Long Term deployment options - 1 & 3 years discounted pricing

AWS OpstHub

- ① Historically, to use snow family devices, you need CU
- ② Today, you can use AWS OpstHub (a software you install on your device) to manage your snowfamily device
 - Unlocking & configuring single or clustered devices
 - Transferring files
 - Launching & managing instances running on snowfamily devices
 - Monitor devices metrics (storage capacity, active instances on your device)
 - Launch compatible AWS services on your devices (Eg - Amazon EC2 Instances, AWS Datasync, Network File System (NFS))

Snowball Edge Pricing

- ① You pay for device usage & data transfer out of AWS network
- ② Data transfer IN to Amazon S3 is \$0.00 per GB
- ③ On-Demand → Includes a one time service fee per job, which
 - Includes
 - 10 days of usage for snowball edge storage optimized 80TB
 - 15 days of usage for snowball edge storage optimized 110TB
 - Shipping days are not counted towards the included 10 or 15 days
- ④ Committed Upfront → Pay in advance for monthly, 1-year & 3 years of usages (edge-computing)
 - Up to 60% discount pricing

Hybrid Cloud for Storage

- ① AWS is pushing for "hybrid cloud"
- ②
 - Part of your infrastructure is on-premises
 - "Data is where you want it on the cloud"
- ③ This can be due to
 - Long Cloud migrations
 - Security Requirements
 - Compliance
 - IT Strategy
- ④ S3 is a proprietary storage technology (unlike EFS | NFS). So, how do you expose S3 data on-premises?
 - (AWS Storage Gateway)

Aws Storage Cloud Native Options

points - 82

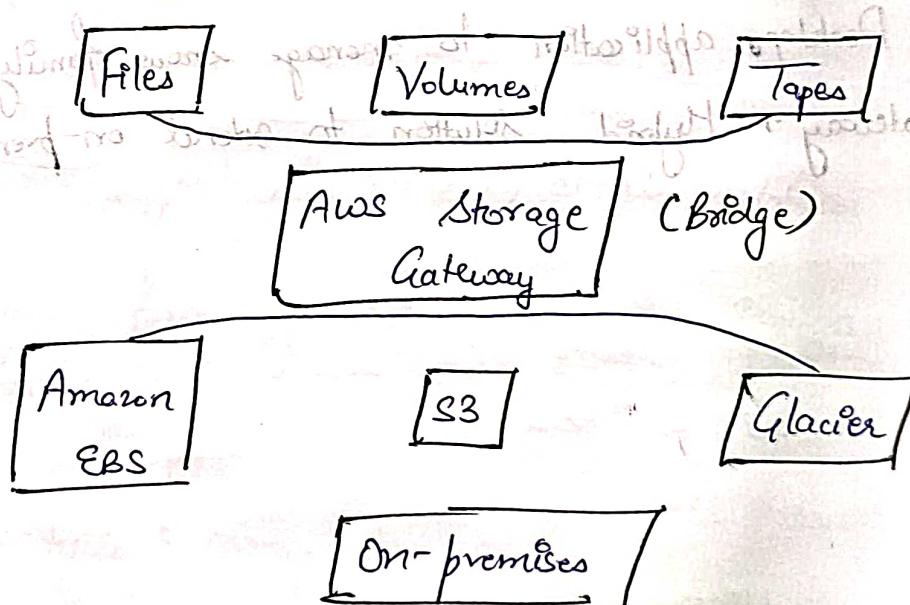
- Block & Amazon EBS, EC2 instance store
- File & Amazon EFS

- Object & Amazon S3, Glacier

Aws Storage Gateway

- bridge b/w on-premises data & cloud data is S3
- Hybrid storage service to allow on-premises to seamlessly use AWS cloud
- Use Cases: Disaster recovery, backup & restore, tiered storage
- Types of Storage Gateways

- File Gateway
- Volume Gateway
- Tape Gateway



S3 - Summary

- ① Buckets vs Objects → Buckets have a globally unique name, tied to a region
Object lives within buckets
- ② S3 Security → IAM policy, S3 Bucket policy (Public access), S3 Encryption
- ③ S3 Websites → Host a static website on Amazon S3 (Ensuring our bucket is public)
- ④ S3 Versioning → Multiple versions for files, prevent accidental deletes
- ⑤ S3 Replication → Same region or cross region, must enable versioning
- ⑥ S3 Storage classes → Standard, IA, 1Z-IA, Intelligent, Glacier (Instant, Flexible, Deep)
- ⑦ Snow Family → Import data on S3 through a physical device, edge computing
- ⑧ Ops-Hub → Desktop application to manage snow family devices
- ⑨ Storage Gateway → Hybrid solution to extend on-premises storage to S3

Databases

① Storing data on disk (EFS, EBS, Instance Store, S3) can have limits

Relational Database

- ① Excel spreadsheets, with links between them
- ② Can use SQL language to perform queries

NoSQL Databases

- ① Non-relational DB's
- ② These are purpose built for specific data models and have flexible schemas for building

Benefits

- Flexibility
- Scalability
- High performance
- Highly functional

Eg → key-value, document, graph, in-memory, search databases

- ① JSON is a common form of data that fits into a NoSQL model

① Data can be nested

① Fields can change over time

① Support for new types: arrays etc

Shared Responsibility Databases on AWS

AWS

User

- Quick provisioning, high availability, v/h horizontal scalability
- Automated backup & restore, operations, upgrades
- OS patching is handled by AWS
- Monitoring, alerting
- Resiliency, backup, patching, high availability, fault tolerance, scaling

Amazon RDS Overview

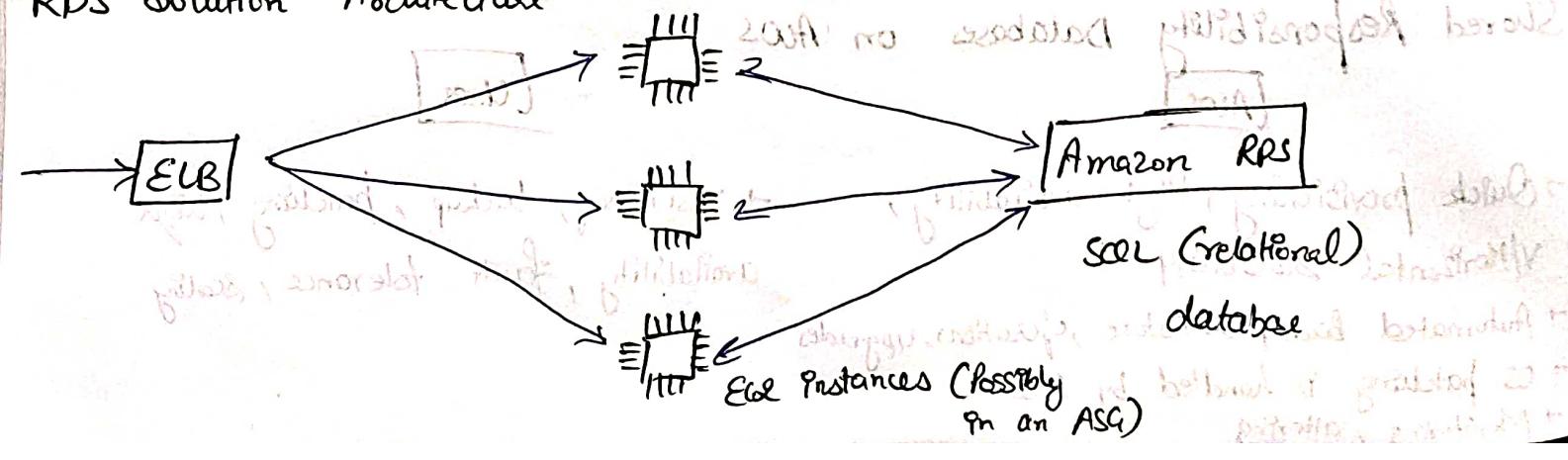
- ① Relational Database Service (RDS) With no db instance
- ② DB that use SQL as a query language read/write partitioned
- ③ Allows us to create databases in cloud that are managed by AWS Database as a service
 - PostgreSQL
 - MySQL
 - MariaDB
 - Oracle
 - Microsoft SQL Server

Why RDS Over Deploying DB on EC2?

- ④ Managed Service
 - Automated provisioning, OS patching
 - Continuous backups & restore to specific timestamp
 - Monitoring dashboards
 - Read replicas for improved read performance
 - Multi AZ setup for disaster recovery (DR)
 - Maintenance windows for upgrades
 - Vertical / Horizontal scaling
 - Storage backed by EBS

- ⑤ But you can't SSH into your instance

RDS Solution Architecture



Amazon Aurora

- ① Aurora is a proprietary technology from AWS (not open-sourced)
- ② PostgreSQL and MySQL are both supported as Aurora DB
- ③ Aurora is "AWS Cloud Optimized" & claims 5x performance improvement over MySQL on RDS, over 3x performance of PostgreSQL on RDS
- ④ Aurora storage automatically grows in increments of 10GB up to 128TB
- ⑤ Costs more than RDS (20% more) — more efficient
- ⑥ Not in free-tier

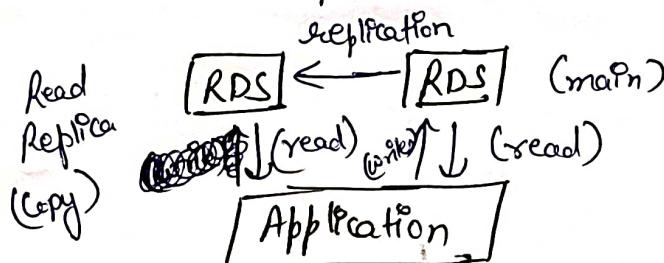
Amazon Aurora Serverless

- ① Automated DB instantiation and auto scaling based on actual usage
- ② PostgreSQL & MySQL both supported
- ③ No capacity planning needed
- ④ Least management overhead
- ⑤ Pay per second, can be more cost-effective
- ⑥ Use Cases → good for infrequent, intermittent or unpredictable workloads

RDS Deployments: Read Replicas, Multi-AZ

① Read Replicas

- ① Scale the read workload of your DB
- ② Can create up to 15 read replicas

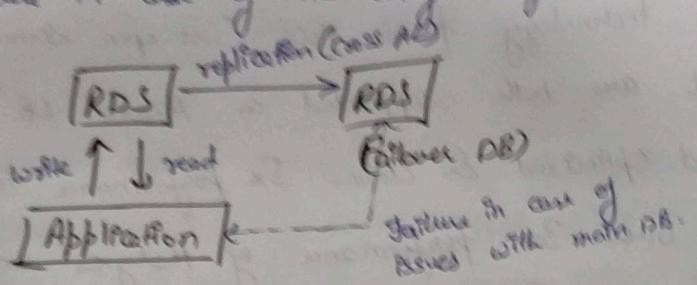


(Applications can read from any RDS)

(Application only has to write to main RDS)

⑤ Multi-AZ

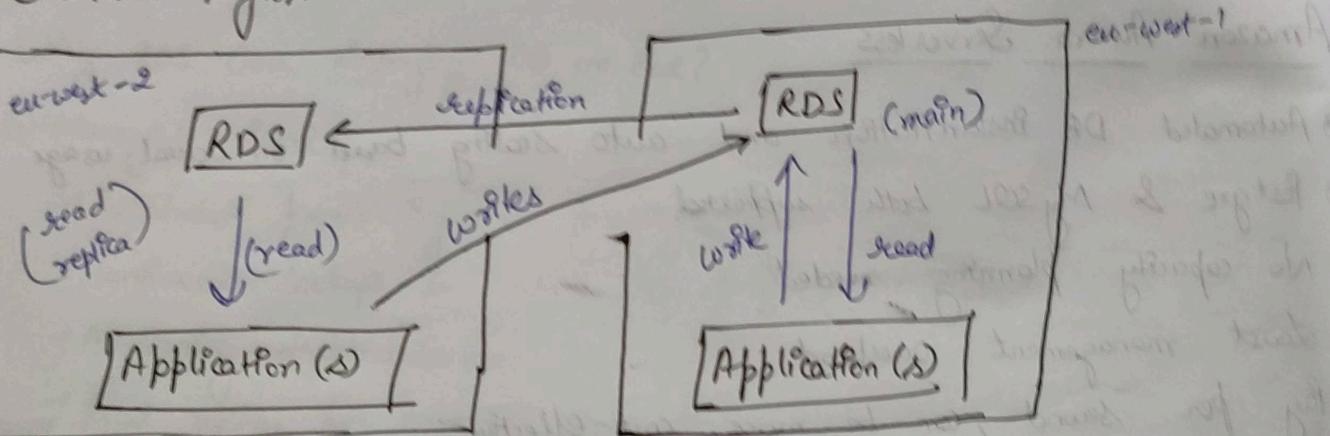
⑥ Failure in one of AZ outage (high availability)



⑦ Data is only read/written to the main DB

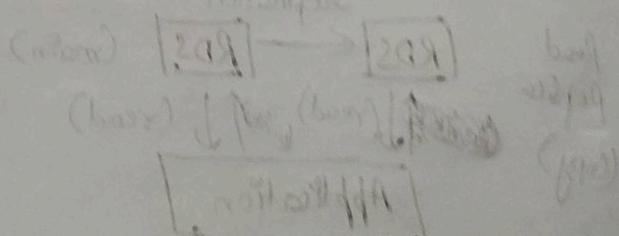
⑧ Can only have 1AZ as failover

⑨ Multi-Region



Multi-Region (Read Replica)

- Disaster recovery in case of region issues
- Local performance for global reads
- Replication cost



(Data must be replicated)

(Data must be replicated)

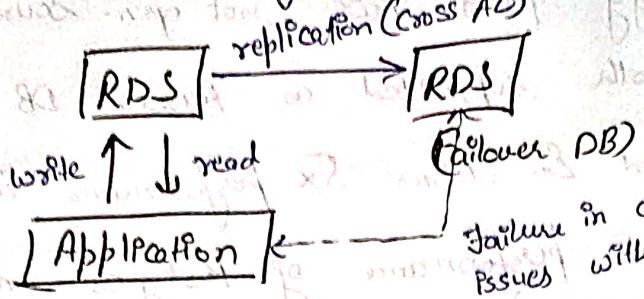
Amazon
 ④ Elasti
 ⑤ Caches
 ⑥ Helps
 ⑦ Acks take
 configuration

Solution

ELB

⑥ Multi-AZ

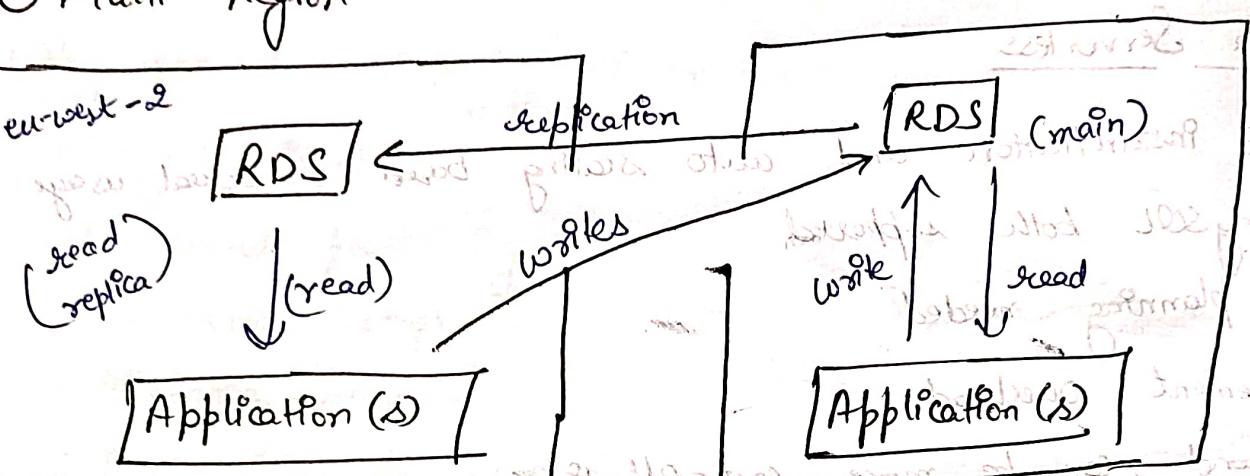
- Failure in case of AZ outage (High availability)



- Data is only read/written to the main DB

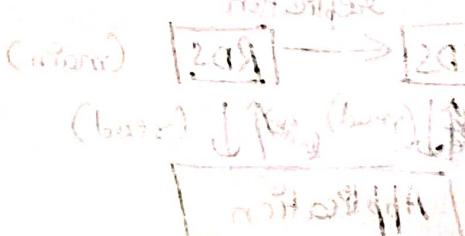
- Can only have 1 AZ as failover

⑦ Multi-Region



Multi-Region (Read Replica)

- Disaster recovery in case of region issues
- Local performance for global reads
- Replication cost



(Read from master) → (Read from replica)

(Copy from master to each of the replicas)

Amazon ElastiCache Overview

- ① ElastiCache is to get managed Redis or Memcached
- ② Caches are in-memory databases with high performance, low latency
- ③ Helps reduce load off databases for read-intensive workloads
- ④ AWS takes care of OS maintenance/patching, optimizations, setup, configuration, monitoring, failure recovery & backups

Solution Architecture - Cache

