

→ AWS Regions :-

→ AWS Regions all around the world

→ How to choose AWS region :-

↳ Compliance with data governance & legal req

↳ proximity to customers

↳ Availability zones <sup>services</sup> within a region

↳ pricing

→ Each region has many availability zones  
( min - 3 / max - 6 )

→ 400 + point of presence

→ 90+ Regional Caches in 90+ countries

→ IAM : Users & Groups

→ IAM :- Identity and Access Management, Global services

→ Root account :- created by default, shouldn't be used or shared.

→ Users are people within your organization, and can be grouped.

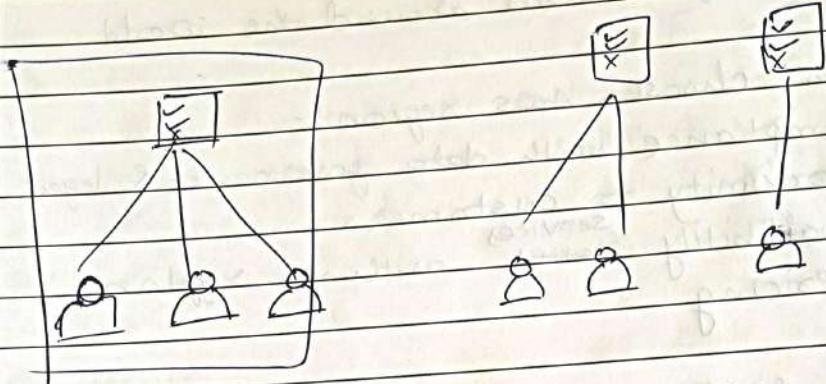
→ Groups only contains users, not other groups.

→ User don't have to belong to a group and user can belong to multiple groups.

→ Users or Groups can be assigned JSON documents called Policies.

→ Policies define permissions.

→ IAM policies Inheritance



→ IAM policy structure

→ consist of

↳ version

↳ Id

↳ statement

↳ Sid

↳ Effect

↳ principal

↳ Action

↳ Condition

→ To access AWS there is three option

↳ AWS Management console

↳ AWS Command Line Interface

↳ AWS Software Development Kit

- AWS CLI
  - Interact with AWS services using commands in your CLISHell.
  - Direct access to public APIs of AWS services
- ⇒ AWS SDK
  - AWS software development kit
  - Enable access and manage AWS services programmatically
  - Supports:- JS, python, .Net, Ruby and etc.
  - Mobile SDKs (Android, iOS)
- IAM Roles for services
  - Some AWS services will need to perform actions on your behalf.
  - So we assign permission to AWS services with IAM Roles
  - Common roles:- EC2 Instance Roles, Lambda Function Roles and roles of CloudFormation.
- IAM Security Tool
  - ↳ IAM Credentials Report (Account level)
  - ↳ IAM Access Advisor (User-level)

## → EC2 sizing & configuration option

- ↳ OS (Mac, Linux & Windows)
- ↳ Computer power & core (CPU)
- ↳ Random-access Memory (RAM)
- ↳ Network-Attached (EBS & EFS)
- ↳ Hardware attached (EC2 Instance Store)
- Network card
- Firewall rules
- Bootscript strip

## → Infer Instance type ! example

- ↳ t2.micro, t2.xlarge, c5d.4xlarge.

→ t2.micro is free on AWS free tier (750 per month)

## → m5.2xlarge

- ↳ m → instance class
- ↳ 5 → generation
- ↳ 2xlarge → size within an instance class.

## → Security Groups

- Security groups are fundamental of network security in AWS.
- Control in and out traffic of our EC2 instances.
- If contain allow rules.
- rules can be referred by IP or security groups.
- can attached to Multiple instances.

→ classic ports

business no 877 < \*

→ 22 → SSH

SSH and telnet not good

→ 21 → FTP ~~too many on the net~~ finding out ports

↳ port scanner) most-pain off

→ 80 → HTTP

→ 443 → HTTPS most-ports not blocked or 99%

→ 3389 → RDP ~~Windows better restricted~~

→ EC2 Instances purchasing option

① On demand Instance:- short workload, predictable pricing, pay by second.

② Reserved (1 & 3) years ~~price negotiations~~  
↳ Reserved Instance → long workloads

↳ Convertible Reserved Instance → long workload with  
flexible instance

③ Saving plan:- (1 & 3 years) commitment to an amount of usage, long workloads

④ Spot Instances - short workloads, cheap, can lose instance

⑤ Dedicated Hosts! - book an entire physical server, control instance placement.

⑥ Dedicated Instances! - no other customers will share your hardware.

⑦ Capacity Reservation - reserve capacity in a specific AZ for any duration.

## \* → EC2 on demand

- pay for what you use
- Has the highest cost but no upfront payment
- No long-term commitment

→ Recommended for short-term and uninterrupted workloads.

## \* → EC2 Reserved Instances

- 72% discount to on-demand pricing
- you reserve specific instance attributes
- Reservation period → 1 year (+ discount 0%) or 3 years (+++ discount)
- payment option (No-upfront (+)) / partial Upfront
- Recommended: All up-front (+) and All up-front (++) for steady state usage
- buy and sell in the Reserved Instance marketplace.

## \* → Convertible Reserve Instances

- can change EC2 instance type, instance family, os, etc.
- up to 66% discount

## \* → EC2 Saving Plans

- Get a discount on long term usage
- Commit to some certain type of usage
- Usage beyond EC2 Saving plan is billed at the on-demand price.

## \* → EC2 Spot Instances

- can get <sup>upto</sup> 90% compared to on-demand server
- Instance can lose at any time
- Most cost-effective instance in AWS.
- Not suitable for critical jobs or databases.

## \* → EC2 Dedicated Hosts

- A physical server with EC2 instance capacity fully dedicated to us
- Allow you to address compliance requirement and use your existing server-bound software licenses
- Purchase option → on demand, reserved
- most expensive
- Useful for SW that have complicated licensing model (BYOL)

## \* EBS Volume :-

- EBS (Elastic Block Store) volume is a network drive you can attach to your instance while they are running.
- Allow to persist data, even after termination.
- Only be mounted to one instance at a time.
- Bound to specific availability zone.
- Network USB stick.
- Use networking to communicate.
- Quickly attach to another.
- Locked to availability zone.

## \* EBS Snapshots

- Make a backup of your EBS volume at a point in time.
- Not necessary to detach volume to do snapshot but recommended.
- Can copy snapshots across AZ or Region.
- EBS Snapshot archive.
- EBS Recycle bin for EBS.
- Fast Snapshot restore.

## \* AMI (Amazon machine Image)

- AMI are a customization of an EC2 instance.
- You can add own SW, config, OS, monitoring.
- AMI is build for specific region (Copied across other regions).

- You can launch EC2 instance from:
  - ↳ A public AMI → AWS provided
  - ↳ Your own AMI → own created
  - ↳ An AWS Marketplace AMI → Someone else made.

### \* AMI process

== Start an EC2 instance and customize it

↓ Stop the instance (data integrity)

Build an AMI (create EBS snapshots)

↓ Launch instance from other AMIs

(combination of EBS based xens)  $\rightarrow$   $\downarrow$   $\rightarrow$   $\rightarrow$

### \* EC2 Instance Store

→ EBS volume are network drive with good limited performance.

→ If you need a high-performance hardware disk

$\rightarrow$   $\downarrow$   $\rightarrow$   $\rightarrow$   $\rightarrow$   $\rightarrow$

↳ Better I/O performance

↳ Good for buffer | cache | scratch | temporary content

↳ Risk of data loss if hardware failure

↳ Backup and replication are your responsibility.

$\rightarrow$   $\rightarrow$   $\rightarrow$   $\rightarrow$   $\rightarrow$   $\rightarrow$

Network (NIA)  $\rightarrow$   $\rightarrow$   $\rightarrow$   $\rightarrow$   $\rightarrow$

Information (NIA)  $\rightarrow$   $\rightarrow$   $\rightarrow$   $\rightarrow$   $\rightarrow$

## \* EBS Volume Types

EBS Volumes comes in 6 types.

- ↳ gp2 / gp3 (SSD)
- ↳ io1 / io2 Block Express (SSD)
- ↳ st1 (HDD)
- ↳ sc1 (HDD)

## \* Amazon EFS - Elastic file System

- Use cases :- Content Management, web serving, data sharing, wordpress.
- Use NFSv4.1 protocol.
- Uses security group to control access to NFS.
- Compatible with Linux based AMI (not windows).
- Encryption at rest using KMS.
- Scale automatically, pay-per-use, no capability planning.

## \* Load balancers

- Load balancers are servers that forward traffic to multiple servers downstream.
- Elastic Load Balancer
- It integrated with many services
- ↳ EC2, EC2 auto scaling, Amazon ECS
- ↳ AWS Certificate Manager (ACM), CloudWatch
- ↳ Route 53, AWS WAF, AWS Global Accelerator

→ Health check are crucial for Load balancer.

→ Types of Load balancer:

↳ Classical Load balancer

↳ Application Load balancer

↳ Network Load Balancer

↳ Gateway Load balancer.

\* Application Load Balancer (ALB)

→ Application Load Balancer is layer 7 (HTTP).

working through two components

→ Load balancing to multiple HTTP application across machines.

→ Load balancing to multiple applications on the same machine.

→ Support redirects

→ best for micro-services & container-based App

→ smart intelligent balancer present in AWS

load balanced on built

\* Network Load Balancer

↳ works at Layer 4 (TCP and UDP)

↳ Forward TCP & UDP traffic to your instance

↳ Handle millions of Request per second.

↳ NLP has 1 static IP / AZ available

→ NLP used for extreme performance TCP or UDP traffic.

Load balancer

Session based routing

at one distributor can have several load balancer

SA in each region

## \* Gateway Load Balancer

- Deploy, scale and manage a fleet of 3rd party network virtual application appliances in AWS.
- Ex:- firewall, IDS, IPS, Deep packet Inspection systems, Play payload manipulation.
- Operate at layer 3.
- combine the following function
  - ↳ transparent network gateway
  - ↳ Load Balancer
- uses the GENEVE protocol on port 6081.

## \* Sticky Sessions

- It is possible to implement stickiness so that the same Client is always redirected to the same instance behind a load balancer.
- works for all the Load balancer.
- make sure that the user doesn't lose his session data.
- Types
  - ↳ Application-based cookies
  - ↳ Duration-based cookies

## \* Load Balancer

cross zone load balancer

Each load balancer instance distributed evenly to all the registered instances in AZ.

## → SSL/TLS

→ An SSL certificate allows traffic between your clients and your load balancer to be encrypted in transit.

SSL → Secure Sockets Layer

TLS → Transport Layer Security

## → SNI :-

→ SNI helps to solve problem to bind multiple SSL certificates onto one web server.

→ Only works with ALB & NLB, Cloudfront

## \* Auto Scaling Group

→ manages load of servers in a base and

→ Scale out (add EC2 Instance) to match increase load

→ Scale in (remove EC2 Instances)

→ Ensure min & Max no of EC2 Instances.

→ Automatically registers new instance to load balancer.

→ Recreate automatically.

→ Force drain sign at beginning of config

→ private address with stop sign from port 80 to

→ configure base specia

## \* Amazon RDS Overview

- RDS → Relational database service
- It's managed DB services for DB use SQL as a query language.
- allows to create db in the cloud allowed by AWS
  - ↳ PostgreSQL, MySQL, MariaDB, Oracle, Microsoft SQL Server, Aurora(AWS)
- RDS also provide other services
  - ↳ Automated provisioning, OS patching, continuous backup and restore at specific time stamp.
  - Monitoring dashboard, read replicas, Multi AZ.
  - Maintenance windows, scaling capability,
  - backed by EBS.

## \* RDS Read Replicas for read scalability

- upto 15 read replicas
- Within AZ, cross AZ, or cross Region.
- Replication is ASYNC, so reads are eventually consistent
- Replicas is promoted to their own DB.
- App must update the connection string to leverage read replicas.

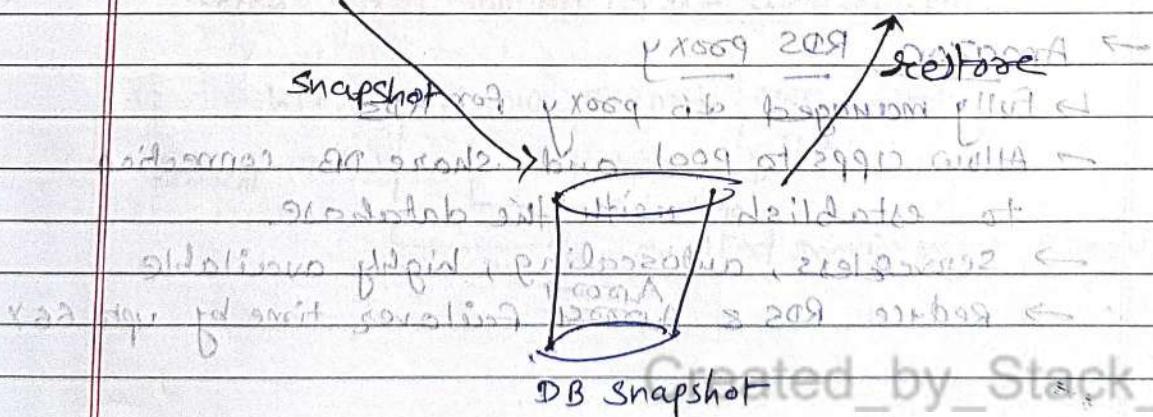
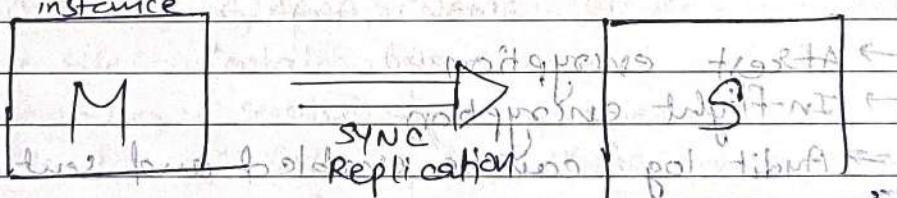
Use case:-

- you have a production db that is taking on normal load
- you want to run reporting application to run some analytics.
- you can create a new read replica to run the new workload there.

- Production workload unaffected
- used only select query
- there is no cost when read replicas happen within same region
- there is a cost when read replicas happen at cross-region.

### \* RDS Multi AZ

- SYNC replication
- one DNS name - automatic failover to stand
- Increase availability
- No manual intervention in apps.
- No need for scaling
- Single AZ to multi AZ is good
- Zero downtime operation



## \* Amazon Aurora

- Aurora is a proprietary technology from AWS.
- PostgreSQL and MySQL are both supported as AuroraDB.
- Aurora is AWS cloud optimized.
- ↳ Claims 5x performance improvement over MySQL.
  - ↳ 3x PostgreSQL on RPS
- Aurora storage automatically grows in increments of 10GB, up to 128TB.
- Aurora can have upto 15 replicas and it faster.
- Failover in Aurora is instantaneous.
- Cost more - but more effective.
- ↳ Features of Aurora
  - ↳ Automatic failover, Isolation and Security
  - ↳ Backup & Recovery, Industry Compliance
  - ↳ Push-button scaling, Advance Monitoring
  - ↳ Backtrack, maintenance windows.
- At-rest encryption
- In-flight encryption
- Audit log can be enabled and sent to Cloudtail.

## → Amazon RDS proxy

- ↳ Fully managed db proxy for RDS
- Allows apps to pool and share DB connection to establish with the database.
- Serverless, autoscaling, highly available
- Reduce RDS & Aurora failover time by upto 66%

## \* Amazon ElastiCache

EE. 3109 \*

- Same way Redis is to get managed Relational Database.
- ElastiCache is to get managed Redis or Memcached.
- In memory database help in high performance
- Low latency. At least 1000x faster than RDBMS.
- Help reduce load off db for read intensive
- make app "stateless".
- take care of os maintenance / patching, optimization, setup config, monitoring, feature recovery
- Turned - heavy app code changes.

AAAA → A.09 - 09T domain

B.09 - 09T subdomain

## \* what is DNS

→ Domain name System which translate the human friendly hostname into the machine IP addresses.

→ backbone of the internet.

→ Pageant domain - www - AAAA

→ Domain Registrar - Amazon Route 53, GoDaddy

→ DNS record - A, AAAA, CNAME, etc

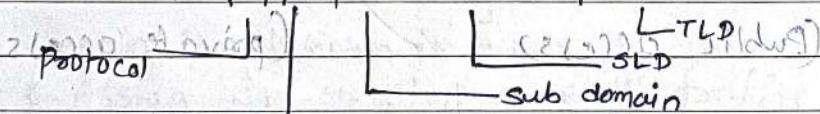
Zone file :- Contain DNS record

Name Server :- resolve DNS queries

Top level Domain (TLD) :- .com, .us, .in, .gov, .org, etc

Second level Domain (SLD) :- amazon.com.

DNS http://api.www.example.com. → Root



Fully Qualified Domain Name (FQDN)

## \* Route 53

- A highly available, scalable, fully managed and Authoritative DNS.
- Route 53 is also known as Domain Registrar.
- Ability to check the health of your resources.
- The only AWS Service which provides 99.9% availability SLA.

- Route 53 - Records :-
- ↳ Domain / subdomain Name :- e.g example.com
- ↳ Record Type :- e.g. A or AAAA
- ↳ Value
- ↳ Routing policy
- ↳ TTL

- DNS type :- A / AAAA / CNAME / NS
- A - maps a hostname to IPv4
- AAAA - maps a hostname to IPv6
- CNAME - maps a hostname to another hostname
- NS - Name servers for the Hosted zone.

- Hosted zones :-
- |                    |                     |
|--------------------|---------------------|
| Public Hosted zone | Private Hosted zone |
| (Public access)    | (Private access)    |

Route 53 - records TTL - similar pointing to

↳ High TTL (24 hr) work out of geographical

↳ less traffic w/ at 24 hr period

↳ Possibly outdated records longer expire

w/ rotation effect no bound in protocol

↳ Low TTL (60 sec)

↳ More traffic on Route 53 need to wait

↳ Records are outdated for less time

↳ easy to change records

→ Mandatory when using CloudFront & TTL

-> resolving from last location <-> first location

\* Routing policies - simple

(current or 2nd, record, 99%)

→ Typically, route traffic to single resources.

→ can specify multiple <sup>value</sup> records in same record.

→ If multiple values are returned, a random one is chosen by the client.

→ When alias enabled, specific only one AWS resource

→ Highly leverages health checks if two or more regions

\* Routing policies - weighted

(100%, 50% - 100%)

→ Control the % of the requests that go to each specific

→ DNS records must have same name and type.

→ associated with health check.

→ Use cases - balancing between region, testing new

application versions, failover configuration

→ Assign no to stop sending traffic to region

## \* Routing policies - Latency-based

- Redirects to the resources that has the least latency close to us.
- Super helpful when latency for user is a priority.
- Latency is based on traffic between user and AWS <sup>Regions</sup>.
- Has ~~is~~ associated with health check.

## \* Route 53 - Health Checks

- HTTP health checks are only for public resources.
- Health checks  $\Rightarrow$  Automated DNS Failover:-
- ① Health checks that monitor an endpoint (app, server, AWS resources)
- ② Health checks that monitor <sup>other health checks</sup> an endpoint
- ③ Health checks that monitor CloudWatch Alarms
- Health checks are integrated with CW metrics.
- About 15 global health checks will check the endpoint health.
- Threshold - 3 (default)
- Interval - 30 seconds
- Supported protocol:- HTTP, HTTPS and TCP
- If  $> 18\%$  healthy
- Configure your route 53 firewall to allow incoming requests from Route 53 Health Checkers.

## \* Route 53 - calculated Health checks

- combine the results of multiple Health checks into a single health check.
- you can use OR, AND or NOT.
- can monitor upto 256 child Health checks.
- Specify how many of the Health checks need to pass to make the parent (pass) - NOT OF
- perform maintenance to your website without causing all health checks to fail.
- Health checks - private Hosted zones
  - health checks outside the VPC in other regions.
  - They cannot access private endpoint.
  - you can create CloudWatch Metrics and associate a CloudWatch Alarm to it. Create a Health Check that checks the alarm itself.

## \* Routing policy - failover - failing pointing

- Routing policies - Geolocation
  - Different from latency based.
  - The routing is based on user location.
  - Specify the location by continent, country or by US state.
  - Should create a default record (if not available specify).
  - Use case :- website localization, restrict content distribution, load balancing.

## \* Geography Routing policy

- Route traffic to your resources based on the geographic location of user and resources.
- Ability to shift more traffic to resource based on the defined bias.
- To change the size
  - To extend - (1 to 99) extend traffic
  - To shrinks - (-1 to -99) shrink traffic
- we use Route 53 traffic flow.

## \* Routing policies - IP-based Routing

- Routing is based on clients IP addresses
- You provide a list of CIDR's for your clients and corresponding endpoint location.
- Use cost :- optimize performance, reduce cost

- Routing policies - Multi Value
  - Use when routing traffic to multiple resources
  - Route 53 return multiple values
  - use & healthy records are returned for each multi value query
  - Multi value is not substitute for having an ELB

## \* VPC & Subnets primer

Primer 29V 3K

- VPC :- private network to deploy your resources
- Subnets :- allows you to partition your network inside your VPC
- A public subnet is a subnet that is accessible from the internet
- A private subnet is a subnet that is not accessible from the internet
- You can use Route Table to access internet & subnets
- Internet Gateway helps our VPC instance connect with the internet. public subnet have route to the internet gateway.
- NAT Gateway & NAT Instance allows your instance in your private subnets to access the internet while remaining private.

## \* Network ACL & Security Groups

- NACL (Network Access Control List) :- a firewall which controls traffic from and to subnet
- Can have allow and deny rules.
- Can be attached at subnet level / only include IP.
- Security Groups
- A firewall that controls traffic to & from an ENI / an EC2 Instance.
- Can only allow rules / Rules includes IP address and other security groups.

## \* VPC Peering

- Connect to VPC privately using AWS network
- Make them behave as if they were in the same network
- Must not have overlapping CIDR

## \* VPC Endpoint :-

- Endpoints allow you to connect to AWS service using a private network instead of the public network
- This gives you enhanced security and lower latency to access AWS service.
- VPC Endpoint Gateway - S3 & DynamoDB
- VPC Endpoint Interface - the restational part
- Only use within your VPC

## site-to-site VPN

- Connect on-premises VPN to AWS
- The connection is automatically encrypted.
- Goes over public internet

\* Direct (connect) ATA 25pin to 25pin cable

- Establish a physical connection between on-premises and AWS
  - The connection is private, secure and fast
  - Goes over private networks
  - Takes at least month to establish.
  - VPC Flow Logs = network traffic logs

\* LAMP stack on AWS

ER Linux, Apache, MySQL, PHP, Python

\* Amazon S3 (DA) till lasten 2020A tills

- main building block of AWS
  - ⇒ use cases :- Backup and storage, disaster recovery, Archive, Hybrid Cloud Storage, Application Hosting, media hosting and many more.
  - 
  - S3 allows people to store objects (files) in buckets
  - Bucket must have a globally unique name
  - Buckets are defined at regional level.
  - Naming convention:- No uppercase, No underscore, 3-63 character long; Not an IP, Must start with lowercase or number, Not start with suffix and not end with prefix suffix.
  - Objects have a key.

- Max object size is 5TB (5000GB)
- If uploading more than 5GB must use "multi-part upload."

## \* Amazon S3 Security

- User-based
  - ↳ IAM policies - which API called should be allowed for a specific user from IAM
- Resource-Based
- Bucket policies - bucket wide rules from the S3 console - allows cross account
- Object Access Control list (ACL) - fine grain
- Bucket Access control list (ACL) - less common
  - ↳ encrypts object in S3 using Encryption Keys.

## \* Amazon S3 - Replication (CRR-SRR)

- Must enable versioning in source and destination buckets
- Cross-Region Replication (CRR)
- Same-Region Replication (SRR)
- Bucket can be in different AWS account
- Copying is asynchronous
- Must give proper IAM permission to S3

Use case :-

HOB \*

CRR:- compliance, lower latency access, replication across accounts.

SRR:- log aggregation, live replication between production & test accounts.

Timeline starts 29A

- only new object is replicated, so we need to use s3 batch replication for existing one.
- can replicate delete markers from source to target
- There is no chaining of replication

## \* EC2 Instance Metadata (IMDS)

AWS EC2 Instance Metadata (IMDS) is powerful but one of the least known features to developers.

Allows EC2 to learn about themselves.

Using URL resolution of metadata of node.

Metadata = info about EC2 instance

User data = launch script of the EC2 instance.

MFA with CLI is not possible by default.

To use MFA with the CLI, you must create a temporary session.

To do so you must run the STS GetSessionToken API.

aws sts get-session-token and run a loop

## \* SDK

- Use to perform action on AWS using Software development kit.
- language support :- Java, .NET, Node.js, Python, Go...
- API Rate limit
- Service Quotas (Service limit)
- Exponential Backoff
  - ↳ If you get ThrottlingException intermittently, use exponential backoff.

## \* Amazon S3 Life Cycle Rules

- Transition actions :- configure objects to transition to another storage class
- Move objects to Standard IA class after creation, at least of 90 days
- Move to Glacier for archiving after 6 months.
- Expiration actions:- configure objects to expire (delete) after some time.
  - Access log file send to delete after 365 days
  - can be used to delete old versions of files
  - delete incomplete multi part uploads
- Rules can be created for the certain prefix
- Rules can be created for certain objects Tags.

## → S3 Event Notification

→ Use case:- generate thumbnails of image uploaded to s3

→ S3 Performance - via origin with  $\Rightarrow$

↳ Multipart Upload

↳ S3 Transfer Acceleration

## $\Rightarrow$ S3 Select & Glacier Select

→ Retrieve less data using SQL by performing server-side filtering.

→ Can filter by rows & columns

→ less network transfer, less CPU cost client side

## \* Amazon S3 - Object Encryption

→ you can encrypt object in S3 using 4 methods

(1) Server-side Encryption (SSE)

→ Server-side Encryption with Amazon S3 Managed keys (SSE-S3)

→ Server-side Encryption with KMS key stored in AWS KMS (SSE-KMS)

→ Server-side Encryption with Customer provided keys (SSE-C)

## (2) Client side Encryption

## \* Cross origin Resource Sharing (CORS)

- Origin = Scheme + Host + Port
- Web browser based mechanism to allow requests to other origins while visiting the main origin.

→

## \* S3 Access log

- For audit purpose, you may want to log all the access to S3 buckets.
- Any type of S3 request is logged into another S3 bucket.
- That data is going to analyze.
- must be in same AWS region.
- Not same logging bucket & monitoring bucket

→

## \* S3 Object Lambda

- Use Aws Lambda function to change the object before it is retrieved by the caller application.

## \* Aws Cloud Cloudfront

- works at monitoring & managing global web sites (IP address)
- content delivery Network (CDN)
- Improve read performance, content is cached at the edge.
- Improves user experience by reducing latency
- 216 point of presence (globally) (edge location)
- DDoS protection, Integration with CloudFront Shield, AWS web app firewall.
- Great for static content that must be available everywhere

## \* Cloudfront Caching

- The cache lives at each Cloudfront Edge Location.
- Cloudfront identifies each object in the cache using the cache key.
- Cloudfront ALB or EC2 as an origin

## \* Cloudfront Signed URLs / Signed Cookies

## \* Cloudfront - Price classes

- You can reduce the number of edge location for cost reduction

- There are three classes

① Price Class All :- all regions

② Price Class 200 :- most regions & not expensive one

③ Price Class 100 :- only the least expensive

## \* Cloudfront - Field Level Encryption

## \* Cloudfront - Real Time Logs

## \* Docker

- Docker is a development platform to deploy apps
- Apps are packed in containers that can be run on any OS.
- Apps run the same, regardless of where they run
- Any machine runs your code
- No compatibility issue to things like
- predictable behaviour
- less work
- easier to maintain and deploy
- work with any language, any OS and any technology

Use case :- microservices architecture, lift and shift apps from on-premises to cloud

## → Docker Image

- Docker Image stored in Docker repositories
- Docker Hub :-
- ↳ public repository
- ↳ private repository
- ↳ public repository
- Docker is sort of a virtualization technology, but not
- Resource are shared with the host ⇒ many containers are on one server

## → Docker Container Management on AWS

- Amazon Elastic container service (Amazon's own container platform) →
- Amazon Elastic Kubernetes Service (EKS) →  
Manage Kubernetes
- AWS Fargate →
- Amazon Serverless container platform →
- work with ECS & EKS, important →
- Amazon ECR →
- store container image →

## ≠ Amazon ECS - EC2 Launch Type

- Elastic Container Services →
- Launch Docker container on AWS = Launch ECS Task on ECS Clusters.
- EC2 Launch type :- you must provision & maintain the infrastructure
- Each instance must run the ECS Agent to register in ECS Cluster
- AWS will take care of starting/stopping containers

## → Amazon ECS - Fargate Launch type

- Launch Docker containers on AWS Lambda
- You do not provision the infrastructure
- It's all serverless
- You just create task definition
- AWS just runs ECS Task for you based on the CPU/RAM you need.
- To scale, just increase the no of tasks.
- EC2 Instance profile
- used by ECS agent

## → Amazon ECS - Data Volume

- Mount EFS file system onto ECS Task
- Work for both EC2 & Fargate
- Fargate + EFS = Serverless
- Use cases :- persistent multi AZ shared storage for your containers

## → ECS Service Auto Scaling

- Automatically increase / decrease the desired number of ECS task

## → Amazon ECS Auto scaling user AWS Application Auto Scaling

- ↳ Target scaling → Scale based on target
- ↳ Step Scaling → Scale based on specified CloudWatch Metrics
- ↳ Scheduled Scaling → scale based on specific date / time.

- Auto scaling group scaling
- ECS Cluster capacity provider

## \* ECS Task Placement

- When a task of type EC2 is launched, ECS must be determine where to place it, with the constraint of CPU, memory and available port.
- Similarly when a service scales in, ECS need to determine which task to terminate.

- So we use task placement and task placement constraints

- Bin packing:  
place task based on the least available amount of CPU or memory.
- The minimize the number of instances in use

- Random:  
Place the task randomly.

↳ Place the task randomly.

- Spread:  
place the task evenly based on specified value

↳ place the task evenly based on specified value

- you can also mix them together.

→ distinctInstance :- place each task on different container image.

→ memberOf :- place task on instance that satisfy an expression.

### \* Amazon ECR = Amazon Elastic Container Registry

→ ECR = Elastic Container Registry

→ Storage and manage docker image on AWS

→ private and public repository

→ Fully integrated with ECS, backed by Amazon S3

→ Access is controlled through IAM

→ Support Image Vulnerability Scanning, versioning image tags, image lifecycle

### \* AWS Copilot = AWS Service for containerizing and running your apps

→ CLI tool to build, release, and operate production-ready containerized apps.

→ Run your apps on Apprunner, ECS and Fargate

→ Help you focus on building apps not more on infrastructure

→ Automated deployment with one command using Codepipeline

→ Deploy to multiple environment

→ Troubleshooting, logs, health status

## \* Amazon EKS

- Elastic Kubernetes service
- way to launch managed Kubernetes cluster on AWS
- Kubernetes is a open source system for automatic deployment, scaling and management of containerized app.
- alternative of ECS, diff API
- use case :- if your company already use Kubernetes on-premises or in another cloud and you want to migrate to AWS using Kubernetes

## \* Elastic Beanstalk

- Elastic Beanstalk is a developer centric view of deploying an application on AWS.
- It uses all the components like EC2, ASG, ELB, RDS...
- Managed services:-
  - ↳ Automatically handle capacity provisioning, load balancing, scaling, application health monitoring,
- Just the application code is responsibility of developer
- have full control over configuration
- Beanstalk is free but you pay for the underlying instances.
- Support many languages & platforms.  
(Custom platform)

- Beanstalk Deployment option for updates
  - All at once - fastest, but instances aren't available to serve traffic for a bit
  - Rolling :- Update few instance at a time and then move onto the next
  - Rolling with additional batches :- like rolling but spins up new instances to move the batch.
  - Immutable :- Spins up new instance in new ASG, deploy version to this and then swap when all the things are healthy.
- Blue Green :- Create a new environment and switch over without ready.
- Traffic Splitting :- Carry testing - send a small % of traffic to new development.
- Elastic Beanstalk CLI (CloudFormation) :-

## \* AWS Cloudformation

- Cloudformation is a declarative way of outlining your AWS Infrastructure, for any resources
- Then cloudformation creates those for you, in the right order with the exact configuration that you specify.
- Infrastructure as code :-
  - ↳ No resources are manually created, which is excellent to control
  - ↳ The code can be version controlled eg. git
    - ↳ changes are reviewed

### Cost :-

- Each resources within the stack is tagged with an identifier so you can easily see how much a stack cost you.
- estimate a cost using (template)
- Saving strategy :- delete template at 5pm and create template at 8am

### Productivity :-

- Ability to destroy & recreate as infrastructure on the cloud
- Automated generation of Diagram for your template
- Declarative programming,

## → Separation of concern

- create many stacks for many apps, and many layers  
Ex:- Svc stacks, Network stacks
- Don't re-invent the wheel
  - Leverage existing templates on the web
  - Leverage the document.

## → YAML File

= =

- YAML & JSON are the languages you can use for Cloud formation.

→ Key Value Pairs

→ Nested objects

→ Support Array

→ Multi-line strings

→ Can include comments

## → Resources

→ Resources are the core of your CloudFormation template.

→ Represent diff AWS components and config

→ Resources are declared and can reference each other.

→ There are over 224 types of resources

→ type identification

AWS:aws-product-name::data-type-name

## → Parameters

→ Parameters are the way to provide input to your AWS CloudFormation template.

## → Outputs :-

→ The output section declares optional output values that we can import into other stacks.

→ You can view the outputs in AWS Console & CLI.

→ Useful for ex:- you define network

CloudFormation and output the variables such as VPC ID and your subnet IDs.

→ It is a best way to perform some collaboration cross stack, as you let except handle their own part of the stack.

→ You can delete CloudFormation stack if its outputs are being referenced by another CloudFormation stack.

→ Fn::Join

→ Fn::Ref

→ Fn::GetAtt

→ Fn::FindInMap

→ Fn::ImportValue

→ Function Fn::Sub

→ Conditions Function

↳ Fn::And

Fn::Equals

Fn::If

Fn::Or

Fn::Not

→ Cloud formation Rollbacks

→ Stacks creation fails :-

↳ everything rollback (get deleted)  
we can look at the log.

→ Stack update Fails :-

↳ The stack automatically rollback to the previous known working state.

→ Ability to see in the log what happened and error message

→ Cloud formation drift

\* Amazon SQS - standard queue service

→ oldest offering

→ fully managed services, decoupled application

Attribute :-

→ unlimited throughput, unlimited no of msg in queue

→ default retention of msg : 4 days, max of 14 days

→ low latency → limitation of 256 KB per msg sent

→ can have duplicate message

→ can have out of order message (best-effort ordering)

→ SQS have unlimited throughput

## \* SQS - Message Visibility Timeout

→ After the message is polled by a consumer, it becomes invisible to other consumers.

→ Default message visibility timeout is 30 seconds.

## → Dead Letter Queue (DLQ)

→ If a consumer fails to process a msg within a visibility timeout, the msg goes back to the queue.

→ We set threshold of how many times it goes back to queue and after that it goes to dead letter queue (DLQ).

→ DLQ is FIFO queue

→ DLQ of a standard queue

→ Make sure to process the msg before expire → expire retention of 14 days in DLQ

## → Amazon SQS - Delay Queue

→ Delay a message upto 15 min (consumer don't see immediately)

→ Default is 0 seconds

→ Can set a default at queue level

→ Can override the default on send using the Deploy DelaySeconds parameter

## \* Amazon SQS - Long polling

- When a consumer requests message from the queue, it can optionally "wait" for message to arrive if there are none in the queue.
- This is called Long polling
- Long polling decreases the number of API calls made to SQS while increasing the efficiency and decreasing the latency of your application.
- The wait time can be 1 sec to 20 sec
- Long polling is preferred to short polling
- Enabled at queue level or API level using `ReceiveMessage.WaitTimeSeconds`

## \* SQS Extended Client

- Message size limit is 256 KB, how to send large message e.g. 1 GB?
- Using the SQS Extended Client (Java library)
- we save in S3 bucket and send metadata message with S3 bucket link.

## \* SQS - Must know API

- `CreateQueue(MessageRetentionPeriod)`, `DeleteQueue`
- `PurgeQueue` :- delete all the message in queue
- `SendMessage(Delay seconds)`, `ReceiveMessage`, `DeleteMessage`

- Max Number Of Messages :- default 1, max 10
- Recieve Message Wait Time Second : Long polling
- ChangeMessageVisibility : change the Msg timeout
- Batch APIs For Send Message, DeleteMessage, ChangeMessage
- Visibility help in decrease in cost

### \* Amazon SQS - FIFO Queue

- FIFO = First In First Out (order in the queue)
- limited throughput : 300 msg/s without batching  
3000 msg/s with
- Exactly once send capability (by removing duplicates)
- Message are processed in order by the consumer

### \* SQS FIFO - Depli Deduplication

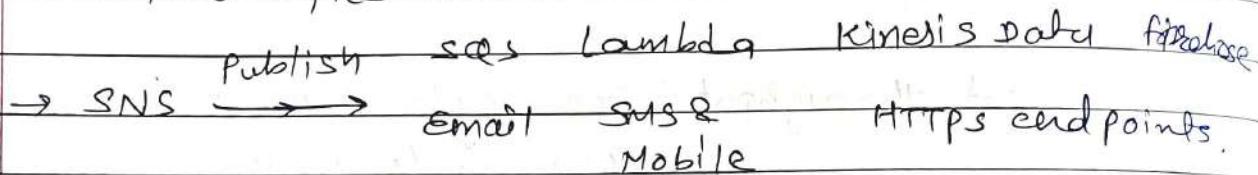
- De-duplication interval is 5 minutes
- Two method :-
- ↳ Content-based deduplication: will do a SHA-256 hash of the message body.
- ↳ Explicitly provide a Message Deduplication ID.

### \* SQS FIFO - Message Grouping

- If you specify the same value of ~~Message Group ID~~ in an SQS FIFO queue, you can only have one consumer and all the messages are in order.
- To get ordering at the level of a subset of msg, Specify different values for Message Group ID
- Each Group ID can have a diff consumer
- ordering across group is not guaranteed.

## \* Amazon SNS

- Pub / sub pattern
- The event producer only send message to one sns topic.
- As many "event receivers" (subscription) as we want to listen to the sns topic notification.
- Each subscriber to the topic will get the msg
- upto 12,500,000 subscription per topic
- 100,000 topics limit.



- Topic publish (using SDK)
- Direct publish (for mobile app SDK)

### SNS Security :-

- Encryption
  - In flight encryption using HTTPS API
  - At rest encryption KMS keys
  - Client side encryption

- Access control → IAM control
- SNS Access policies.

\* SNS + SQS : Fan out

- Push once in SNS, receive in all SQS queues that are subscribers
- fully decoupled, no data loss
- SQS allow for :- data persistence, delayed, processing and priorities of messages
- Ability to add more SQS subscribers over time
- Make sure your SQS queue access policy allows for SNS to receive.
- Cross Region delivery : SQS queues in other region

- S3 Events to multiple queues
- SNS to Amazon S3 through Kinesis Data Firehose
- SNS have FIFO capability
- SNS FIFO + SQS FIFO : Fan out
- SNS Message Filtering

\* Kinesis

- Make it easy to collect, process and analyze streaming data in real time.
- Ingest real-time data such as Application logs, Metrics, website clickstream, IoT telemetry data
- Kinesis Data Stream → Capture, Process & store
- Kinesis Data Firehose → load data into AWS data store
- Kinesis Data Analytics → analyze data stream using SQL or Apache Flink
- Kinesis Video Streams → Capture, process and store video streams.

## Kinesis Data Streams

- Retention between 1 day to 365 day
- Ability to Reprocess data.
- Once data is inserted in Kinesis, it cannot be deleted.
- Data that shared the same partition goes to the same shard.
- Producers :-
- Producers :- AWS SDK, Kinesis producer library (C/C++)
- Consumers :- write your own - KCL, AWS SPK
- Managed :- AWS Lambda, KDF, KDA

## Capacity Mode

- Provisioned Mode :-
- You can choose no of shards provisioned, scale manually or using API.
- Shards get 1 MB/s in (1000 records/ps)
- Shards get 2 MB/s out
- Pay per hours.

## On demand mode :-

- No need to provision or manage the capacity
- Default capacity provisioned (4 MB/s in or 4 MB/s out)
- Scaled automatically based on last 30 days
- Pay per stream per hour & data in/out (22 MB/s)

## Kinesis producers

→ put data records into data streams

→ Data record consist of:-

↳ sequence no., partition key, Data blob

→ Producers:-

↳ AWS SDK

↳ Kinesis producer library

↳ Kinesis Agent

→ write throughput:- 1 MB/sec or 1000 records/sec per shard

→ put record API

→ Use Batching with putRecords API to reduce costs & increase throughput.

## Kinesis consumers

→ Get data records from data streams and process them

→ AWS Lambda

→ Kinesis Data Analytics

→ Kinesis Data Firehose

→ Customer Consumers (AWS SDK)

→ Kinesis client library (KCL)

## Shared classic fan-out consumer - Pull

→ low number of consuming application

→ Read throughput: 2 MB/s per shard across all consumer

→ Max 5 GetRecords API call per sec

→ latency ~200ms

→ Minimum cost (\$)

→ Consumer pull data from Kinesis using GetRecords API

→ Return upto 10 MB or 10000 records

→ Enhanced Fan-out Consumers - Push

→ Multi consuming application for the same stream

→ 2 MB/sec per consumer per shard

→ Latency ~ 70 ms

→ Higher costs (\$\$\$)

→ Kinesis pushes data to consumers over HTTP/1.2

→ soft limit of 5 consumers app per data stream

⇒ Kinesis Client Library (KCL)

→ A Java library that help read records from a Kinesis data stream with distributed app sharing the read workload.

→ Each shard is to be read by one KCL instance

→ 4 shards → 4 KCL

6 shards → 6 KCL

→ process is checkpointed into DynamoDB

→ Track other workers and share the work amongst shared using DynamoDB

→ KCL can run on EC2, Elastic Beanstalk and on-premi

→ Records are reads in order at shard level

→ version → KCL 1.x (shared consumer)

KCL 2.x (shared & enhanced fanout)

## \* Kinesis Data Firehose

- Fully managed service, no administration, automatic scaling, serverless → AWS - Redshift, S3, OpenSearch
- 3<sup>rd</sup> party - Splunk, MongoDB, DataDog.
- Send to any HTTP endpoint
- pay for data going through firehose.
- near real time
- 60 sec latency
- Supports many data formats, conversions, transformations, compression.
- support custom data transformation using AWS Lambda
- can send failed or all data to a backup s3 bucket

## \* AWS Monitoring, Troubleshooting & Audit

- AWS CloudWatch :-
- Metrics :- Collect and track key metrics
- Logs :- collect, monitor, analyze and store log files.
- Events :- send notifications when certain event happen in AWS
- Alarm :- React in real-time to metrics / events.
- AWS X-Ray :-
- Troubleshooting application performance and errors
- Distributed tracing of microservices.
- AWS CloudTrail :-
- Internal monitoring of API calls being made
- Audit changes to AWS Resources by user.

## \* AWS CloudWatch Metrics

- CloudWatch provides metrics for every services in AWS
- Metrics is a variable to monitor
- Metrics belongs to namespaces
- Dimension is an attribute of an metric
- Metrics have timestamps
- CloudWatch dashboard for metrics
  
- Cloudwatch Custom Metrics
  - Possibility to define and send your own custom metrics to CloudWatch.
  - Ex:- RAM usage, disk space, number of logged in users
  - Use API calls putMetricData
  - Ability to use dimensions → Instance id.
    - Environment.name
  - Metric resolution → Standard
    - High Resolution
  - Important: Accept date point two weeks in past and two hours in the future.

## → CloudWatch logs

- Log groups :- arbitrary name, usually representing an application.
- Log stream :- instances within application / log files / containers
- can define log expiration policy
- log send to → S3, KDS, HDFS, Lambda, opensearch
- Logs are encrypted by default.
- can setup KMS-based encryption with your own keys
- Source :- SDK, Cloudwatch log agent

## → Cloudwatch logs for EC2

- By default, no logs from your EC2 machine will be sent to Cloudwatch
- You need to run a Cloudwatch agent on EC2 to push the log file you want.
- Make sure IAM permission are correct
- The Cloudwatch log agent can be setup on-premise too.

## → Cloudwatch Unified Agent

- CPU, Disk Metrics, RAM, Netstat, Processes, Swap Space

## → Cloudwatch Log Metric Filter

- Cloudwatch logs can use filter expression ex:- find specific IP, logs, errors

- generate point after filter <sup>Created</sup>

- ability to specify upto 3 dimensional metric filter.

## → Cloudwatch Alarms

- Alarm are used to trigger notification for any metric
- various option
- Alarm status :- OK, INSUFFICIENT DATA, ALARM

## → Period :-

- Length of time in second to evaluate the metric
- 10 sec, 30 sec, multiples of 60 sec.
- send notification to sns

## → Cloudwatch Synthetics Canary

- Configurable script that monitors your APIs, URLs, website
- Reproduce what your customer do programmatically to find issue before customers are impacted.

→ Integrated with Cloudwatch alarm

→ Script written in Node.js or Python

→ Programmatically access to handles @ chrome

→ Run on schedule

## → Heartbeat Monitor

### → API Canary

### → Broken Links checker

### → Visual Monitoring

### → canary Reorder

### → GUI Workflow Builder

## → AWS EventBridge

- Schedule cron jobs (scheduled scripts)
- Event pattern :- Event rules to react to a service doing something.
- Trigger lambda function, send SNS/SQS message
- Event bus can be accessed by other AWS accounts using Resource based policies.
- You can archive event
- Ability to archive & replay archived events.
- Schema - Registry
  - Event bus can analyze the events in your bus and infer the schema.
- It helps to generate code
- It is versioned
- Resource based policy
  - Manage policy permissions for a specific event bus.
  - Ex:- allow/deny events from another AWS account or AWS region.
- Use case:- aggregate all the events from AWS organization.

## \* AWS X-Ray

- visual analysis of your application
- advantage :-
  - Troubleshooting performance
  - understand dependencies in a microservice architecture
  - pinpoint service issue
  - Review request behavior
  - Find error & exception
  - identified user that are impacted.
- Compatibility :-
  - Lambda, Elastic Beanstalk, ECS, ELB, API Gateway, EC2 Instance
- AWS X-ray leverage Tracing
  - Tracing is an end to end way to follow a request
  - Tracing is made up of Segment (sub segment)
- X-Ray Instrumentation in your code
  - Instrumentation means the measure of product performance, diagnose errors and to write trace information.
  - To instrument your application code, you use the X-Ray SDK.

## X-Ray Concept

1. Segment :- each application / service will send them

→ SubSegments :- if you need more details in segments

→ Trace :- segment collected together to form an end-to-end trace.

→ sampling :- decrease the amount of request send to X-Ray reduce cost

→ Annotations :- key value pairs used to index traces and use with filter.

→ Metadata :- key value pairs, not indexed, not use for searching.

→ X-Ray daemon / agent has a config to send traces cross account.

## AWS Distro for Open Telemetry

→ secure, production-ready AWS-supported distribution of the open source project Open Telemetry project.

→ provide single set of API, library, agent and collector services.

→ Auto-instrumentation Agents to collect traces without changing your code.

## \* AWS Cloud Trail

- Provide, governance, compliance and audit for your AWS account.
- CloudTrail is enabled by default
- Get an history of events / API calls made <sup>within</sup> by your AWS account
- Console, SDK, CLI, AWS services
- Can put logs from CloudTrail to CloudWatch logs or S3
- A trail can be applied to All Regions or Single region
- If all the resources deleted in AWS, investigate CloudTrail first.
- Events :-

- Management Events :-
- Data Events :-
- CloudTrail Insights
- CloudTrail Insights to detect unusual activity in your accounts
- Stored for 90 days normally
- Can more then send to S3

## \* AWS Serverless

- serverless is new paradigm in which the developers don't have to manage servers more.
- They just deploy code, just function
- initially serverless == FaaS (Function as a service)
- before there is only AWS Lambda but now database, messaging, storage etc.
- It does not mean no server it means you do not need to manage it.

## → AWS Lambda

- virtual functions - no server to manage
- limited by time - short executions
- Run-on demand
- Scaling is automated.

### → Benefits :-

- Easy pricing → Pay per use
- Integrated with the whole AWS services
- Many programming language support
- Easy monitoring
- Easy to get more resources.
- RAM ↑ → CPU and network ↑
- language support
  - Python, Node.js, Java, C#, Ruby, (with) Runtime API

## → Lambda Container Image

- The container image must implement by Lambda API
- ECS / Fargate

## → Lambda - Synchronous Invocations

- Synchronous :- CLI, SDK, API Gateway, ALB
- Return its result result is returned right away.
- Error handle must happen client side.
- User Invoked :-
  - ↳ GLB, API Gateway, CloudFront,
- Service Invoked :-
  - ↳ Amazon Cognito, Step Function
- Other services :-
  - ↳ Lex, Alexa, Kinesis Data Firehose

## → Lambda Integration with ALB

- To explore a lambda function as an HTTP endpoint
- you can use ALB or API Gateway.
- Lambda fun<sup>n</sup> must be register in target grp.

## → Lambda - Asynchronous Invocations

→ S3, SNS, CloudWatch Events

→ The event is placed at event queue.

→ Lambda attempts 3 tries total, 1 min after 1<sup>st</sup>  
2 min after 2<sup>nd</sup>

→ Make sure the processing is idempotent

→ Services →

→ S3, SNS, CloudWatch Event, EventBridge, CodeCommit  
CodePipeline

## → Lambda - Event and context objects

### → Event Objects

- JSON formatted document contain data from the ~~function~~ function to process
- Contain "info" from the invoking service
- Lambda runtime converts events to an object.
- Ex:- input arguments, invoking series arguments

### → Context object

- Provides properties and methods that provide the information about the invocation, function and runtime environment.
- passed to your ~~f~~ "f" by lambda at runtime

## → Lambda - Destination

=      =

- can config to send results to destination
- Asynchronous invocation: can define <sup>destination</sup> definition for successful and failed event :-
  - AWS SQS, SNS, Lambda & EventBridge bus
- Now AWS recommends you use destination instead of DLQ now
- Event source mapping : for discarded events batches
- SQS, SNS
- you can send event to a DLQ directly from SNS

## → Lambda Execution Role (IAM Role)

→ Grants the lambda function permission to access services/resources.

### → Policies :-

→ AWSLambdaBasicExecutionPolicy <sup>Role</sup>

→ " " KinesisExecutionRole

→ " " DynamoDBExecutionRole

→ " " SQSEueueExecutionRole

→ " " VPCAccessExecutionRole

→ " xRay Daemon writeAccess

→ When we use event source mapping, Lambda uses execution role to read event data.

→ Create Lambda execution Role per function

## → Lambda Resource Based policies

→ Use resource based policies to give other accounts and AWS services to use Lambda resources

→ give access to all the other services

## → Lambda Environment Variable

- Environment Variable = key / value pair in string
- Adjust the fn behavior without updating code to run
- environment variable are available at code.
- Lambda services add its own system environment variable as well.
- Helpful to store ~~keys~~ <sup>secrets</sup>.
- Secrets can be encrypted

## → Lambda logging & Monitoring

### → Cloudwatch logs :-

- AWS Lambda execution logs are stored in AWS Cloudwatch log.

### → Cloudwatch Metrics :-

- AWS Lambda metrics are displayed in the AWS Cloudwatch Metrics.

### → Lambda Tracing with X-Ray.

- Enable in Lambda configuration
- Run the X-Ray Daemon for you
- environmental variable to communicate with X-Ray
  - X\_AMZN\_TRACE\_ID
  - AWS\_XRAY\_CONTENT\_MISSING
  - AWS\_XRAY\_DAEMON\_ADDRESS

## → Lambda@Edge

- Lambda function is written in Node.js & Python
- Scales to 1000s of requests/second.
- Used to change CloudFront request and response.
- Author your function in one region and replicate to its location.

## \* Lambda in VPC

- You must define VPC ID, the subnet and the security.
- Lambda will create an ENI (Elastic Network Interface) in your subnet.
- AWSLAMBDA VPC Access Execution Role.
- Deploying a Lambda function in private subnet gives it internet access if you have NAT gateway/Instance.
- Use VPC endpoint to privately access to AWS services.

## → Lambda Function Configuration

- If your application is CPU bound (computation heavy) increase RAM.
- Timeout:- default 3 second, maximum is 900 sec.

→ Lambda layers

→ Custom Runtimes

→ Externalize Dependencies to reuse them

→ Lambda - File System Mounting

→ Lambda function can access EFS file system if they are running on VPC.

→ Config Lambda to mount EFS file system to local directory during initialization

→ Must leverage EFS access point.

→ Limitation:- watch out of the EFS connection limit

→ Lambda Concurrency & Throttling

→ Concurrency limit :- upto 1000 concurrent execution

→ can set "reserved concurrency" at the function level

→ Each invocation over the concurrency limit will trigger a throttle.

→ If you need higher limit, open a support ticket

→ Cold Start :-

→ Provisioned Concurrency :-

## → Lambda and CloudFormation - inline

→ inline functions are very simple.

→ Use the `Code.ZipFile` property

→ you cannot include function dependencies with inline functions.

→ other option using S3.

→ you must store the lambda zip in S3

→ you must reflect the S3 zip location in the CloudFormation code.

→ S3 Bucket

→ S3 Key: full path to zip

→ S3 Object Version: if versioned bucket

→ if you update code in S3 but not ~~in~~ the other then CloudFormation won't update for it.

## → Lambda Container Image

→ Deploy Lambda function as container image upto 10 GB from ECR.

→ Pack complex dependencies, large dependencies in container

→ Base images are available for Python, Node.js, Java, etc.

→ can create your own image as long as it implements the Lambda Runtime API.

→ Test the container locally using the Lambda Runtime Interface Emulator

→ unified workflow to build apps

## → Best practices

- Use AWS-provided BaseImage
- Use Multi-stage Build
- Build from stable to frequently changing
- Use a Single Repository for function with Large layers
- Use them to upload large lambda function (10GB)

## → AWS vs Lambda versions

- \$LATEST
- pointer to lambda function version :- Aliases

## → Lambda & CodeDeploy

- CodeDeploy can help you automate traffic shift for lambda aliases.
- Features is integrated to SAM framework
- Linear :- grow traffic every N minutes until 100%.
  - linear 10 percent every 3 minutes.
  - " 10 " " 10 minutes"
- Canary :- try x percent than 100%.
  - canary canary 10 percent 5 minutes
  - canary 10 percent 30 minutes
- All at once - immediate
- Can create pre & post traffic hooks to check the health of the lambda function

## → Lambda - Function URL

- Dedicated HTTPS endpoint for your Lambda function
- A unique URL endpoint is generated for you (never changes)
- Invoke by web browser, Postman, curl, or HTTP client
- Access through public & Internet only
- Supports Resource-based policies & CORS configuration.
- Can be applied to any function Alias or \$Latest.
- Create and config using AWS console & CLI.
- Throttle using Reserved Concurrency.

## → Security :-

- Resource based policy
- CORS

- AuthType NONE - allow public and Unauthenticated access.
- AuthType AWS-IAM - IAM is used to authenticate and authorize request

## → Lambda and codefury profiling

- Gain insights into runtime performance of your Lambda function using CodeFury profiler.
- CodeFury creates a profiler Group for your Lambda function.
- Support Java & Python

→ AWS Lambda limit

→ Execution :-

→ Memory allocation - 128 MB - 10GB (1MB increment)

→ Maximum execution time - 900 second (15 minutes)

→ Environment Variable (4 KB)

→ Disk capacity in the function container (in /tmp)  
↳ 512 MB to 10GB

→ Concurrency executions :- 1000 (can be increased)

→ Deployment :-

→ Lambda function deployment size (.zip) → 50MB

→ Size of Uncompressed deployment (code+deps) : 250 MB

→ Can use the /tmp directory to load other files at startup.

→ Size of environment variable - 4KB

→ AWS Lambda Best practices

→ Perform heavy duty work outside of your function handler

→ Use environment variable

→ Minimize your development package size to its runtime necessities.

→ Avoid using recursive code, never have a lambda function call itself.

## \* AWS DynamoDB

- NOSQL databases scales horizontally
- Fully managed, highly available with replication across multi-AZs.
- No SQL database
- Scales to massive workloads, distributed database
- Millions of request per second, billions of row
- 100s of TB storage
- fast & consistence in performance
- integrated with IAM securities, authorization and administration.
- Enable event driven programming with DynamoDB streams
- low cost and auto-scaling capability
- Standard & Infrequent Access (IA) Table class
- DynamoDB Made up of tables
- Each tables have primary key
- can have infinite no of items
- Each item has attribute
- Max size of an item is 400KB
- Data type:-

Scalar type :- String, Number, Binary, Boolean, Null

Document type :- List, Map

Set type :- String set, Number set, Binary set

## Primary Key

→ Option 1 :- Partition Key (HASH)

→ Partition Key :- Key must be unique for each item

→ Partition Key must be diverse so data is distributed

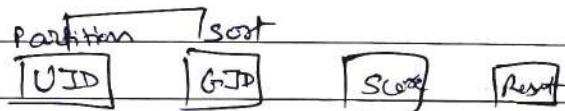
→ Ex:- "USER-ID" for user Table.

Option 2 :- Partition Key + sort key (Hash + Range)

→ The combination must be unique for each item

→ Data is grouped by Partition key

→ Ex:- User-game table "User-ID" for partition key and "Game-ID" for sort key.



→ DynamoDB - Read/write capacity Modes

→ Control how you manage your table capacity

→ Provisioned Mode (default)

↳ you specify the number of read/write per second

↳ you need to pay capacity beforehand

↳ pay for provisioned read & write capacity unit

→ ON demand Mode

→ Read/write automatically scale up/down with your workloads

→ No capacity planning needed

→ Pay for what you use, more expensive

## R/W Capacity Modes - provisioned

- Table must <sup>provisioned</sup> provide read and write capacity unit
- Read capacity unit (RCU)
- Write capacity unit (WCUs)
- option to setup auto-scaling
- Throughput can be exceeded temporarily using "Burst capacity"
- advice to do exponential backoff retry.

### ⇒ RCU

- one write capacity unit represent one write per second for an item upto 1 KB is size.
- If item size more it consume more WCU.

Ex-1 :- 10 item per second, with item size 2 KB

$$\rightarrow 10 \times \frac{2 \text{ KB}}{1 \text{ KB}} = 20 \text{ WCU}$$

Ex-2 :- 6 item per second, with item size 4.5 KB

$$6 \times \frac{4.5 \text{ KB}}{1 \text{ KB}} = 30 \text{ WCU}$$

Ex-3 :- 120 item per minutes, with item size 2 KB

$$= \frac{120}{60} \times \frac{2 \text{ KB}}{1 \text{ KB}} = 4 \text{ WCU}$$

→ Eventually Consistent Read

→ If you <sup>wrote</sup> <sub>read</sub> just after a write, its possible we'll get some stale data because of replication

→ Strong Consistent Read

- If you read just after write we will get correct data.
- get consistentRead parameter True
- Consume twice the RCU.

→ Read Capacity RUnits

- One read capacity unit represent "one" strongly consistent Read per Second and "two" Eventually Consistent Reads per Second for an item upto 4KB in size
- If the item are larger than 4KB more RCU are consumed.

Ex:- 10 Strongly Consistent Read per sec, with item Size of 4 KB

$$\rightarrow 10 \times \frac{4 \text{ KB}}{4 \text{ KB}} = 10 \text{ RCU}$$

Ex:- 16 Eventually Consistent Read per sec, with item Size of 12 KB

$$\rightarrow \left(\frac{16}{2}\right) \times \frac{12 \text{ KB}}{4 \text{ KB}} = 24 \text{ RCU}$$

Ex:- 10 Strongly consistent Read per sec, with item Size of 6 KB (6 → 8 KB)

$$\rightarrow 10 \times \frac{8 \text{ KB}}{4 \text{ KB}} = 20 \text{ RCU}$$

→ If we exceed provisioned RCUs and WCU's we get "provisionedThroughputExceededException"

→ Reason:-

→ HOT keys, HOT partition, very large item

→ Solutions!:-

→ Exponential backoff, Distribute partition keys  
→

→ ON Demand - R/W Capacity Mode

→ Read/write automatically scale up/down with your workloads

→ No capacity planning needed.

→ 2.5X more expensive

→ Dynamo DB - Writing Data

→ PutItem!:-

→ Create a new item or fully replace an old item (same Primary key)

→ Consume WCU's

→

→ UpdateItem!:-

→ Edits an existing item, attributes or add new item if it doesn't exist.

→ can be used to implement Atomic counters!:-

→ Conditional writes!:-

→ Accept write / update / delete only if conditions are met otherwise written error

→ Helps with concurrent access to item

→ No performance impact.

## → Reading Data

- GetItem :-
  - Read based on primary key
  - Primary key can be HASH or HASH+Range
  - Eventually consistent read
  - ProjectionExpression Can be specified to retrieve only certain attributes.
- Scan ! - the entire table and then filter out data
  - for fast we use parallel scan
  - consume a lot of RCU.
  - return upto 1MB data.
- Can use ProjectExpression & filter Expression.

## → Deleting Data

- DeleteItem ! -
  - Delete an individual item
  - Ability to perform a conditional delete
- DeleteTable ! -
  - Delete a whole table and all its items
  - Much quicker deletion than called DeleteItem on all items

## → Batch operation

- Allow you to save in latency by reducing the no of API calls
- Operation are done on parallel

### → BatchWriteItem :-

- upto 25 putItem and/or DeleteItem in one call
- upto 16 MB of data written, 400 KB per item
- can't update item
- UnprocessedItem for file write operation.

### → BatchGetItem :-

- Return item from one or more tables.
- upto 100 items, upto 16 MB of data
- Unprocessed Keys for failed operation

### → Partial

- Allow you to insert, select, update and delete data into a dynamoDB using SQL

## \* → DynamoDB - Local Secondary Index (LSI)

### → Alternative Sort Key for your table

- The sort key consist of one scalar attribute (String, NoP)
- upto 5 local secondary index per table
- Must be define at table creation
- Attribute projection! - can contain some or all the attributes of the base table

## → Global Secondary Index (GSI)

- Alternative primary key (Hash or Hash + Range key) from the base table.
- Spread up queries on non-key attributes.
- The index key consist of scalar attribute.
- Attribute projection :- same or all the attribute of the base table.
- Much provision RCVs & focus for the index.
- can be added / modified after table creation.

## → DynamoDB DA Accelerator (DAX)

- fully managed, highly available, seamless in-memory cache for DynamoDB.
- Microsecond latency for cache reads & queries.
- Don't required application logic modification.
- Solve Hot key issue.
- 5 minute TTL for cache.
- upto 10 nodes in cluster.
- Multi AZ.
- Secure.

## → DynamoDB Stream

- ordered stream of item-level modification in a table.
- Stream records can be
- sent to KDS
- send Read by Lambda
- Read by Kinesis Client library application
- Data retention 24 hours
- made up of shards.
- Records are not retroactively populated in a stream after enabling it.

→ DynamoDB - Time to live (TTL)

→ Automatically delete items after an expiry timestamp

→ DynamoDB CLI

→ Project-expression : one or more attributes to retrieve

→ --filter-expression : filter item before return to you.

→ AWS CLI Pagination option

— page-size

— max-items

— starting-token

→ DynamoDB Transaction

→ Coordinated, all-or-nothing operations to multiply item across one or more tables.

→ provides ACID

→ Two Operations :-

→ Transact GetItems

→ Transact WriteItems

## → DynamoDB Transactions - capacity computations

Ex:- 3 Transactional writes per second with item size 5 KB  
 $= 3 \times \frac{5 \text{ KB}}{1 \text{ KB}} \times 2 = 30 \text{ RCU}$

Ex:- 5 Transaction read per second with item size 5 KB  
 $= 5 \times \frac{8 \text{ KB}}{4 \text{ KB}} \times 2 = 20 \text{ RCU}$

## → DynamoDB write Sharding

→ Imagine we have a voting app with two candidate A and B

→ A suffix to partition key value

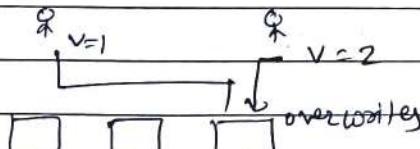
→ Two methods:-

↳ Sharding using Random suffix

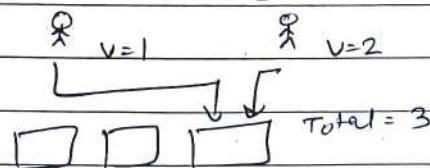
→ Sharding using calculated suffix

## → write Types

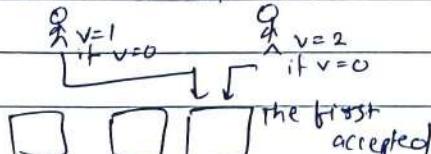
→ Concurrent writes



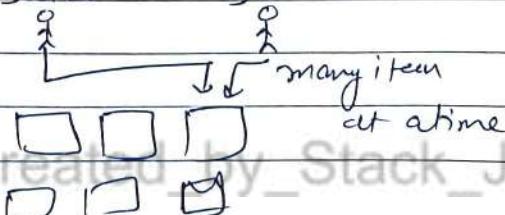
→ Atomic writes



→ Conditional writes



→ Batch writes



## \* AWS API Gateway

- AWS Lambda + API Gateway: no infrastructure to manage
- Support for the websocket protocol.
- Handle API versioning (v1, v2...)
- Handle different environments (dev, test, prod...)
- Handle security
- Create API keys, handle request throttling
- Swagger / Open API import to quickly define API
- Transform and validate request and response
- Generate SDK and API specifications
- Cache API response
  
- Lambda Function
  - ↳ Invoke Lambda function
  
- HTTP
  - ↳ Expose HTTP endpoint in the backend
  
- AWS Service :-
  - Expose any AWS API through the API Gateway.
  
- Endpoint Types
  - Edge-optimized :- for global client
  - Regional :- for client within same region
  - Private :- can only be used from your VPC

→ API Gateway - Security

→ User Authentication through

→ IAM Roles

→ Cognito

→ Custom Authorizer

→ Custom Domain Name HTTPS

→ API Gateway - Deployment stages

=====

→ Changes are deployed to "stages"

→ You the name you like for stages (dev, test, prod)

→ Each stage have its own parameters

→ Stages can be rolled back as a history of deployment is kept.

→ Stage Variable

=====

→ Stage Variable are like environment variable for API Gateway.

→ Use them to change often changing configuration value

→ Used in Lambda function ARN, HTTP endpoint

Use Case

→ Configure HTTP endpoint your stages talk to  
(dev, test, prod)

→ Pass configuration parameter to AWS Lambda

→ Pass to context object in Lambda

→ Format: \${stageVariable.variableName}

## → API Gateway - Canary Deployment

- possibility to enable Canary deployments for any stage. (usually prod)
- Choose % of traffic the canary channel receive
- Metrics and logs are separated.
- Possibility to override stage variable for canary
- This is blue/green development with AWS Lambda & API Gateway.

## → API Gateway Integration Type

- Integration Type Mock
  - API Gateway return to the response without sending the request to the backend.
- Integration Type HTTP / AWS (Lambda / AWS services)
  - You must configure both the integration request and integration response.
  - setup data mapping using mapping templates for the request & response.
- AWS\_PROXY (Lambda proxy)
  - incoming request from the client is the input to lambda.
  - The function is responsible for the logic of request/response.

## → HTTP-PROXY

- No Mapping template
- The HTTP request is passed to the backend.
- The HTTP response from the backend is forwarded to by API Gateway.

## → Mapping Templates

→ Mapping template can be used to modify request/response.

- Rename/modify query query string parameter.
- Modify body content, add header.
- content type set to application/json or application/xml.

## → JSON to XML with SOAP

→ SOAP API are XML based.

→ Extract data from the request! either path, payload or header

→ Build SOAP Message based on request data.

→ Call SOAP Service and receive XML response.

→ Transform XML response to desired format and respond to the user.

## → API Gateway - Open API Spec

- common way of defining REST APIs using API definition as code
- Can export current API as openAPI spec
- Can be written in YAML & or JSON
- we can generate SDK for your application

## → Caching API responses

- Caching reduces the no of calls to backend
- TTL is 300 second
- Caches are defined per stage
- possible to override cache setting per method.
- Cache encryption method
- Capacity → 0.5 GB to 237 GB
- expensive, make sense in prod, not in dev, test.

## → API Gateway - logging & tracing

- Cloudwatch Logs
- log contains info about request / response body
- Enable Cloudwatch logs at the stage level
- can override settings on a per API basis

### X-Ray

- Enable tracing to get extra info about requests in API Gateway
- X-Ray + AWS Lambda give you the full pic.

## → Cloudwatch Metrics

- Metrics are by stage, possibility to enable details metrics;
- CacheHitCount & CacheMissCount: efficiency of the cache.
- Count, IntegrationLatency, latency
- 4XXError & 5XXError  
(client side)      (server side)

## → Throttling

### → Account Limit :-

→ API Gateway throttle request at 1000 rps

→ In case of throttling:- 429 Too many requests

→ Can set Stage limit & Method limit to improve the performance.

→ You can define Usage plan.

## → CORS

→ CORS must be enable when you receive API calls from another domain

→ Option pre-flight request contain in header

## → API Gateway - Security

→ IAM permissions

→ Create an IAM policy authorization and attach to User/role.

→ Authentication = IAM

→ Authorization = IAM Policy

→ Leverage "Sig v4" capacity capability where IAM credential are in header.

## → Resource policies

→ Allow for cross Account Access

→ Allow for a specific source of IP Address

→ Allow for the VPC Endpoint.

## → Cognito user pool

- Cognito fully manage user life cycle, token expiring automatically
- API Gateway verifies identity automatically from AWS Cognito.
- Authentication :- Cognito user pools  
Authorization :- API Gateway method.

## → Lambda Authorizer

- Token based authorizer (bearer token)
- A request parameter based Lambda authorizer
- Lambda must return an IAM policy for the user, result policy is cached.  
Authentication - External  
Authorization - Lambda function

## → WebSocket API

- Two way interactive communication between a user browser and a server
- Server can push information to the client
- → This enable stateful application like course
- real time application such that chat app, collaboration platform, multiplayer games and trading platform.

## \* AWS CI/CD :-

AWS CodeCommit → storing our code

AWS CodePipeline → automating our pipeline from code to <sup>Elastic</sup> Bean

AWS CodeBuild → building & testing our code

AWS CodeDeploy → deploy code to EC2 Instance

AWS CodeStar → manage S/W development Activity in one place

AWS CodeArtifact → store, publish

AWS CodeBuild → automated code review using ML

### → CI

- Developer push the code to a code repository often
- A testing / build server checks the code as soon as its pushed
- The developer get feedback about the tests and checks that have passed/failed.
- find bugs early, then fix bugs.
- delivery faster as the code is tested / deployed often

### → CD

- Ensure that the S/W can be released reliably whenever needed.
- Ensure deployments happen often and are quick.
- Shift away from "one release every 3 months to 5 new releases a day."

## → AWS Code Commit

- version control is ability to understand the various changes that happen to the code over time
- Version Control - Git
- A Git repository can be synchronized on your computer, but it is online.
- Benefits
  - collaborative with other developer
  - Make sure the code is backed up somewhere
- Private Git Repo
- No size limit
- fully managed and highly available
- Integrated with Jenkins, AWS CodeBuild, and other CI tools.

## → AWS Code pipeline

- Visual workflow to orchestrate your CI/CD
- Source - Codecommit, ECR, S3, Bitbucket, GitHub
- Build - Codebuild, Jenkins, CloudBees, Teamcity
- Test - codebuild, AWS Device Farm
- Deploy - code deploy, Elastic Beanstalk, Cloudform, ECR, S3
- Invoke - Lambda, Step function
- Stage
- Build → Test → Deploy → Load testing →

- Each pipeline stage can create artifacts
- Artifacts stored in an s3 bucket and passed on to the next stage
- Code convert convert in artifact
- > Code Build

- Source → Code commit, S3, GitHub
- Build instructions :- code file buildspec.yml or <sup>insert</sup> manually in console.
- output log can be stored in S3 & CloudWatch Log.
- CloudWatch metrics to monitor
- EventBridge to detect failed builds and trigger notifications
- Build project can be defined within CodePipeline or CodeBuild.
- buildspec.yml file must be the root of your code
- env → define environment variable
- phases → specify command to run
- artifacts → what to upload to S3
- cache → files to cache to S3

### Code Deploy

- Deployment Service that automates application deployment.
- Deploy new application version to EC2 instance, on-premises server, Lambda function, ECS services.
- Automated Rollback capability in case of failed deployment.
- Gradual development control.
- A file named appspec.yml defines how the deployment happens.

- can deploy to EC2 or on-premises servers
- perform in-place deployment or blue/green-deployment
- development speed! -
- All at once, Half AT A Time, Once AT A Time, Custom
- AWS CodeStar
- An integration solution that gops : GitHub, codecommit, CodeBuild, codedeploy, CloudFormation, codepipeline, Cloudwatch.
- Quickly create CI/CD-ready project for EC2, Lambda, elastic Beanstalk
- Issue tracking integration with JIRA / GitHub issue
- Ability to integrate with Cloud9 to obtain a web IDE
- One dashboard to view all the components
- Free services, pay only for the underline usage of other services.
- limited customization

### AWS Artifact

- SW packages depend on each other to be built and new ones are created.
- Storing & referencing these dependencies is called artifact management.
- Traditionally, you need to setup your own artifact management system.
- Work with common dependencies management tool such as Maven, Gradle, npm, yarn, twine, pip and NuGet.

## → Amazon Code Guru

- An ML Powered Service for automated code review and application performance recommendation
- provide two functionalities
- Code Reviewer :- automatic code review for static code analysis
- Code Guru profiler :- visibility & recommendation about application performance during runtime.
- Max Stack Depth.
- Memory Usage limit Percent
- Minimum Time for Reporting In Milliseconds → sending reports
- Reporting Interval In Milliseconds → report profile
- Sampling Interval In Milliseconds → sample profile

## → AWS Cloud 9

- Cloud based Integrated Development Environment
- Code editor, debugger, terminal in a browser
- works on your project from anywhere
- prepacked with essential tools for popular programming languages.
- Share your development env with your team.
- fully integrate with AWS SAM & Lambda to easily build serverless application.

## \* AWS SAM

- SAM = Serverless Application Model
- framework for deploying serverless application.
- All the configuration in YAML code.
- Generate complex CloudFormation from simple SAM YAML file.
- Support anything from CloudFormation: output, mapping, parameter and resource.
- only two commands to deploy to AWS.
- SAM use CodeDeploy for Lambda function
- can help to run Lambda, API Gateway, DynamoDB

### SAM - CLI Debugging

- Locally build, test, and debug your serverless application that are defined using AWS SAM templates.
- provide Lambda like execution env locally
- SAM CLI + Toolkits ⇒ step through and debug your code.
- Support IDE :- Cloud9, VSCode, Jetbrains
- AWS Toolkits :- IDE plugins which allow you to build, test, debug, deploy and invoke Lambda function built using AWS SAM

### SAM Policy Template

- list of template to apply permission to your Lambda
- IMP example :-
- S3ReadPolicy :- Give read only permissions to object in S3
- SQSPollerPolicy :- Allow to poll an SQS queue
- DynamoDBCrudPolicy :- CRUD = Create, read, update, delete

## → SAM cmd codeDeploy

- SAM framework natively uses codeDeploy to update lambda function.
- Traffic shifting features
- pre-post traffic hooks features to validate deployment
- Easy and automated rollbacks using CloudWatch

## → AutopublishAlias :-

- Detect when new code is being deployed.
- Create and ~~up~~ publishes an update version of that function with the latest code.
- points the alias to the updated version of the lambda function

## → Deployment Preference

- Canary, linear, AllAtOnce

## → Alarm

- Alarm that can trigger a rollback

## → Hooks

- pre and post traffic shifting lambda function to test your deployment.

## → SAM - Local Capabilities

- Local start AWS Lambda
- Start a local endpoint that emulates AWS Lambda
- can run automated test against this local endpoint

- locally invoke lambda function
- Sam local invoke
- Invoke lambda function with payload once and quit after invocation complete.
- Helpful for generating testcase.
- If function call API to AWS then use --profile option.
  
- locally start an API Gateway endpoint
  - Sam local start-api
  - Start a local HTTP server that host all your f
  - automatic reloaded
  
- Generate AWS event for lambda function
  - Sam local generate-event
  - Generate local sample payloads for event types

### Summary

- SAM is built on Cloud formation
- SAM required the Transform and Resources Section.
- Command to know:-
  - Sam build : fetch dependencies and create local deployment artifacts
  - Sam package : Package and upload to S3, generate CF template
  - Sam deploy : Deploy to Cloud formation
  
- SAM policy templates for easy IAM policy definition
- SAM is integrated with CodeDeploy to do deploy to lambda aliases.

## \* AWS Cloud Development Kit

- Define your cloud infrastructure using similar language - Javascript / typescript, python, java and .NET
- Contains high level components called constructs
- The code is compiled into a CloudFormation templates (json / YAML)
- you can therefore deploy infrastructure and application runtime code together.
  - Great for lambda function
  - " " Docker containers in Fargate, EKS.
- CDK + SAM
  - You can use SAM CLI to locally test your code apps.
  - You must first run CDK Synth.

## CDK Constructs

- CDK Construct is a component that encapsulates everything CDK need to create the final CloudFormation Stack.
- Can represent a single AWS resource or multiple related resources.
- AWS construct library
  - A collection of constructs included in AWS CDK which contain constructs for every AWS resources.
  - Contains 3 different levels of construct available (L1, L2, L3)
- Construct Hub - contains additional constructs from AWS, 3rd parties, and open-source CDK community.

- Layer 1 → CFN Resources which represent all resources directly available in CF.
- Layer 2 → higher level AWS resources but with a higher level.
- Layer 3 → can called patterns which represent multiple related resources.

### → CDK Bootstrapping

- The process of provisioning resources for CDK before you can deploy CDK apps into an AWS environment.
- AWS environment = account & region
- Cloudformation stacks called CDKToolKit is created
  - S3 Bucket
  - IAM Role

### → Testing

- CDK Assertions Module
- Two types!
  - Fine-grained Assertions
  - Snapshot Tests
- Template.fromStack(myStack) → inside
- Template.fromString(myString) → outside.

## \* Amazon Cognito

→ Give user an identity to interact with our web or mobile application.

→ Cognito user pool :-

→ Sign in functionality for app user.

→ Integrate with API Gateway & Application load balancer.

→

→ Cognito Identity pool (federated identity).

→ Provide AWS credential to user so they can access AWS resources directly.

→ Integrate with Cognito User pools as an identity provider.

→ Cognito vs IAM :-

→ "hundred of user", "mobile user", "authenticate with SAML"

→ Cognito User pools (CUP) - User features

→ Create a serverless database of user for your web and mobile apps.

→ Simple login : Username & Password

→ Password reset, phone no & Email verification, MFA

→ Federated Identities

→ Feature : block user if their credential are compromised elsewhere.

→ login send back to json web token.

→ Integrate with API Gateway & Application load balancer.

## → Cognito Identity Pool

- Get identities for "users" so they obtain temporary AWS credentials.
- your identity pool include:
  - public provider, user of Cognito, OpenID Connect provider
  - Cognito Identity pool allow for unauthenticated access.
  - Then user can access AWS services directly through API Gateway.

## \* AWS Step Functions

- Model your workflow as state machines (one per workflow)
  - order fulfillment, data processing
  - web application, Any workflow
- Visualization visualization of workflow, execution of workflow, as well as history
- start workflow with SDK call, API Gateway, Eventbridge
- Error handling
  - Any state can encounter runtime error for various reasons:
    - State machine definition issue, Task failure, Transient issue
  - use Retry and catch
  - predefined error code: states.ALL, states.Timeout, state.TaskFailed, state.permissions

- ErrorEquals
- IntervalSeconds
- Backoff Rate
- MaxAttempts

→ max attempts then enter ~~catch~~ catch

- Error Equals

- Next

→ Resultpath define the input for next state/field

### → Activity Task

→ Enable you to have the task work performed by an Activity worker.

→ Activity worker apps can be running on EC2, Lambda, mobile device.

→ Activity worker pool for a task using GetActivityTask API.

→ After Activity worker complete its work, it send a response of its success/failure using SendTaskSuccess or SendTaskFailure

→ wait time → TimeoutSecond

~~heartbeat from activity worker~~ → SendTaskHeartBeat HeartBeatSeconds

→ define time using TimeoutSecond and Max time → 1 year.

### → AWS AppSync.

- AppSync is managed service that uses GraphQL
- GraphQL make it easy for application to get exactly the data they need.
- Combination of data.
- Start with one GraphQL Schema

- Amplify studio
- Amplify CLI
- Amplify Libraries
- Amplify Hosting

## \* AWS STS - Security Token Service

- Allow to grant limited and temporary access to AWS resources.
- AssumeRole :- Assume role within your account or cross account.
- AssumeRoleWithSAML :- Credential of user logged in with SAML.
- AssumeRoleWithWebIdentity :- Get token for user logged in with an IPP.
- GetSessionToken :- for MFA.
- GetFederationToken :- obtain temp cred for federated user.
- GetCallerIdentity :- return details about the IAM User or Role used in API calls.
- DecodeAuthorizationMessage :- decode error message when AWS API denied.

## → AWS Directory Services

### → AWS managed Microsoft AD

- ↳ Create your own AD in AWS, manage users locally, support MFA.



### → AD Connector

- Directory gateway to redirect to an on-premise AD, support MFA.
- User are managed on the on-premise AD

→ Simple AD

→ AD-compatible managed directory on AWS

→ Cannot be joined with on-premise AD

## \* AWS Security and Encryption

### → AWS KMS

→ AWS manage encryption key for us

→ Fully integrated with IAM for authorization

→ Easy way to control access to your data.

→ Able to audit KMS key usage using CloudTrail

→ Seamlessly integrated into most AWS service

→ Never ever store your secrets in plaintext, especially in your code.

→ also available through API calls.

→ Encryption <sup>secret</sup> key can be stored in code / environment

### → KMS Keys Types

#### → Symmetric (AES-256 Keys)

→ Single encryption key that use for Encrypt & decrypt

→ AWS service that are integrated with KMS use symmetric <sup>KMS</sup> keys

→ You never get access to the KMS keys unencrypted CMKs  
(through API only)

#### → Asymmetric (RSA & ECC Key pairs)

→ Public (Encrypt) and Private key (Decrypt) pair

→ Used for Encrypt / Decrypt or sign / verify operation

→ The public key is downloadable, but you can't access the private key unencrypted.

→ encryption outside AWS

## → Types of KMS Key

- AWS owned keys (free)
  - AWS managed key (free) (eg:- aws/service-name)
  - Customer managed keys created in KMS: \$1/month
  - Customer managed keys imported
  - + pay for API calls to KMS
- 
- Automatic key rotation
  - AWS managed KMS key : auto every 1 year
  - Customer-managed KMS Keys = 11
  - Imported KMS Keys : manual rotation

## → Envelope Encryption

- KMS Encryption API calls has limit of 4KB.
- If we need to encrypt > 4KB we need to use Envelope Encryption.
- The main API helps is GenerateDataKey API

for exam :- anything over 4KB of data that needs to be encrypted must use the Envelope Encryption  
= GenerateDataKey API

- The AWS Encryption SDK implemented envelope encryption for us!
- Data Key Caching
  - Reuse keys instead of creating new all the time for encryption
- Use Local Cryptomaterials cache

→ S3 Bucket key for SSE-KMS encryption.

→ New setting to decrease

→ No of API call made to KMS from S3 by 99%.

→ Cost of overall overall KMS encryption with Amazon S3 by 99%.

→ A S3 Bucket key is generated

→ That key is used to encrypt KMS object with new data keys

→ You will see less KMS CloudTrail events in CloudTrail

## Cloud HSM

→ AWS manage the SW for encryption.

→ CloudHSM → AWS provision encryption hardware

→ HSM is tamper resistant

→ Support both type of encryption

→ Good for SSE-C encryption

## SSM Parameter Store

→ Secure storage for Configuration and secrets

→ Optional seamless encryption using KMS

→ Serverless, Scalable, durable, easy SDK

→ Version tracking of configuration / secrets

→ Security through IAM

→ Notification with Amazon EventBridge

→ Integration with CloudFormation

## \* AWS Secrets Manager

- Newer service, meant for storing secrets
- Capability to force rotation of secrets every x days
- Automatic generation of secrets or rotation
- Integration with Amazon RDS and many
- Secrets are encrypted using KMS
- Replicate secrets among multiple AWS Regions
- Make read replicas sync with primary secrets

## \* AWS SES Simple Email Service

- send email to people using:-
- SMTP interface → or AWS SDK
- Ability to receive msg <sup>email.</sup> integrates with S3, → SNS → Lambda
- Integrated with IAM for allowing to send emails

## \* Amazon OpenSearch Service

- <sup>OpenSearch</sup> Amazon Elasticsearch is successor to Amazon Elastic Search
- In DynamoDB, queries only exist by primary key or indexes/indexes. With OpenSearch, you can search any field, even partially matches.
- Two modes:- managed cluster or serverless cluster
- Does not natively support SQL.
- Ingestion from Kinesis Data Firehose, AWS IoT, and CloudWatch logs.

- Security through Cognito & IAM, KMS encryption, TLS
- Comes with open search Dashboards

## \* Amazon Athena

- serverless query service <sup>to</sup> that analyze data stored in S3.
- Uses standard SQL query language to query the files. (built on Presto)
- Support CSV, JSON, ORC, Parquet
- Commonly used with Amazon Quicksight for reporting / dashboards.
- Use case :- business Intelligence, Analytics, reporting, analyze & query VPC flow logs, ELB logs.
- Exam Tip :- analyze data in S3 using serverless SQL, use Athena.
- Use columnar data for cost saving
- Compress data
- partition datasets
- large file > 128 KB
- use datasource connector that run on AWS Lambda to run federated queries

## \* Amazon Managed streaming for Apache Kafka (Amazon MSK)

- Alternative to Amazon Kinesis
- Fully managed Apache Kafka on AWS
- Allow to create, update, delete cluster
- MSK create and manage Kafka brokers nodes
- Deploy the MSK Cluster in your VPC, multi AZ
- Automatic recover from failure.
- Data is stored on EBS volume for long as you want.

### MSK Serverless

- Run Apache Kafka on MSK without managing the capacity.
- MSK automatically provisions resources and scale Compute & storage.

## \* AWS Certificate Manager (ACM)

- Let's you easily provision, manage, and deploy SSL/TLS certificates.
- Used to provide in-flight encryption for website.
- Support both public and private TLS certificates.
- Free of charge for public TLS certificates.
- Automatic TLS certificate renewal.
- Integration with
  - ELB, CF Distribution, API on API Gateway

## → AWS private Certificate Authority (CA)

- Manage service allows you to create private certificate authority (CA) including root and subordinates CA.
- can issue and deploy end-entity X.509 certificates
- certificates are trusted by only your organization
- works for AWS services which are integrated with ACM.
- use case: Encrypted TLS communication, Cryptographically signing code, Authenticate user, Compute, API, IoT.

## \* Amazon Macie

- Amazon Macie is fully managed data security and data privacy service that uses ML and pattern matching to discover and protect your sensitive data in AWS.
- Macie help identify and alert you to sensitive data such as personally identifiable information (PII).

## \* AWS AppConfig

- Configure, validate and deploy dynamic configuration.
- Deploy dynamic configuration changes to your application independently of any code deployment.
- Feature flag, application tuning, allow block listing use in EC2, Lambda, ECS, Fargate.
- Gradually deploy the config changes and roll back if issue occurs.
- Validate configuration changes before deployment :-
  - JSON Schema
  - Lambda function

## \* → CloudWatch Evidently

- Safely validate new features by serving them to a specified % of your user.
- Launches (=feature flag) : enable and disable features for a subset of users.
- Experiments (=A/B testing) : compare multiple versions of the same features.
- Overrides ! - pre-define a validation for a specific user.
- Store evaluation events in CloudWatch Logs or S3.