

IAM (Identity and access management)

- Pay-as-you-go-basis
- creating users and allowing them to access particular servers
- Third party application slack, to prevent eavesdropping

1. Principal
2. Authentication
3. Request (API-call behind)
4. Authorization
5. Actions
6. Resources



User

- can be a person or an application
- each IAM user is attached with only one AWS account

Groups

- collection of IAM users
- permission applied to group, applied to all users in the group

Policies

- Policy sets permissions and controls access
- stored as JSON documents

There are 2 types

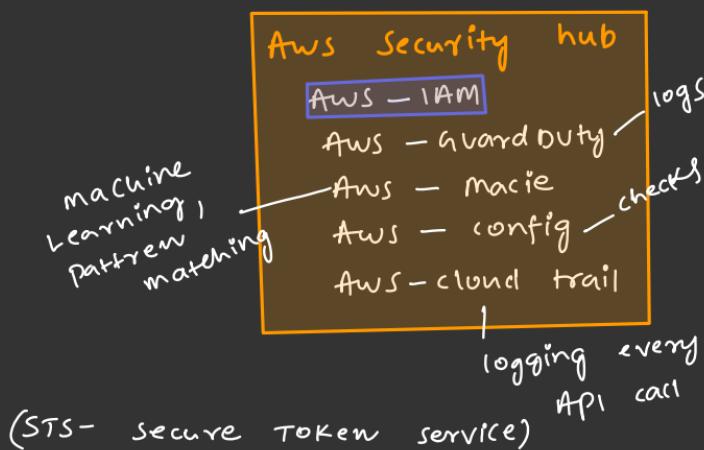
managed policy - multiple entities
inline policy - single entity

Roles

- set of permissions that define what actions are allowed
- role permissions are temporary credentials
 - use resources but don't want to save the password or credential

Features

- Shared access to AWS account
- granular permissions (least privileged action)
- MFA (Multi-factor Authentication)
- Identity Federation (Facebook, Google)
- Free to create, use - no charge
- PCI DSS compliance
- password policy (password changes etc--)



Benefits

- control over data
 - fine grain identity
 - continuous monitoring for real time security information
 - reduce human config errors
 - wide variety of deeply integrated solutions
 - regularly achieves third-party validation for thousands of global compliance
 - help you meet security and compliance standards
 - most secure global infrastructure
 - ability to encrypt data move it and manage retention
- STS generates **git credentials** for IAM users (Temporary credentials - implies)
- AWS partner to carry out audit on AWS account - **cross account IAM role**
- Assume role with SAML - AWS STS service

Billing

- pay-as-you-go
- There will be a forecast of budget graph based on our past resource consumption.
- **cost explorer** - analyze the cost and usage
 - set time interval and granularity
 - filter/group the data
 - forecast future costs and usage
 - save the progress

→ **Budgets** - track cost usage from simplest one to most complex use cases

→ stay informed with alerts and reports

→ customize budget as per needs

→ granular budget time periods

There are 4 types here

cost budget - receive alerts when the user-defined thresholds are met

usage budget - monitor one or more specified usage types and alert when the threshold is met (for each service)

savings plan budget - track the utilization and receive alerts when the percentage drops below the threshold

Reservation budget - track the utilization with reservation and alerts when the percentage drops below the threshold

alerts - email recipient
Amazon sns alert (email, sms --)
Amazon chatbot alerts (chat channels like slack and chime)

cost and usage reports (CUR)

→ accurate budget plans (precise)

↳ product code, usage type, operation

Benefits

→ Access comprehensive AWS and usage info

- Track your reserved instance (RI) and savings plan usage
- Leverage integrations with other analytic services
- public and private pricing
- cost category and cost allocation tags
- we need to create or add existing S3
bucket. It will send reports to that bucket
- Parquet - easy to extract data using SQL

Pricing calculator

- view bills, and pay invoices in preferred currency
- configure cost estimate that fits business or personal needs
- we have to choose the resource type
- By giving specifications and resource types we can get the estimated budget

Key management service (KMS)

- encryption for both in-transit and at-rest.
- two methods for encryption
 - client-side Encryption
 - server-side Encryption
- we encrypt our own data and manage the keys
- AWS encrypts and manages keys for you (S3, TBS, Redshift provide encryption using Kms behind the scenes)

- Service that manages Encryption keys
- It uses Hardware security module (HSM) to store only CMKs
- uses AWS CloudTrail to log usage logs for audits and compliance needs

Customer master keys (CMK)

- A CMK never leave the region
- can only encrypt max of 4KB
(if 4KB↑ we use data keys)
- we will ask the CMK to generate data key
- we call KMS API for data keys (2 will be sent, plaintext data key, encrypted data key)
- After encryption we will delete plaintext data key (if present easy to decrypt)
- we can use OpenSSL or AWS Encryption SDK (outside AWS)
- To decrypt I take encrypted data key and ask KMS API for decryption algorithm and get the decrypted text

Hardware security module (HSM) - within a VPC

- computing device that processes cryptographic operations and provides secure storage for cryptographic keys
- enables to generate and use own encryption keys (compliant with FIPS 140-2 level)
- automates administrative tasks like hardware provisioning, software-patching, high availability and back-ups
- Cloud HSM can scale quickly on demand with

no up-front costs

- AWS Cloud HSM interacts VPC instance with SSL
- in clusters
- each cluster is logical HSM
- To interact HSM cluster, we need AWS cloud HSM client software
- Operations on one HSM, other HSM's in the same cluster gets updated
 - Offload the SSL/TLS processing for web servers
 - protect the private keys for an issuing certificate authority
 - enable transparent data encryption for oracle database

Benefits

- easy for load balance and scale
- pay by the hour

Simple Storage Service (S3)

- infinitely scaling storage
- Object storage

use cases

- Backup and storage
- Disaster recovery
- Archive
- Hybrid cloud storage
- Application hosting

- media hosting
- data lakes and big data analytics
- software delivery
- static website

Buckets

- files are stored in buckets
- defined at region-level not global → name
- no uppercase, no underscore.

Key \leftarrow s3://my-bucket/my-folder | file.txt
Prefix Object

S3 - Security

- user-based (IAM policies)
- resource-based
 - bucket policies in bucket itself
 - object access control list (ACL)
 - bucket access control list (ACL)
- encrypt object in S3 using encryption keys

S3 - Bucket policy

- JSON based policies
- S3 bucket for policy
- user access to S3 - IAM permission
- ec2 instance access - IAM roles
- cross account access
- public access
- To make the objects public we need to uncheck block access and add bucket policy to objects to make them public

S3 - static website hosting

- we need to make the bucket public (for no issues)
- enable static website hosting
- if 403 error comes, (policy allow: public reads)

S3 - versioning

- enabled at bucket level
- it is best practice to version your buckets
 - protect against unintended deletes
 - easy roll back to previous version
- no versioning will be represented as null
- suspending version doesn't delete previous version
- same name files can be uploaded to update the web page
- when we accidentally delete a object, we can restore it by deleting the object that is of type delete marker

S3 - Replication (Asynchronous replication)

CRR - cross region replication

SRR - same region replication

CRR - compliance, lower latency access, replication across accounts

SRR - log aggregation, live replication between production and test accounts

- After you enable Replication, only new objects are replicated
- For existing objects we can perform batch operations
- can replicate delete markers from source to target
- Deletions with a version ID are not replicated
- NO chaining of replication
- To enable replication, versioning should be applied on both the buckets

S3 - storage classes

- Standard - General purpose
- Standard - Infrequent access
- One zone - infrequent access
- glacier instant retrieval
- glacier flexible retrival
- glacier deep archive
- Intelligent tiering
- can move between classes manually or using S3 life cycle configurations

Durability or Availability

S3 standard - general purpose

- frequently accessed data
- low latency high throughput
- sustain 2 concurrent facility failures

use cases big data analytics, mobile & gaming applications, content distribution

S3 - Infrequent access

- less frequently accessed, require rapid access
- less cost than above
- use cases: disaster recovery, backups

one-zone infrequent access

- high durability in single AZ, data lost when AZ is destroyed
- use cases: storing secondary backup copies of on-premises data, or data to recreate

Glacier

- low-cost object storage meant for archiving/backup
- price for storage + object retrieval cost

Instant retrieval

- millisecond retrieval (90 days in a year)
- low latency retrieval

flexible retrieval

- expedited (1-5 mins) (90 days)
- standard (3-5 hrs)
- Bulk (5-12 hrs) - free

Deep archive - long term storage - **cheapest**

- Standard (12 hrs), Bulk (48 hrs) (150 days)

Intelligent tiering

- automatically move between classes
- have to pay for monthly monitoring and auto-tiering
- - IA tier - object not accessed for 30 days
 - Archive instant - object not accessed for 90 days
 - Archive Access - (config) 90 - 700+ days
 - Deep Archive - (config) 180 - 700+ days
- alternative to Intelligent tiering is manual life cycles (free cost)

	Standard	Intelligent-Tiering	Standard-IA	One Zone-IA	Glacier Instant Retrieval	Glacier Flexible Retrieval	Glacier Deep Archive
Storage Cost (per GB per month)	\$0.023	\$0.0025 - \$0.023	\$0.0125	\$0.01	\$0.004	\$0.0036	\$0.00099
Retrieval Cost (per 1000 request)	GET: \$0.0004 POST: \$0.005	GET: \$0.0004 POST: \$0.005	GET: \$0.001 POST: \$0.01	GET: \$0.001 POST: \$0.01	GET: \$0.01 POST: \$0.02	GET: \$0.0004 POST: \$0.03 Expedited: \$10 Standard: \$0.05 Bulk: free	GET: \$0.0004 POST: \$0.05 Standard: \$0.10 Bulk: \$0.025
Retrieval Time	Instantaneous					Expedited (1 – 5 mins) Standard (3 – 5 hours) Bulk (5 – 12 hours)	Standard (12 hours) Bulk (48 hours)
Monitoring Cost (pet 1000 objects)	\$0.0025						

	Standard	Intelligent-Tiering	Standard-IA	One Zone-IA	Glacier Instant Retrieval	Glacier Flexible Retrieval	Glacier Deep Archive
Durability	99.999999999 == (11 9's)						
Availability	99.99%	99.9%	99.9%	99.5%	99.9%	99.99%	99.99%
Availability SLA	99.9%	99%	99%	99%	99%	99.9%	99.9%
Availability Zones	>= 3	>= 3	>= 3	1	>= 3	>= 3	>= 3
Min. Storage Duration Charge	None	None	30 Days	30 Days	90 Days	90 Days	180 Days
Min. Billable Object Size	None	None	128 KB	128 KB	128 KB	40 KB	40 KB
Retrieval Fee	None	None	Per GB retrieved	Per GB retrieved	Per GB retrieved	Per GB retrieved	Per GB retrieved

S3 - life cycle rules

Transition actions: config obj to transition to another storage class

Expiration actions: config objects to expire after some time

→ can be used to delete old versions of files

→ can be used to delete incomplete multi-part uploads

S3 - Analytics

- helps you decide when to transition objects to the right storage class
- does not work for one zone-IA or Glacier
- Report updated daily
- 24 to 28 hrs to start seeing analytics

S3 - Requester pays

- if a person (requester) has to pay networking / transfer costs then Requester pays bucket
- helpful to share large datasets
- the requester must be authenticated in AWS

S3 event notifications (SNS, SQS, TFn)

use case: generate thumbnails of images uploaded to S3

S3 → put object → TFn → thumbnail

→ can be created as desired (sec/min ↑)

event Bridge (18 AWS Services)

- multiple destinations
- event bridge capabilities - Archive, Replay events, Reliable delivery

S3 performance

- multi-part upload ($>100\text{MB}$ $>5\text{GB}$)
- S3 Transfer Acceleration
 - Increase transfer speed by transferring file to AWS edge location which will forward data to S3 bucket
 - compatible with multi-part upload
- Byte-Range fetches (partial access of file)

S3 Select & glacier-select

- Retrieve less data using SQL by performing server-side filtering
- less network transfer, less CPU cost client-side
- same thing on many obj's use batch operation
- can use S3 inventory to get object list

Elastic compute cloud (EC2)

- infrastructure as a service
- Renting virtual machines (EC2)

- storing data on virtual drivers (EBS)
- distributing load across machines (ELB)
- scaling services using auto-scaling (ASG)

config options

- OS → CPU → RAM → EBS / EFS
- EC2 Store → Network card → Firewall rules
- Bootstrap script

User data - EC2

- Bootstrap instances using EC2 user data script
- Bootstrapping: launching commands when machine starts
- EC2 user data is used to automate
 - installing updates
 - installing softwares
 - downloading files from internet
- EC2 data script runs with root user
- public IPV4 will change frequently

EC2 - instance types

- General purpose → Storage optimized
- Compute optimized → Instance features
- Memory optimized → Measuring instance performance
- Accelerated computing
- General purpose
 - great for diversity of workloads such as web servers or code repositories

- balance between compute, memory and networking
- compute optimized
 - good for computing tasks that require high performance
 - Batch processing workloads
 - media transcoding
 - high performance web servers
 - high performance computing (HPC)
 - scientific modeling
 - machine learning
 - dedicated gaming servers
- Memory optimized (volatile)
 - fast performances for workloads that process large data sets in memory
- use cases: High performance, relational / non-relation databases
 - distributed web scale cache stores
 - in memory databases optimized for BI
 - Applications performing real-time processing of big unstructured data
- storage optimized (storing, non-volatile)
 - great for storage-intensive tasks that require high, sequential read and write access to large data sets on local storage

- use cases: High frequency OLTP systems
- Relational & NoSQL databases
 - Cache for in-memory databases (Redis)
 - Data warehousing applications
 - Distributed file systems

Instance	vCPU	Mem (GiB)	Storage	Network Performance	EBS Bandwidth (Mbps)
t2.micro	1	1	EBS-Only	Low to Moderate	
t2.xlarge	4	16	EBS-Only	Moderate	
c5d.4xlarge	16	32	1 x 400 NVMe SSD	Up to 10 Gbps	4,750
r5.16xlarge	64	512	EBS Only	20 Gbps	13,600
m5.8xlarge	32	128	EBS Only	10 Gbps	6,800

EC2 - Security groups

- It controls how traffic is allowed **into or out of** our EC2 instances
- They only contain rules (rules can be ref by IP or by security group)
- **Firewall**
- **Access to ports**
- **Authorised IP ranges**
- one security group can be attached to multiple instances
- Regional service
- If application gives "connection refused" error, it's an application error or it's not launched

- All inbound traffic is blocked by default
- All outbound traffic is authorised by default

SSH - 22 — log into a linux instance

FTP - 21 — file share

SFTP - 22 — using SSH

HTTP - 80 — access unsecured websites

HTTPS - 443 — access secured websites

RDP - 3389 — log into windows instance

- windows < 10 SSH X (SSH using putty only)

Ec2 Instance connect

- within the browser (no key pairs)
- works only with Amazon-Linux-2 AMI
- make sure port 22 is opened.

Instance purchasing - EC2

- on demand instance — short workload, predictable pricing, pay by second
- Reserved (1 \$ 3 years)
 - Reserved instances — long work loads
 - convertible Reserved instances — long workloads with flexible instances
- Savings plans (1 \$ 3 years) — commitment to an amount of usage, long workload

- Spot instances - short workloads, cheap, can lose instances (less reliable)
- Dedicated hosts - book an entire physical server, control instance placement
- Dedicated instances - no other customers will share your hardware
- capacity Reservations - reserve capacity in a specific AZ for any duration
- Ec2 on demand
 - highest cost, no upfront payment
 - short term, uninterrupted workloads
- Reserved instances
 - 72% discount
 - discount based on time period and upfront payment
 - scope is regional or zonal
 - steady state usage applications like databases
- Convertible Reserve instance
 - can change Ec2 instance type, instance family, OS, scope and tenancy
 - upto 66% discount
- Savings plan
 - discount on long-term usage (72%)
 - commit to a certain type of usage

→ locked to a specific instance type & AWS region

→ flexible with instance size, OS, tenancy

→ spot instances

→ most discounted (~90%)

→ can lose at any point of time if your max price is less than current spot price

→ workloads that are resilient to failure

→ Batch jobs

→ Data analysis

→ Image Processing

→ Any distributed workloads

→ workloads with a flexible start and end time

→ no suitable for critical jobs or databases

→ Dedicated hosts

→ physical server fully dedicated

→ pay per second / reserved (1-3 years)

→ most expensive

→ useful for software that have complicated licensing model

→ companies that have compliance needs

→ Dedicated instances

→ Instances run on hardware that's dedicated to you

→ may share hardware with other instances

in same account

→ no control over instance placement

→ capacity reservations

→ on-demand instance in specific AZ for any duration

→ no billing discount

→ always have EC2 capacity when needed

→ create/cancel anytime

→ on-demand rate whether you run or not

→ suitable for short-term, uninterrupted workloads that needs to be in a specific AZ

→ you can only cancel spot instance requests that are open, active or disabled.

→ cancelling a spot request does not terminate instances

→ you must first cancel a spot request, and then terminate the associated spot instances

Spot fleets

→ spot fleet = set of spot instances + on-demand instances

→ spot fleet will try to meet the target capacity with price constraints

→ private IP machines connect to www using a NAT+ internet gateway

Elastic IPs

- if fixed IP for instance is needed, we need an Elastic IP
- can attach it to one instance at a time
- not a very good practice to use Elastic IP (Best is Load balancer)

Placement group

- control over EC2 instance placement strategy
 - cluster
 - spread
 - partition

cluster placement group

- high performance and great throughput with low latency
- great network (10Gbps)
- if rack fails, all instances fail at the same time

Spread placement group

- spans across diff AZ
- reduced risk
- maximize high availability
- disaster management, each instance is isolated from the other

partition placement group

- can span across multiple AZ in same region
- one rack failure affects many instances but not the other rack
- HDFS, HBase, cassandra, Kafka.

Elastic network Interface (ENI)

- logical component in a VPC that represents virtual network card
- **attributes**
 - private IPv4, one or more 2nd IP
 - elastic IP for every IPv4
 - public IPv4
 - one or more security groups
 - A MAC address
- Bound to specific AZ

Hibernation - EC2

- The RAM state is preserved

use cases

- Long-running processing
- Saving the RAM state
- services that takes time to initialize
- cannot be hibernated more than 60 days

EBS (Elastic block store)

- It cannot be used directly
- can be combined with EC2, RDS, ECS, RDS

- For EC2 we need hard disk to the Server
- hard disk will act as EBS
- It is a block-level storage
- OS launching can only be done by EBS
- recommended for data that requires frequent and granular updates
- An EBS volume can be attached to single instance at a time
- EBS allows encryption using Amazon EBS encryption feature
- EBS volumes can be backed up by creating snapshot of the volume, which is stored in S3
- EBS created in one AZ can attach to an instance in the same AZ
- To make volume available outside of AZ create a snapshot and restore the snapshot to new volume anywhere in that region

EBS volume types

- general purpose SSD volumes (gp2)
- provisioned IOPS SSD volumes (io1)
- magnetic volumes (standard)

EBS SSD Volume Types

Characteristic	General Purpose (SSD)	Provisioned IOPS (SSD)
Use Cases	<ul style="list-style-type: none"> System boot volumes Virtual desktops Small to medium DBs Dev and test 	<ul style="list-style-type: none"> I/O intensive Relational DBs NoSQL DBs
Volume Size	1 GB - 16 TB	4 GB - 16 TB
Max. Throughput	160 MB/s	320 MB/s
Max IOPS / volume	10,000	20,000
Max IOPS / instance	65,000	65,000
Max. Throughput / instance	1,250 Mbps	1,250 Mbps
API Name	gp2	io1



EBS HDD Volume Types

Characteristic	Throughput Optimized HDD	Cold HDD	Magnetic
Use Cases	<ul style="list-style-type: none"> Infrequent Data Access Streaming Big Data Logs Cannot be a boot volume 	<ul style="list-style-type: none"> Throughput-oriented for large volumes of data Lowest storage cost is important Cannot be a boot volume 	<ul style="list-style-type: none"> Infrequent Data Access
Volume Size	500 GB - 16 TB	500 GB - 16 TB	1 GB - 1 TB
Max. Throughput	500 MB/s	250 MB/s	40 – 90 MB/s
Max IOPS / volume	500	250	40 – 200
Max IOPS / instance	65,000	65,000	48,000
Max. Throughput / instance	1,250 Mbps	1,250 Mbps	1,250 Mbps
API Name	st1	sc1	standard



General purpose SSD volumes (gp2)

- cost-effective storage that is ideal for broad range of workloads
- single digit millisecond latencies

Magnetic volumes (standard)

- lower cost per gigabyte
- ideal for workloads performing sequential reads, where data is accessed infrequently
- They can be striped together in a RAID configuration for larger size and greater performance

EBS creation

- create new volumes
- restore volumes from snapshots

EBS deletion

- will wipe out data
- can be backed up before deletion using EBS snapshot
- snapshots can be used to create new volumes, increase the size of volumes or replicate data across AZ's

Benefits

- Data availability

- EBS volume is automatically replicated in an AZ to prevent data loss of failure of any single hardware component
- Data persistence
 - EBS volume persists independently of running life of EC2 instance
 - Root EBS volume is deleted, by default, on instance termination but can be modified by changing the Delete on Termination flag
- Data Encryption (EBS encryption feature)

EBS

- Access to raw unformatted block level storage
- Persistent storage
- Detailed metrics captured via CloudWatch Metrics, Applications, gaming, Social, mobile, education, marketing analytics

Standard volumes

- workloads with low or moderate IOPS needs and occasional bursts
- File server, log processing, website, analytics, boot etc.

Provisioned IOPS volumes

- transactional workloads requiring consistent IOPS
- Business applications, MongoDB, SQL server, MySQL, PostgreSQL, Oracle
- 90% of IOPS performance 99.9% of time over a given year
- good latency after PIOPS

consistent snapshot performance

- minimize impact of snapshots on performance
- create snapshots from a read replica of your data
- plan snapshots during off-peak usage
- use cross-region snapshot copy to keep distant regional copies for low latency access
- For workloads targeting above 100 IOPS use PIOPS.

Elastic Load balancer (ELB)

- It distributes and manages the incoming traffic load among several devices to improve network performance.
- distributes client traffic across servers
- improves performance of application

- Routing protocols for load balancer are
 - Round-Robin
 - Server with least no. of requests

Types of Load balancer

- Classic Load balancer
 - traffic across multiple instances in multiple AZ
 - supports both EC2 classic EC2-VPC
 - increases availability of application by sending traffic to healthy instances
 - supports HTTP, HTTPS, TCP and SSL listeners
 - supports sticky sessions

Limitations

- max 20 per region
- 50 listeners
- 8 security groups
- Registered instances 1000
- Subnets per AZ is 1.

Network Load balancer

- It handles sudden and violates traffic across EC2 instances in order to avoid any latency
- new layer 4 load balancing platform
- connection based load balancing
- supports TCP Protocol

→ can handle millions of requests/sec

Limitations

→ 20 LB per region

→ 3000 target groups per region

→ 50 listeners

→ 200 targets per AZ with cross-zone load balancing disabled

→ 200 targets with cross-zone load balancing enabled

→ subnets per AZ 1

Application Load balancer

→ will distribute traffic based on port no (layer 7 of OSI) (content)

→ Reduces hourly cost

→ supports web sockets, HTTP and HTTPS

→ supports microservices and container based application including deep integration with EC2 container service

Limitations

→ 20 LB per region

→ ↑ 3000 per region

→ https 25 certificates (remaining all same)

Key benefits of Load balancer

→ Better performance

→ supports monitoring high check independently

- supports registering target IP
- supports host and path based routing
- route requests to multiple application hosted on single EC2 instance
- Gateway Load balancer
 - easy to deploy, scale and manage third-party virtual appliances
 - eliminates potential point failures and increases availability
 - Scale your virtual appliances automatically
 - monitor continuous health and performance metrics (firewall)
 - runs within one AZ

Limitations

- 20 LB per region
- LB per VPC - 10.

Web application Firewall (WAF)

- It helps address critical attack surface in Application security Ecosystem
 - Application vulnerability
 - zero-days
 - Denial of service
 - BOT attacks
- multi-layered security

- no operational overhead
- default and customizable security
- frictionless set up



- monitors HTTP and HTTPS requests that are forwarded
- web ACL capacity unit (WCUs) is a dimension that is used to calculate and control the operating resources that are used to process your rules within a web ACL. WCUs does not represent a rule for pricing consideration
- can use Amazon Cloud Watch / Amazon Kinesis Firehose (request details)
 - Amazon Cloud Watch - incoming traffic metrics
- AWS managed rules address issues like the Open Web Application Security Project (OWASP) security risks.
- charges per Web ACL's created, no. of rules added per Web ACL and the no. of web requests that you receive

Virtual private networks (VPC)

CIDR - subnet mask

→ It allows part of underlying IP to get additional next values from base IP

Public vs private IP (IPv4)

→ Private IP → 10.0.0.0 - 10.255.255.255

Big networks

→ 172.16.0.0 - 172.31.255.255

AWS default VPC IP range

→ 192.168.0.0 - 192.168.255.255

home networks

→ All the rest are public IP's

Default VPC on AWS

→ has internet connectivity and all EC2 instances inside it have public IPv4 address.

→ we also get a public and private IPv4

→ VPC's are regional in nature.

→ can have multiple VPC's in AWS region (5)

→ max CIDR range per VPC is 5

min size is 128 (16 IP addresses)

max size is 16 (65536 IP addresses)

→ your VPC CIDR should not overlap with other networks

→ 5 IP address are reserved (first 4 and last 1)

→ For CIDR block 10.0.0.0/24

10.0.0.0 - Network Address

10.0.0.1 - VPC router

10.0.0.2 - mapping to amazon-provided DNS

10.0.0.3 - AWS future use

10.0.0.255 - Network Broadcast address

→ Web servers should be in public subnet
and databases in private subnet

Internet Gateway (IGW)

- Allows resources in a VPC connect to internet
- Scales horizontally, highly available and redundant
- One VPC $\xleftarrow{\text{attached}}$ one gateway

Bastion Hosts

- can be used to SSH into private EC2 instances
- Bastion should be in public
- Bastion host security group must allow inbound from internet on port 22 from restricted CIDR

NAT Gateway (Network Address Translation)

- AWS managed NAT, higher bandwidth, high availability, no administration
- Pay per hour for usage and bandwidth

- NATGW is created in a specific AZ, uses an elastic IP
- Can't be used by EC2 instance in the same subnet
- Requires an IGW (private subnet → NATGW → IGW)
- No security groups required
- NAT is resilient within single AZ
- multiple NAT's in multiple AZ's for fault tolerance

Security groups and NACLs

(statefull)

(stateless)

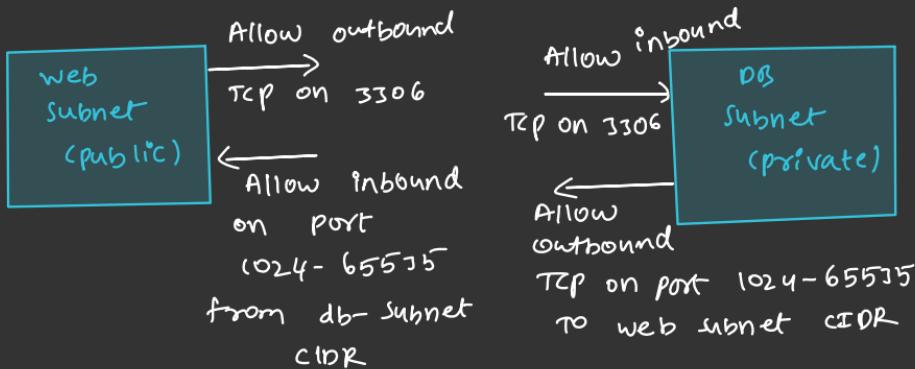
Network access control list (NACL)

- like firewalls from and into subnets
- one NACL per subnet
- higher precedence with lower number
- NACL are a great way of blocking a specific IP address at the subnet level

Ephemeral ports

- For any two endpoints to establish a connection, they must use ports
- clients connect to a defined port and expect a response on ephemeral port

NACL with Ephemeral ports



→ one VPC can have multiple subnets

VPC Peering

- privately connect two VPC's using AWS network
- must not have overlapping CIDR's
- it is not transitive
- must update route table
- can create between VPCs in diff AWS accounts / regions
- can reference SG in peered connection

VPC Endpoints (AWS Private Link)

- accessing objects in S3 bucket for an EC2 instance in private subnet. It has to be through internet and we will charged for that. Is there any way to access objects of S3 on a private network? Yes using VPC Endpoints.

Interface Endpoints (Amazon SNS)

- provisions an ENI (must attach sg)
- supports most AWS services
- \$ per hour + \$ per GB of data processed

Gateway Endpoints

- provisions a gateway and must be used as a target in a route table (no sg)
- supports both S3 and Dynamo DB
- free

Lambda in VPC accessing Dynamo DB

- Access from public internet
 - cuz lambda is in a VPC, it needs a NAT gateway in a public subnet and an internet gateway
- Access from the private VPC network (free)
 - Deploy a VPC gateway endpoint for Dynamo DB
 - change the route tables

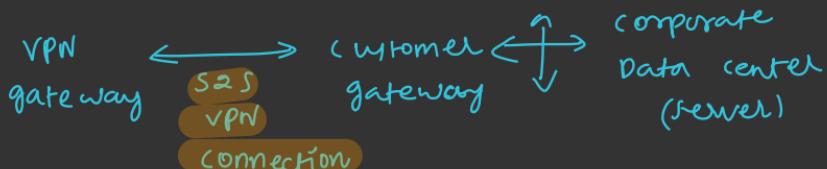
VPC Flow Logs

- capture about IP traffic going into interfaces
 - VPC flow logs
 - Subnet flow logs

→ ENI Flow logs

- helps monitor and troubleshoot connectivity issues
- Flow logs data can go to S3 Cloud Watch Logs
- ELB, RDS, Elasticache, Redshift, workspaces, NAT GW, transit gateway

Site to site VPN



Direct connect (DX) — 30 days to connect

- fastest and costliest way of communication
- need to setup a virtual private gateway on your VPC
- Access S3 and EC2 on same connection.
- use cases
 - increased bandwidth throughput
 - working with large datasets
 - lower cost
 - more consistent network experience
 - applications using real-time data feed
 - hybrid environments (prem + cloud)

Types

→ Dedicated connections

100Gbps capacity

Physical ethernet port

→ Hosted connections

10Gbps capacity

can be added or removed on demand

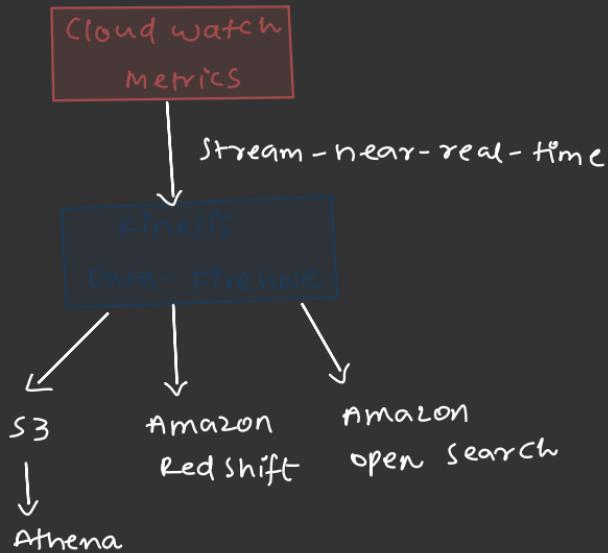
→ can have backup connection in case direct connect fails (expensive)

Amazon CloudWatch Metrics

- It provides metrics for every service in AWS
- Metric is a variable to monitor
- Dimension is an attribute of metric
- Metrics have timestamps
- can create CloudWatch dashboard or custom metrics

CloudWatch Metric Streams

- continually stream CloudWatch metrics to a destination of your choice, with near-real-time delivery and low latency
 - Amazon Kinesis Data Firehose
 - 3rd party service providers:
Datadog, Dynatrace, New Relic,
Splunk, Sumo, Logic...



Cloudwatch Logs

- Log groups: arbitrary name, usually representing an application
- Log stream: instances within application
- Can define log expiration policies
- These logs can be sent to
 - Amazon S3
 - Kinesis Data Streams
 - Kinesis Data Firehose
 - AWS Lambda
 - Elastic search

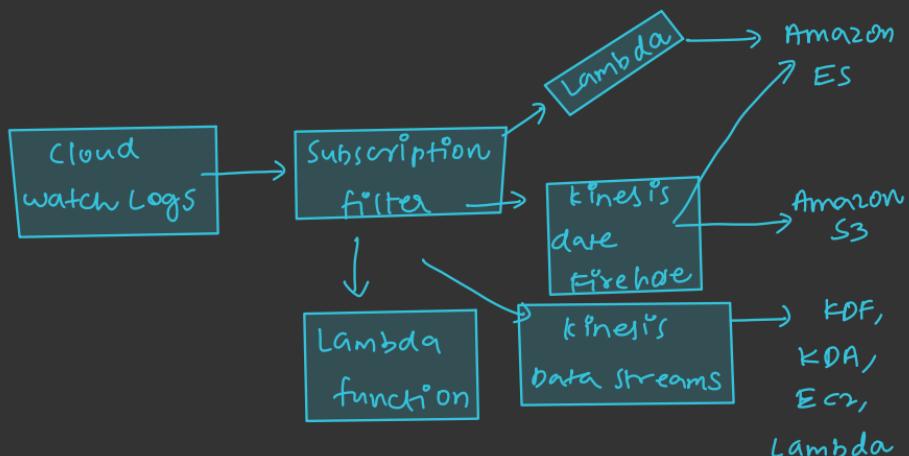
Cloud watch Logs - Sources

- SDK, Cloudwatch Logs Agent, Cloudwatch Unified Agent

- Elastic Beanstalk: collection of logs from application
- ECS: collection from containers
- AWS Lambda: collection from function logs
- VPC flow logs: VPC specific logs
- API gateway
- CloudTrail based on filter
- Route 53: log DNS queries

Cloudwatch Logs - S3 Export

- Log data can take up to 12 hrs to become available for export
- the API call is `createExportTask`



- Kinetic Data Firehose can only delivery to S3, redshift, open search (Amazon ES)

Cloudwatch Logs for EC2

- we need to run a cloudwatch agent on EC2 to push the log files you want (IAM permissions)

Cloudwatch Logs Agent

- old version of agent
- can only send to Cloudwatch Logs

Cloudwatch Unified Agent

- collect additional system-level metrics such as RAM, processes etc..
- collect logs to send Cloudwatch Logs
- centralized configuration using SSM Parameter Store
- CPU, Disk Metrics, Disk I/O, RAM, netstat, processes, Swap Space

Cloudwatch Alarms

- To trigger notifications for any metric states:
 - OK
 - INSUFFICIENT DATA
 - ALARM

Cloudwatch Alarm Targets

- Stop, Terminate, Reboot or Recover an EC2 instance
- triggers Auto Scaling Action

→ Send notification to SNS

Composite Alarms

- For monitoring the states of multiple other Alarms
- we use AND and OR conditions
- Helpful to reduce alarm noise

EC2 Instance Recovery



- To test alarms and notifications, set the alarm state to Alarm using CLI

Cloudwatch Container Insights

- collect aggregate, summarize metrics and logs from containers
- Amazon ECS, Amazon EKS, Kubernetes platform on EC2, Fargate
- In Amazon EKS and Kubernetes, CloudWatch Insights is using a containerized version of the CloudWatch Agent to discover containers

Cloudwatch Lambda Insights

- monitoring and troubleshooting solution

for serverless applications running on AWS Lambda

- Lambda insights are provided as a Lambda layer

Cloudwatch contributor Insights

- see metrics about the top-N contributors
- helps you find top talkers and understand who or what is impacting system performance
- Cloudwatch Application insight is powered by sage maker
- Findings and alerts are sent to Amazon Event bridge and ssm ops center.

Cloud Trail

- provides governance, compliance and audit for your AWS account
- get an history of events / API calls made within AWS account
- if a resource is deleted in AWS, investigate Cloud Trail first
- can send to Cloud watch logs or S3 bucket or (Event bridge event)
- can enable Cloud Trail insights to detect
 - inaccurate resource provisioning

- hitting service limits
 - bursts of AWS IAM actions
 - gaps in periodic maintenance activity
- events are stored for 90 days in cloud trail
- to keep events beyond this period, log them to S3 and use Athena
- Cloud Trail records
 - management Events
 - Data Events
 - Insights Events

AWS config

- Helps with auditing and compliance of AWS account
- can receive SNS notification
- per region service
- data in S3 analyzed by S3
- custom config rules must be defined in AWS Lambda
- AWS config rules does not prevent actions from happening
- No free tier, a bit expensive.

CloudWatch vs CloudTrail vs Config

- CloudWatch
 - Performance monitoring (metrics, CPU, network, etc...) & dashboards
 - Events & Alerting
 - Log Aggregation & Analysis
- CloudTrail
 - Record API calls made within your Account by everyone
 - Can define trails for specific resources
 - Global Service
- Config
 - Record configuration changes
 - Evaluate resources against compliance rules
 - Get timeline of changes and compliance

Cloud Front

- managed AWS service that provides solution for content delivery network (CDN)
- A CDN is a geographically distributed network of proxy servers and their data centers
 - placing datacenters closer
 - high global availability
 - content caching at locations
 - static content, high read performance
 - reducing bandwidth costs
 - improving page load times
- Lambda @ edge to run custom code
- Integrates with AWS Web Application Firewall

- Pay-as-you-go method
- can be used to secure and accelerate web socket traffic as well as API calls
- Allows you to communicate with external HTTPS and talk to internal HTTP backends
- **Regional Edge cache** are global cloud Front servers that is a mediator between host and the edge location
- **Edge locations** are placements of proxy server

Cloud Front - origin

- origin is a location where content is stored, and from which CloudFront gets content to serve the viewers

S3 origin config

- use this type to specify an Amazon S3 bucket that is not configured with static website hosting

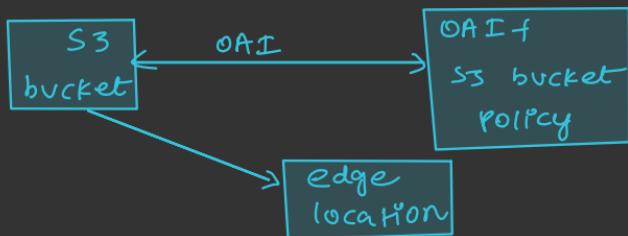
Custom origin config

- S3 for static website hosting
- elastic load balancing load balancer
- AWS Elemental media package
- AWS Elemental media store container
- HTTP server, running on an Amazon

EC2 or any kind of other host

Origin Access Identity (OAI)

- used for sharing private content via Cloud Front



- Adding the public IP of edge location to the security group of EC2 instance
- Adding the public IP of edge location to the security group of ALB and then Add the load balancer to the security group of EC2 instance

Restriction

- use Cloud Front geo restriction

white list: Allow your users to access your content only if they're in one of the countries on whitelist of approved countries (Allow list)

(block list) Black list: prevent users to access ---

process :

→ add a whitelist that contains only the

name of location

- when you send a request, DNS routes the request to the CloudFront edge location
- edge location determines that the user is not allowed
- CloudFront returns an HTTP 403 (Forbidden) to the user
- use a third party geo location service
 - with this you can restrict based on city, zip or postal code, or even longitude or latitude
 - it is recommended to use CloudFront signed URLs when using 3rd party
- low latency, high bandwidth, redundant, scalable, global, cost-effective
- Log analyzer for elastic map reduce
 - generate usage reports containing total traffic volume, object popularity, a breakdown of traffic by client IPs and edge location
- standard wsc format
 - Identify performance bottlenecks caused by slow loading content

- Tiered pricing, rates go down as volume increases
- Reserved cloud front capacity pricing reduces with a longer term commitment
- RTMP distributions for streaming data
- DDOS protection
- Price class All: best performance, expensive
- Price class 200: most regions excluding expensive ones
- Price class 100: only the least expensive regions

Cloud Front - cache Invalidations

- Forcing an entire partial cache refresh to reset the cache

unicast IP vs Anycast IP

- unicast: every server has diff IP address
- AnyCast IP: every server has same IP and the nearest one is chosen

AWS Global Accelerator

- 2 Anycast IP addresses are created for your application
- The Anycast IP send traffic directly to

edge locations

- The edge locations send the traffic to your application
- works with elastic IP, EC2 instance, ALB, NLB, public or private
- consistent performance (low latency, NO cache)
- health checks (disaster recovery)
- security (DDoS protection, only 2 IPs) client

CloudFront vs global accelerator

- Improved performance for both cacheable content and dynamic content
 - content is served at the edge
- global
- performance for a wide range over TCP or UDP
 - proxying packets running in one or more AWS regions
 - good fit for non-HTTP use cases, such as gaming (UDP), IoT (MQTT) or voice over IP
 - good for HTTP use cases that require static IP address
 - good for HTTP use cases that require deterministic, fast regional failover

AWS Databases - RDS

- Divided into SQL and No-SQL
- we can have databases running on EC2 instance
- webserver application in one AZ, and db on another AZ (extra cost for communication between different AZ)
- we run EC2 because
 - Some db's require OS level access
 - Advanced DB option tuning (DB Root)
 - vendor demands
 - DB version that AWS doesn't provide
- we shouldn't run on EC2 because
 - Admin overhead
 - Backup and disaster recovery
 - EC2 running on a single AZ
 - will miss out features from AWS DB
 - skills and setup time to monitor
 - performance will be slower

Relational database Service (RDS)

- Database server - as - a - Service
 - MySQL
 - MariaDB
 - PostgreSQL
 - Oracle
 - Microsoft SQL
 - Amazon Aurora

RDS database Instance

- DD connects with a CNAME.
- db.m6 general, db.r5 memory, db.t3 burst
- EBS is located in the same AZ as RDS
- RDS is vulnerable to failures in that AZ
- Billing per instance and hourly rate for that compute. You will be billed for storage allocated
- We can connect to the database using the database endpoint. And it doesn't change and unique for every RDS instance

Migrating DB from EC2 to RDS

- Get the dump of your existing DB on EC2
- Connect to your RDS DB instance
- Migrate the DB dump that you created to RDS
- Verify if the data is available

RDS multi-AZ (High-availability)

- RDS access only via DB CNAME. The CNAME will point at the primary instance. You cannot access the standby replica for any reason via RDS

- The standby replica cannot be used for extra capacity
- Synchronous replication
- The standby replica cannot be accessed directly unless a failure occurs
- Failover is highly available but not fault tolerant

RDS Backup and restores

RPO - Recovery point objective

- Time between the last backup and when the failure occurred
- Amount of max data loss
- Influences technical solution and cost
- Business usually provides an RPO value

RTO — Recovery time objective

- Time between the DR event and full recovery
- Influenced by process, staff, tech and documentation.

RDS Backups

- First snap is full size of consumed data and next snap will be taking of what's new.

- Snapshot need to manually deleted
- Automatic snapshots
- Every 5 min transaction logs are saved to S3. A database can then be restored to a 5 min snapshots in time
- when you delete the db, they can be retained but they will expire based on their retention period.

RDS - Read Replicas

- Asynchronous replication.
- first stored in primary instance. once it is stored - on disk, it is pushed to the replica. so a small lag
- 5 Read Replicas per db instance
- can provide global performance improvements
- cross-region replication doesn't change RTO.

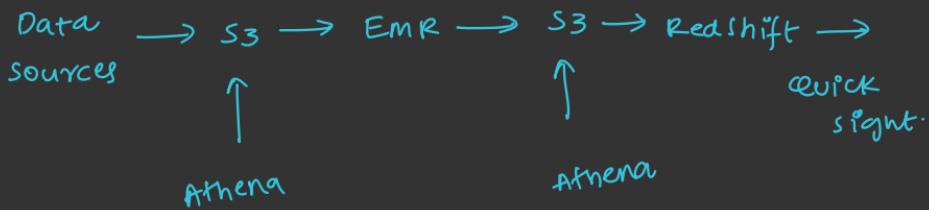
Athena

- It is an interactive query service that makes it easy to analyze data directly from Amazon S3 using standard SQL
- Serverless, no infrastructure, no spin-up time, transparent upgrades.

→ Highly available

Query data from Amazon S3

- no loading of data
- query data in its raw format
- convert to an optimized form like ORC or Parquet for the best performance and lowest cost.
- no ETL required.
- ANSI / SQL
- Presto is used for SQL queries, in-memory distributed query engine ANSI-SQL compatible with extensions
- Hive is used for DDL functionality
- Athena supports multiple data formats
- to improve query performance
 - compress your data
 - use columnar formats
- pay per query, \$5 per TB scanned from S3
- high degree of flexibility
- Hive allows other concepts such "external tables" and partitioning data
- meta data will be stored in metadata store.



Metadata

- Highly available and durable
- Access via DDL statements

Converting ORC to PARQUET

- can use hive cat's
- can use spark
- 20 lines of PySpark code, running on EMR
 - converts 1TB data into 130GB (\$5)
- ways to save costs
 - compress
 - convert to columnar format
 - use partitioning

Use cases

Quicksight $\xleftarrow{\text{Athena}}$ Amazon S3

Macie

- A fully managed data security and data privacy service that uses machine learning

ng and pattern matching to discover and protect sensitive data

- lowers the cost of protecting data
- works on S3 to identify and alert if any sensitive data, such PII is found
- can be sent to Amazon CloudWatch events
- it can help you meet HIPAA and GDPR
- allows to define our own custom sensitive data types
- multi-account support
- charges on the basis of no. of Amazon S3 buckets that are evaluated and the quantity of data processed
- configure a repository if you want to store the info of movie for more than 90 days (repository - S3 bucket)

AWS Kinesis

- The CloudFormation template will launch
 - A VPC
 - A Cloud9 instance to run Kinesis client
 - two S3 buckets store your source taxi trip dataset

- A Kinesis Data Analytics studio application and an associated glue database
- A Lambda function to process data from Kinesis Data Streams

Kinesis

- Real time streaming solution which allows you to ingest buffer and process data.
- Fully managed
- Scalable and can handle any amount of data with low-latencies.
- Kinesis Stream consists of shards, shards consists of records and records consist of data
- Data records are composed of
 - Sequence number
 - Partition key
 - Data blob (will not change - 1 MB)
- A partition key is used by group data by shard within a stream. It uses the partition key that is associated with each data record to determine which shard a given data record belongs to.

- each data record has a sequence no. that is unique partition-key within its shard.
- consumers get records from Amazon Kinesis Data Streams and process them. These consumers are known as Amazon Kinesis Data Streams Application.

SDK vs KPL

- Put record operation allows writing multiple records to multiple shards per request
- only available in Java
- will increase latency

Lambda consumer

- It will consume data from Kinesis data stream
- it inspects the incoming message for unclean records with missing fields and filters them out
- it sends clean records to dynamo DB.
- A throttling error (all failed or only some)

Kinesis Client Library

- connects to the data stream
- instantiates a record processor for every

shard it manages.

- pulls data records from the data stream
- Balances shard-worker associations
- checkpoints processed records
- handles instance failure, automatic load balancer

Kinesis studio notebook

- allows you to interactively query data streams in real time, and easily build and run stream processing applications using standard SQL, Python and Scala.
- data visualization
- exporting data to files
- converting output format for easier analysis

Kinesis Firehose

- easiest way to load streaming data into data stores and analytic tools
- near real-time analysis (by loading data into various other AWS services)
- Fully managed service and requires no ongoing administration
- can also batch, compress and encrypt data before loading, minimizing data storage

- video streams (live video, stores and encrypts at rest, specify custom retention period)
- Data stream (data availability in milliseconds, store data for a min of 24 hours)
- Kinesis Analytics (process data by SQL)

Data Transfer

Data sync

- Data transfer service for moving large amounts of data into AWS
- typically stored on-premises but can also be moved between AWS storage services
- Built in security capabilities (encryption in transit)
- AWS Data Sync agent on a server that connects you to a file system - it copies data and writes to AWS
- Can connect with S3, EFS or FSX

Snowball

- Physical device used for petabyte-scale data transport that can securely transfer large amounts of data
- Addresses lot of common challenges

large scale data-transfer

including high network costs

long transfer times

security concerns

- very secure, has multiple security layers, is tamper resistance and once the transfer is complete everything is completely wiped from snowball device
- two sizes 50TB & 80TB

Snowball Edge

- 100TB data transfer device
- has both storage and compute capabilities
- supports specific amazon EC2 instance type and AWS Lambda functions
- secure as data is encrypted and secure in transit.
- temporary storage tier in remote area locations

Snow mobile

- 45 foot long container pulled by a truck for moving extremely large amounts of data to AWS

- can transfer 100PB per Snowmobile
- can be used for complete datacenter migrations and can import to S3
- fast, secure and cost effective

ENI - Elastic Network Interface

- An ENI is a logical component in a VPC that represents a virtual network card.
- When we move a network interface from one instance to another, network traffic is redirected to the new instance

Why do we need network interfaces?

- Creating a management network
- Using security and network appliances in our VPC
- Creating dual homed instances with workloads / roles on distinct subnets
- Creating low-budget, high-availability solution
- No additional charges of ENI, but your EC2 instances need to support it
- Use EN instead of ENI

- used when you need a reliably high throughputs and speeds between 10 - 100 Gbps
- using single root I/O virtualisation to provide high performance networking
- provides higher bandwidth and packets for second