# Notes on AWS Solutions Architect Associate Exam SAA-C03

Reference: https://www.udemy.com/course/aws-certified-solutions-architect-associate-saa-c03

Reference Page to lecture pdf version 37, also use with DVA notes

Last update: 2024-09-09

# EC2

## EC2 Monitoring

- cpu, network, disk are in cloudwatch default metrics
- memory is custom metric

## EC2 Reservation

- capacity reservation is flexible, can set schedule to stop any time
- start billed when capacity is provisioned.
- zonal/ regional reserved need 1y/3y commitment and only zonal can reserve capacity

## P.73 Spot Instance

- define max price, get instance when spot price < max then

- spot request

    - max price, target capacity (vCpu, ram, number of instance), launch specifications
    - type: one-time/ persistent
    - can set terminate on request end
    - can set which types of instances (vCpu, ram,...)

- when interrupted, set to stop or terminate instance

- when not needed, cancel the spot request and terminate instance (otherwise persistent request will get new instance)

## EBS-backed vs Instance Volume Store

- an EBS-backed EC2 can attach instance store volume (temporary storage)
- only EBS-backed EC2 can be stopped

## Spot Fleet

- set of spot instances + (on demand)

- can have multiple launch pools

- allocation strategy:

    - `lowestPrice`: from the pool with the lowest price (cannot see this)
    - `diversified`: distributed across all instance pools selected
    - `capacityOptimized`: pool with the optimal capacity for the number of instances
    - `priceCapacityOptimized` (recommended): consider Capacity then LowestPrice within the pools

---

## P.79 Private vs Public vs Elastic IP

- private IP, connect to www using internet gateway(a proxy)
- stop EC2 and start again may change IP
- elastic IP can assign to one machine, but can change to another any time
- 1 Account 5 Elastic IP
- using elastic IP to EC2 instance is a bad practice

## P.85 EC2 Placement Groups

- cluster
    - all instance into a low-latency group in single AZ
    - all instance will fail at same time if AZ fail (low availability)
    - for low latency, high network throughput
    - use in big data analytics
    - should launch all instance of same type at once, increase instance afterwards may fail and need to stop and restart all
- spread
    - one instance in partition
    - 7 instance per AZ per partition group
    - max availability but cannot scale large
    - use in most critical system
- partition
    - many instance in one partition
    - 7 partitions per AZ, many AZ in 1 region
    - can support 100s instance
    - usage, e.g. HDFS, HBase, Cassandra, Kafka

## Elastic Network Interfaces (ENI)

- one EC2 instance can have many ENI
- construct extra ENI, and using the extra private IP
- easy for quick network failover in case a instance fails

## EC2 Hibernate

- hibernate = put RAM into EBS then stop the instance
- when start, load RAM from EBS
- EBS have to be encrypted and have enough space to store RAM
- public IP may change after reboot
- can hibernate for at most 60 days

- Hibernate function can only be set at launch, cannot change afterwards

## EC2 Dedicated Host vs Dedicated Instance

- dedicated host support Bring-Your-Own-License/ server-bound software, instance stay on same machine
- provide visibility to CPU, memory, socket, host ID, so can support per-socket/thread/cpu software
- pay for entire physical server
- dedicated instance may be moved to another machine
- can be expensive if high utilization in long term, but can launch as needed

# Database Extra

### P.169 RDS Custom for Oracle and Microsoft SQL Server

- RDS = full managed database and OS
- RDS Custom = with admin access to DB and OS
- support Oracle and MS SQL Server only
- connect to underlying EC2 instance using SSH/SSM session manager
- Recommendation: Take a snapshot and deactivate automation mode when performing customization

### Database SSL

- set `rds.force_ssl` = true and reboot RDS instance
- download the AWS RDS root CA cert, put in EC2 instance which need to use SSL

# Amazon Aurora

- Aurora support serverless and auto scaling (default is Provisioned)

- to migrate from server to serverless, use DMS

- Read Replica

    - add read replica will auto extend read endpoint
    - Cross Region Read Replica
        - simple
    - Global Database
        - 1 primary region for R/W
        - up to 5 secondary region for Read, 16 read replica per region
        - failover take usually < 1s

- Custom Endpoint

    - create custom endpoint for different purpose
    - e.g. an endpoint for large instance (faster computation)
    - usually not to use default read endpoints after custom endpoint

- use standby replica in another AZ for auto fail over

- read replica can be promoted to primary manually

- Aurora replica give sub-ms latency while RDS replica give seconds level latency

- also Aurora clone is much faster then RDS clone

- Proxy fleet

    - reduce number of open connections to improve efficiency
    - integrate with IAM auth
    - auto scales in behind

- Machine Learning

    - support SageMaker and Comprehend
    - use SQL query to do recommendations

## Database Backup

- Manual
    - backup can store forever
- Automated
    - store 1 to 35 days (can disable in RDS)
    - RDS: transaction log every 5 mins
    - Aurora: point-in-time recovery
- to restore from a backup, create a new database
    - MySQL support snapshot and restore on S3

## Aurora Clone

- faster than snapshot + restore
    - use `copy-on-write`, initially same data volume as original DB cluster
    - good use in create staging env

## RDS Proxy

- fully managed
- reduce stress on db resource and minimize open connections
- also reduce RDS/ Aurora failover time
- enforce IAM auth for DB, store credentials in Secrets Manager

## Common Ports

- PostgreSQL: 5432
- MySQL/ MariaDB: 3306
- Oracle: 1521
- MS SQL: 1433
- Aurora: PostgreSQL = 5432, MySQL = 3306

# ElastiCache

- support IAM auth for Redis Enterprise clusters

- IAM policy on ElastiCache used for API-level security

- Redis

    - support sorted set, unique element + ordering

- Redis AUTH

    - password/token
    - support SSL in-flight

- Memcached

    - support SASL-based authentication

## Common Patterns

- Lazy loading
- Write through, no stale data
- Session State

---

# S3

## S3 Batch Operation (SAA)

- S3 Inventory get object list -> S3 Select filter -> S3 Batch Operation
- usage example:
    - encrypt all unencrypted files
    - mass modify metadata
    - invoke lambda on many objects
    - restore objects from S3 Glacier
- for better management and auditing

## S3 Glacier Retrieval

- Bulk retrieval: large-scale data in several hours
- Expedited retrieval: small data within 1-5 mins
- use Provisioned Retrieval Capacity for faster expedited retrieval

## S3 Glacier Vault Lock (SAA)

- WORM model (Write Once Read Many), apply to entire Glacier
- vault lock policy, lock for future edits (cannot be changed/deleted)

## S3 Object Lock (SAA)

- WORM model per object, S3 bucket must have versioning enabled
- Retention mode - Compliance:
    - no change/delete
    - retention mode cannot be changed
    - retention period cannot be shortened
- Retention mode - Governance:
    - most user cannot overwrite/delete
    - admin have special permissions to change retention/ delete the object
- Retention Period: protect the object for a fixed period, it can be extended
- Legal Hold:
    - protect object indefinitely
    - can placed/ removed by `s3:PutObjectLegalHold` IAM permission

## S3 Lifecycle

- S3 -> S3 standard IA/ one-zone IA, minimum need to store 30 days first
- S3 -> Glacier can have 0 days
- min storage duration for IA: need to paid for 30+ days

## S3 sync command

- aws s3 can sync between local/s3 or s3/s3 or s3/local

## S3 replication

- may take up to 15 mins

## S3 access control on other accounts

- IAM: grant user in own account
- ACL: grant other aws account
- Bucket policy: own account or other account, attach to users and group

## S3 logs

- use server access logging for detailed log into another bucket

# CloudFront

## P.339 Price

- the more data flow used, the cheaper the price is (lowest after 5PB)
- Price Classes
    - Price Class All
    - Price Class 200, most region
    - Price Class 100, cheapest region

## Cache Invalidation

- default cache is 24 hours
- can enforce refresh entire / partial cache
- invalidate CloudFront then invalidate edge location

## AWS Global Accelerator

- anycast IP: multiple server having the same ip
- 2 static anycast ip for an app
- route traffic to edge location using Anycast IP, then go to actual server through private AWS network
- can improve streaming through TCP or UDP
- work with elastic IP, EC2, ALB, NLB
- CloudFront/ Global Accelerator integrate with AWS Shield for DDOS
- good for non-HTTP usage/ require static IP/ fast regional failover

## CloudFront with S3

- use presigned url or cookies for private file
- use OAI/ OAC to secure S3 file access from CloudFront, other access will be blocked

## CloudFront SSL Cert

- put SSL cert in AWS Certificate Manager/ IAM certificate store

## CloudFront Geo Block

- geo block applies to whole distribution
- use third party geo location service if need blocking partial content

---

# P.349 AWS Storage Extras

## AWS Snow Family

- physical devices data transport solution

- main usage: transfer data in/out from offline to S3, with some edge processing power

- but can also from AWS to offline

- used when connectivity, bandwidth is limited or unstable

- Snowcone

    - small physical disk HDD/SSD with limited computation, 8-14TB

- Snowball Edge

    - larger storage and computation power, 80-210TB

- request device

- install snowball client/ AWS OpsHub

- connect device from offline server

- ship back device, and load to S3

- Snowball -> S3 -> S3 lifecycle policy -> Glacier

    - cannot go to Glacier directly

# Amazon FSx

- fully managed high performance file system

- support SSD/ HDD type

- can be accessed in on-premise using VPN/ direct connect

- Windows

    - SMB, windows NTFS
    - mainly used in Windows, support Linux EC2 instance
    - integrate with Active Directory(AD)
    - can have multi-AZ
    - daily backup to S3

- Lustre = Linux + Cluster

    - POSIX-compliant
    - for High Performance Computing(HPC), machine learning, media processing
    - Scratch File System
        - temp storage, may suffer data loss
        - higher burst
    - Persistent File System
        - data replicated in same AZ
    - integrate with S3
        - hot data: put in Lustre, expensive but higher performance
        - cold data: put in S3, keep a reference on Lustre for lower cost

- NetApp ONTAP

    - SMB, NFS, *iSCSI*
    - most board support
    - usually lower cost
    - storage auto scaling
    - snapshots, point in time instantaneous cloning
    - use SnapMirror to continue replicate data, archive RPO of 5 mins, RTO within hour
    - good for hybrid environment

- OpenZFS

    - NFS

- board support
- usually lower cost
- high IOPS, low latency
- snapshots, point in time instantaneous cloning

| Feature | FSx for Windows File Server | FSx for Lustre | FSx for OpenZFS | FSx for NetApp ONTAP |
|---|---|---|---|---|
| File System | Windows-native | Lustre | ZFS | ONTAP |
| Protocol Support | SMB | Lustre | NFS | NFS, SMB, iSCSI |
| Primary Use Cases | Windows workloads, home directories | HPC, ML, rendering | Low-latency file access, Oracle DB | Multi-protocol access, DR |
| Performance | Up to 2 GB/s throughput | Up to 100s GB/s throughput | Up to 1 million IOPS | Up to 2 GB/s throughput |
| Capacity | Up to 65,536 GB | Up to 100s of PB | Up to 512 TB | Up to 192 TB |
| Access Control | AD integration | POSIX permissions | POSIX permissions | AD integration, UNIX-style permissions |
| Data Protection | Backups, replication | Data compression | Snapshots, replication | Snapshots, replication |
| Deployment Options | Single-AZ, Multi-AZ | Single-AZ | Single-AZ | Single-AZ, Multi-AZ |
| Integration | AWS Managed AD, Self-Managed AD | Amazon S3 | EC2, ECS, EKS | AWS Managed AD, Self-Managed AD |
| Ideal For | Microsoft applications | High-performance workloads | High IOPS, low latency needs | Versatile enterprise workloads |

## Storage Gateway

- useful in hybrid solution (cloud + on-premise), e.g. compliance issue, security

- storage

  - block: EBS, EC2 instance store
  - file: EFS, FSx
  - object: S3, Glacier

- need on-premise virtualization or hardware appliance(FSx, Volume, Tape)

### S3 File Gateway

- on premise <> *NFS/SMB* <> S3 File Gateway <> HTTPS <> S3

- LRU cache in gateway
- need IAM roles for each gateway
- SMB protocol can integrate with AD for authentication
- use lifecycle policy to send to Glacier

## FSx File Gateway

- SMB client <> FSx Gateway <> FSx Windows File Server
- LRU cache
- Windows native compatibility
- useful of group file share

## Volume Gateway

- back up on-premise server volume through volume gateway with iSCSI protocol

- backed by EBS snapshot in S3

- cached volume/ stored volume

## Tape Gateway

- back up data with tape-based process
- Virtual Tape Library backed by S3, Glacier

## AWS Transfer Family

- file transfer through FTP(unencrypted)/ FTPS/ SFTP
- fully managed, integrate with S3/ EFS
- can integrate with any external authentication system, e.g. Active Directory
- pay for provision endpoint per hour + data transfer
- usage: share files, public dataset, CRM

## DataSync

- move large amount of data
    - AWS <> AWS, no agent
    - on-premise/other cloud <> AWS, need agent
- sync to S3, EFS, FSx
- can schedule task for replication (hourly, daily, non real time)
- can set bandwidth limit
    - if on-premise <> AWS and bandwidth is expensive, use snowball
- file permission and metadata is preserved using SMB protocol

---

# SQS

P.397 SQS + Auto Scaling Group

SQS + EC2

- In SQS queue, set up CloudWatch Metric for Queue Length using `ApproximateNumberOfMessages`
- fire CloudWatch Alarm to ASG to provide more EC2 instances to congest the messages

SQS + DB

- transaction in DB may fail if there is sudden rush
- use SQS as a buffer

SQS + Apps

- SQS can be used for buffer/ decouple two applications

## Amazon MQ

- managed broker for RabbitMQ and ActiveMQ
- support many open protocols such as MQTT, AMQP, STOMP, OpenWire
- cannot scale as large as SQS, SNS
- run on server, can run in multi-AZ with failover
    - one active one standby server in another AZ, both connected to same EFS

## P.451 AWS App Runner

- fully managed service for deploy web app and API at scale
- support container and source code
- configure CPU, memory, auto scaling, health check
- with VPC access support to connect to DB, cache, message queue
- cost
    - per vCPU, per GB memory per second, minimum charge 1 min
    - storage cost
    - provisioned concurrency cost
    - data transfer cost
    - build minutes

# Serverless

## Lambda

### P.462 AWS Lambda SnapStart

- 10x lambda performance for Java 11+
- function invoked from a pre-initialized state
- lambda take snapshot of memory and disk state of initialized function

### Lambda IP

- if non specify VPC, use AWS-managed VPC and do not consume IP address to, have internet access
- use custom VPC, need an IP address

# Database P.522

- RDS
  - relational DB, support SQL
- Aurora
  - PostgreSQL, MySQL compatible
  - more performance, more expensive version
- ElastiCache
  - in-memory cache
  - managed Redis, Memcached
  - store <key,value> and for frequent read
- DynamoDB
  - NoSQL, document(JSON)
- S3
  - key/ value store for large objects
- DocumentDB
  - MongoDB compatible (like Aurora)
  - store index JSON data, in Binary JSON format(more compact)
  - fully managed, replicate across 3 AZ
  - auto scaling and storage scaling in 10GB increment
- Neptune
  - graph database
  - 3AZ, 15 read replicas
  - Neptune Stream, like DynamoDB stream
    - real-time, ordered, no duplicates
    - accessible in RESTful API
- Amazon Keyspaces
  - Cassandra compatible NoSQL, serverless
  - 3AZ
  - on demand or provisioned with auto scaling
  - store IoT device info, time series data
- QLDB
  - Quantum Ledger Database
  - book for recording financial transaction
  - fully managed, serverless, 3AZ
  - manipulate data using SQL-like query
  - immutable, cryptographically verifiable (blockchain like but centralized)
- Timestream
  - time series table
  - much faster than relational DB
  - auto scaling, serverless
  - scheduled query
  - built-in time series analytics functions

# Section 22 P.537 Data and Analysis

## Athena

- serverless query data stored in S3 with SQL

- create the schema of log data first, query, then store result in S3

- support csv, json, orc, ..., charged by scanned volume

- integrate with QuickSight for report/ dashboard

- also commonly used in CloudTrail improvement

- performance

    - columnar data(with Parquet, ORC format)
        - use Glue to convert
    - compress data, partition dataset, use single large files

- support federated query (lambda Data Source Connectors) to other db or on-premise

## Redshift

- PostgresSQL based, used in OLAP(online analytical processing) in data warehouse (PB type scale)
- columnar storage for data & parallel query engine, better performance to Athena
- integrate with BI tools, e.g. QuickSight
- cost by provisioned instance
- auto/manual point-in-time recovery, backed by S3, use snapshot and restore into new cluster

### Redshift Operations

- Redshift is for data warehouse

- when query to Redshift cluster, client communicate with leader node in cluster

- leader node will send to compute nodes to perform

- to load data to Redshift

    1. through Kinesis Firehose
    2. S3 copy with/without private VPC (need enhanced VPC Routing)
    3. EC2 instance with JDBC driver (faster)

- Redshift support 2 AZ, and cross-region data sharing for read

### Redshift Spectrum

- query on S3 data without loading redshift tables
    - redshift cluster -> leader node -> compute node -> *redshift spectrum* -> get data from S3
- Redshift is for data warehouse but Redshift Spectrum can be used for smaller data

# OpenSearch

- search anything instead of key/ index
- support managed cluster or serverless
- use OpenSearch as a complement to another database (need data store on OpenSearch)
- support SQL via plugin

## OpenSearch Pattern

```
DynamoDB Table -> Stream -> Lambda -> store on OpenSearch
```

use API to search item in OpenSearch, then get full item from DynamoDB

```
CloudWatch Log -> Subscription Filter -> Lambda(real)/ Data Firehose(near real) -> store
on OpenSearch
```

```
Kinesis Data Stream -> Lambda / Data Firehose with Lambda for data transform ->
```

# Amazon EMR (Elastic MapReduce)

- create Hadoop clusters to analyze data

- cluster may have hundreds of EC2 instance, auto-scaling and integrated with spot instance

- support Apache Spark, HBase, Presto, Flink...

- master node: manage cluster, long running

- core node: run task and store data, long running

- task node(optional): run task, use Spot Instance to save cost

- can have long-running(buy reserved instance) or transient

# QuickSight

- machine learning BI dashboard
- in-memory computation if data is loaded in QuickSight
- access control: user/ group(enterprise version) (not IAM)
- enterprise version support column-level security
- dashboard = published version, read-only snapshot of analysis

# AWS Glue

- serverless, managed extract, transform, load(ETL) service

- use job bookmark as persist state information

- usage:

    - S3/ RDS -> Glue ETL -> Redshift
    - convert to Parquet format (use columnar data) for better Athena use

- glue data catalog: make metadata of different data source

- glue job bookmarks: prevent reprocess data

- glue elastic view: combine data across multiple source, leverage virtual tale

- databrew: data cleaning

- studio: GUI to run and monitor ETL jobs

- streaming ETL: integrate with Data Stream, Kafka, MSK

## Apache Spark vs AWS glue

- Apache Spark need manual setup but more flexible

# AWS Lake Formation

- full managed service to create data lake
  - data lake = central place to have all data for analytics
  - store in S3
- out-of-the-box source blueprint
- fine grained control on row and column data using data filter (security!)
- also support tag-based access control

# Kinesis Data Analytics (SQL)

- real time analytics on Data Stream/ Firehose + reference data in S3
- fully managed, auto scaling
- output: Data Stream/ Firehose
  - data stream: to Lambda, EC2 instance
  - firehose: S3, redshift, ...

# Kinesis Data Analytics (Flink) (Managed Service for Apache Flink)

- use flink(java, scala, sql)
- more powerful then SQL
- input: data stream/ MSK

# Amazon MSK (Kafka)

- fully managed Apache Kafka, support multi AZ

- provisioned or serverless

- automatic recovery

- data stored on EBS forever

- MSK producer/ consumer

- MSK can have 10MB message(default 1MB) while Data Stream at most 1 MB

- MSK use adding partition Streams use shard splitting/ merging

- MSK can have plaintext transmitting, or TLS, both support KMS at rest

## Serverless Data Pipeline

```
IoT device-> Kinesis Data Stream -> Firehose -> S3(ingestion bucket) -> SQS ->
Lambda -> trigger Athena -> output S3 -> Quicksight, Redshift
```

# Section 23 P.572 Machine Learning

### AWS Rekognition

- object detection, text, person in images and video
- facial analysis/search
- can create database of faces
- content moderation(age control), if high confidence of offensive then push to Amazon Augmented AI(A2I) for manual review

### Amazon Transcribe

- ASR, speech to text, support multi language (auto language detection)
- support streaming
- support removal of PII automatically

### Amazon Polly

- text to speech
- support pronunciation lexicons, e.g. AWS -> Amazon...
- support SSML for customization

### Amazon Translate

- translate

### Amazon Lex & Connect

- Lex is ASR service with natural language understanding
- Connect is a cloud contact center, use with ASR and integrate with CRM
- used in chatbot/ phone bot

### Amazon Comprehend

- NLP, extract keywords, feelings, positive/negative
- Comprehend Medical for clinical text

### SageMaker

- fully managed service to build ML models

## Forecast

- integrate with S3 to give fully managed forecast

## Kendra

- document search service, extract answer from document (text, pdf, HTML, word...)

## Personalize

- provide real-time personalized recommendations, usually work with S3 for data

## Textract

- OCR, extract text/ data from scanned documents

---

# Section 24 P.588 Monitoring & Audit

## P.611 CloudWatch Insights

### Container Insight

- collect, aggregate and summarize metrics and logs
- support ECS, EKS, K8S on EC2, Fargate
- need CloudWatch Agent in EKS, K8S

### Lambda Insight

- provide as lambda layer
- include cold starts/ worker shutdown

### Contributor Insight

- works for any aws generated logs, e.g. VPC, DNS
- create time series about top-N contributors and their usage thought CloudWatch Log
- identify bad user (heaviest network user)
- use prebuilt rules or custom rules

### Application Insight

- automated dashboard to show potential problems with app, help troubleshoot
- the app must run on EC2 instances with selected technology only, with other aws associated resources
- powered by SageMaker
- fire alarm to EventBridge, SSM OpsCenter

## P.623 AWS Config

- help audit and record compliance

- per region service, can aggregated across regions and accounts

- store configuration data in S3 and analyze in Athena

- support aws managed config rules or custom rules (define in lambda)

- rules are triggered when config change/ in interval

- charged by configuration item and rule evaluated

- if non compliant, go to CloudTrail for details

- AWS config do not prevent non compliance, but can trigger SSM Automation Document(Auto-Remediation Action) to disable access keys

  - use remediation retries if still non-compliant

### Notification

- configuration changes -> trigger SNS
- non compliance -> trigger EventBridge -> Lambda

### CloudTrail

- CloudTrail Console store 90 days only, for longer access, use Athena
- CloudTrail store logs with SSE-S3
- CloudTrail log file integrity check with SHA-256 hash, and make a reference every hour
- CloudTrail Lake: tools to query and analyze CloudTrail logs, no need to use Athena

### Others

- EventBridge has archive and replay

# Section 25 IAM Advanced

## Organization Unit

- one management account in an organization, has full access to anything no matter the SCP

- OU can have sub OUs, one member in one OU only

- central billing, share saving plan

- with API for account creation

- CloudTrail and CloudWatch logs are send to central account

- use Service Control Policy(SCP) to control access, default is deny, any explicit deny will deny

  - control the max permission in every account under OU

# IAM policy

- ban IP: `NotIpAddress` with `aws:SourceIp`
- ban region: `aws:RequestedRegion`
- tag, `aws:PrincipalTag/...`,
- MFA, `aws:MultiFactorAuthPresent`
- `aws:PrincipalOrgId` applied if account in inside organization

## IAM role vs Resource based policy

- both can allow S3 cross-account access
- when assume another role, the permission granted to original role is dropped
- resource based policy supports S3, SNS, SQS, Lambda...
- IAM roles: Kinesis, ECS task...

## IAM permission boundary

- support user and roles only, cannot set the group
- set max permission a user can do, need IAM permission + within permission boundary -> pass
- can be used with SCP
- order: any deny -> SCP -> resource based policy -> IAM policy -> IAM permission boundary -> session policy

# IAM Identity Center

- one single login for all AWS account in Organization, business app, SAML2.0, EC2 windows instance
- identity provided by IAM identity store/ AD/ ...
- in multi-account, with single login, can add Permission Set to the actual AWS account
- Attribute-Based Access Control (ABAC): fine grain permission control based on attribute on Identity Center Store

# Active Directory

- use AD with IAM roles

- Microsoft AD = database of objects, eg user, printer, file shares

- object organized in trees and trees form a forest

- any machine communicate with domain controller to check authentication

- AWS managed MS AD: support MFA, large users

    - two-way trust connection with on premise AD, i.e. user can login in any side

- AD Connector: support MFA, proxy to small on-premise AD

    - proxy: user login in on-premise AD only

- simple AD: no on-premise AD

## Control Tower

- setup and govern secure and compliant multi-account AWS env using organizations
- Preventive Guardrail: SCP
- Detective Guardrail: AWS config

---

# Section 26 Security

## KMS

- Key Management Service use CloudHSM behind

**Copy snapshot across region**

- make snapshot, use KMS ReEncrypt with another KMS key, then restore the volume from the new snapshot in another region

**KMS Key policy**

- default policy: allow root, deny all
- custom policy: define user, roles can access the key
- to copy snapshot (using CMK) in cross account, use KMS key policy to authorize the cross-account access, then create a copy with new CMK key

**KMS Multi-Region key**

- replica in other regions, same id, but manage independently

- use in global client site encryption, global DynamoDB, global Aurora (encrypted in client, get data from another region need key in another region)

- in S3 replication, SSE-KMS encrypted objects can be replicated in different regions

- need decrypt from source KMS key and encrypt with destination KMS key

- even if multi region key is used, S3 will decrypt and encrypt agin

**AMI sharing**

- add launch permission in AMI, share KMS key with key policy
- in target account, add IAM role/user to assume the role in source account

## SSM Parameter Store

- support plain text or secure string(backed by KMS)
- security through IAM
- can access secret manager or public parameter(e.g. common AMI id)
- support hierarchy, e.g. /dev/db/password

- advanced parameter (8KB) support TTL with EventBridge Event, but charged

## AWS Secret Manager

- can force rotation of secret every X days with auto generation using lambda
- good integration with RDS
- secret encrypted with KMS
- secret can be replicated across multiple regions, or promote a read replica secret to standalone secret
- multi region secret manager used in multi region DB

## AWS Certificate Manager

- manage and deploy TLS cert, provide inflight encryption
- support public/private cert, public TLS is free
- support auto renew if using aws generated cert, otherwise configure EventBridge rules for daily expiration and invoke SNS
- can integrate with load balancer, cloudfront, API Gateway, cannot use in EC2 instance

**Check cert expiry**

- for AWS managed cert, will renew automatically
- for 3rd party cert:
    - AWS Config managed rule `NON_COMPLIANT` to check cert is about to expire (recommended)
    - monitor `days to expiry` Amazon CloudWatch metric

## ACM with API Gateway Endpoint Types

- edge-optimized(default), request routed by CloudFront Edge location, gateway in one region
    - create custom domain name in api gateway (CNAME or better Alias record)
    - TLS cert attached to CloudFront = us-east-1
- regional, can manually combine with CloudFront to have more control, but more work
    - TLS cert in same region
- private, access from VPC only using ENI and resource policy

## AWS Web Application Firewall (WAF)

- Layer 7 (HTTP), not layer 4 (TCP)
- deploy on ALB(not NLB), API gateway, CloudFront, AppSync GraphQL API, Cognito User Pool
- define web ACL (10000 IPs) for SQL injection, XSS(cross site scripting), rate limiting, geo blocking
- web ACL are regional, except CloudFront is global
- rule group = reusable set of rules
- no NLB? use ALB + Global Accelerator to get fixed IP

## AWS Shield

- AWS Shield Standard: free, protect against DDoS at layer 3, 4
- AWS Shield Advanced: costly, with dedicated support, protect against DDoS at layer 7, 3, 4
- for ALB, CLB, NLB, Elastic IP, CloudFront

## AWS Firewall Manager

- manage rules in all accounts in organization
- charged by rule
- security policy are auto deployed to new resources
    - WAF rules, AWS Shield, Firewall, Security Group

## GuardDuty

- intelligent threat discovery to protect AWS account
- machine learning with *CloudTrail Logs, VPC flow log, DNS log*, (advanced) EBS, Lambda...
    - VPC flow log: suspicious network data
    - CloudTrail log: unauthorized action
    - DNS log: check any malicious IP/domain
- integrate with EventBridge
- can protect against cryptocurrency attacks

## Amazon Inspector

- automated security assessments
- in EC2, need AWS SM agent to analyze OS vulnerabilities
- in container image, scan when pushed
- in lambda, inspect at deployment
- report to Security Hub, integrate with EventBridge

## Macie

- detect sensitive data in AWS using machine learning and pattern matching

# Section 27 VPC

## CIDR

- base IP/ subnet mask
- 0.0.0.0/0 = all IP
- 192.168.0.0/24 = 192.168.0.0 - 192.168.0.255 (256 address)

## Private IP list

- 10.0.0.0/8 = 10.X.X.X (big network)

- 172.16.0.0/12 = 172.16.0.0-172.32.255.255 (AWS default VPC)

- 192.168.0.0/16 = 192.168.0.0-192.168.255.255 (eg home)

- all others are public IP

- at most 5 VPC per region(soft limit)

- at most 5 CIDR in one VPC, min size /28(16 addresses) and max is /16(65536 addresses)

- AWS reserve first 4 + last 1 IP in each subnet

- VPC = public subnets + private subnets

- VPC span across all AZ in one region, subnet span in one AZ only

## P.713 Internet Gateway

- allow resources e.g. EC2 in VPC connect to Internet (with route table)
- created separately from VPC
- 1 VPC - 1 IGW pair
- `EC2 in public subnet <> Route Table <> Router <> IGW <> Internet`

## Bastion Host (public subnet)

- allow SSH to private EC2 instance
- scale horizontally, highly available, redundant
- Bastion Host SG must allow port 22 for SSH
- Private EC2 instance SH must all SH of Bastion/ private IP of Bastion Host
- `SSH to Bastion Host <> SSH to EC2 in private subnet`

## NAT Instance (Network Address Translation, public subnet)

- must disable EC2 setting: source/ destination check
- must have Elastic IP
- NAT instance will change the source/ destination IP, then forward to private EC2/ send from NAT instance IP
- NAT instance support bastion server, port forwarding, SG and NACL while NAT Gateway use NACL only

## NAT Gateway

- AWS managed, created in AZ level, use Elastic IP
- resilient within single AZ, need multi NAT Gateway in multi AZ for fault-tolerance
- charged by hour + network flow
- no SG required
- cannot be used by EC2 in same subnet (for private EC2 + NAT gateway in public)

## P.725 SG/ NACL

- SG stateful while NACL stateless

- SG provide allow rules only

- one NACL per subnet, new subnet with use default NACL(allow all)

- SG evaluate all rules, NACL evaluate in order

- NACL has precedence number(1-32766), will allow/ deny starting from the lowest number with explicit rule

- new NACL will deny all by default

## Ephemeral Ports

- client connect to defined port, and expect response from ephemeral port(random port, system based), need to allow outbound ephemeral port
- linux ephemeral port range: 32768-60999
- windows ephemeral port range: 49152-65535

## VPC Peering

- connect 2 VPC privately, can cross ac/region, not transitive
- require non-overlapping CIDR
- need update rout table in each VPC subnet
- can reference SG in peered VPC (if same region)

## VPC Endpoint/ AWS Private Link

- connect AWS services using private network instead, remove the needs of IGW, NATGW...

- Interface Endpoint

    - provide elastic IP, need SG, support most AWS services
    - charge by hour + data transfer

- Gateway Endpoint

    - connect to S3, DynamoDB
    - free of charge, use with route table
    - use gateway endpoint except need on-premise connection or different VPC

- lambda <> DynamoDB

    - `lambda <> NAT <> IGW <> DynamoDB`
    - `lambda <> VPC Endpoint for DynamoDB <> DynamoDB`

## P.740 VPC Flow Logs

- VPC, subnet, ENI flow log

- send to S3 (with Athena), Cloudwatch (with insights/ alarm), Data Firehose, also capture aws managed service flow log

- flow log syntax:

    - `srcaddr, dstaddr, srcport, dstport`
    - `action: ACCEPT/REJECT`

## Site-to-Site VPN

- Virtual Private Gateway (VGW) in VPC

- Customer Gateway (CGW) in on-premise, physical device or software

- use CGW public IP or NAT public IP (with NAT-T enabled) if CGW is private

- need to enable Route Propagation for VGW in route table

- if ping EC2 from on-premise, need to allow ICMP protocol in SG

- VPN CloudHub: multiple VPN connection to one VGW, allow communication between on-premise

- set up dynamic routing + route table in VGW

## Direct Connect(DX)

- dedicated private connection from remote network to VPC, for better bandwidth, ultra low latency

- support IPV4, IPV6, lower cost if have large data transfer

- not encrypted, can DX+VPN for IPsec-encrypted private connection

- connection `On-premise DC <> AWS direct connect location <> VPC`

- on-premise: customer route/ firewall

- AWS Direct Connect Location:

    - customer/ partner cage: customer/ partner router
    - aws cage: AWS Direct Connect Endpoint

- VPC: Virtual Public Gateway

- private virtual interface: access to aws private subnet through `AWS Cage <> private cable`

- public virtual interface: access to S3, Glacier through `AWS Cage <> Public`

- direct connect gateway: use for multi VPC

    - `On-premise <> Direct Connect Location <> Direct Connect Gateway <> VGW`

- Connection type:

    - dedicated connection : 1/10/100 Gbps, physical ethernet port, request to AWS
    - hosted connection: 50 Mbps - 10 Gbps, on demand, request to Direct Connect Partner if low bandwidth
    - take >1 month

- use site-to-site VPN for backup plan for Direct Connect(DX)

- high resiliency: 2 on-premise <> 2 AWS direction connection

- max resiliency: 2 on-premise, each has 2 router to connect to 2 separate device in connect location

## P.757 Transit Gateway

- transitive peering between VPCs, on-premise, hub-and-spoke connection

- also `transit gateway <> VPN <> customer gateway/ TGW <> Direct Connect Gateway`

- regional resource, can cross-region

- can have cross account using resource access manager(RAM)

    - transit gateway split a direct connection gateway into different VPC

- need route table to limit which VPC-VPC

- only feature to support IP multicast

- charged by network flow

- Site-to-Site VPN ECMP (equal cost multi path routing)

    - increase bandwidth of VPN to connect AWS by multiple site to site VPN
    - VPN - VGW, 1.25Gbps (2 tunnels in VPN)
    - VPN - Transit Gateway, 1=2.5Gbps

## VPC Traffic Mirroring

- capture traffic in VPC, from: ENI, to: ENI or NLB
- same VPC or cross-VPC with VPC peering
- usually set the destination to security apps in EC2 in ASG with NLB

## IPv6

- all IPv6 in AWS are public, IPv4 cannot be disabled(dual stack mode, private/public IPv4 + public IPv6)

- if cannot launch EC2 in subnet, check if all IPv4 address are occupied

- use Egress-only Internet Gateway with IPv6 to have outbound IPv6 connection

- IPv4: `private EC2 <> NAT <> IGW <> internet`

- IPv6: `private EC2 <> Egress only gateway -> outbound`, no inbound

## AWS Network Firewall

- provide layer 3-7 protection, in any direction to protect entire VPC (WAF on layer 7 only)

- internally use managed gateway load balancer

- rule managed by firewall managers

- filter by IP, port, protocol, domain level, regex

- allow, drop or alert if rule matched

- active flow inspection

- alert to S3, cloudwatch, kinesis data stream

# Section 28 P.777 Disaster Recovery and Migration

- Disaster Recovery Type:

- on-premise -> on-premise: very expensive
- on-premise -> cloud: hybrid
- cloud region A -> cloud region B

- RPO = recovery point objective (how long is data lost)

- RTO = recovery time objective (how long to resume service = downtime)

## DR Strategy

- Backup and Restore

  - hight RPO = how frequent the snapshots are, also high RTO
    - storage gateway -> S3
    - snowball -> glacier, long
    - EBS, RDS, redshift -> snapshot

- Pilot Light

  - small version, most critical part are on cloud
  - data replication to RDS, EC2 set but not running

- Warm Standby

  - full system running on cloud, with minimal scale
  - scale to prod load if diaster (EC2 auto scaling, default minimum) + RDS slave

- Hot Site/ Multi Site

  - full system with scale running on cloud

- all can use Route 53 to route problematic server to healthy AZ

## DMS Database Migration Service

- support homogeneous/ heterogeneous (need SCT) db migration

- support EC2 instance with DMS service/ serverless to perform replication

- supported database:

  - on-premise/ EC2 instance DB
  - Azure SQL DB
  - RDS and Aurora
  - S3, DocumentDB
  - Kafka
  - Amazon Neptune
  - Redis

- in heterogeneous migration, use DMS + Schema Conversion Tool

  - in `on-premise -> cloud`, SCT is put on on-premise and DMS instance put on VPC public subnet

- support multi-AZ synchronous replication

## RDS and Aurora MySQL/ PostgreSQL Migration

- RDS MySQL/PostgreSQL -> Aurora MySQL/ PostgreSQL
    - DB snapshot, has downtime
    - create Aurora read replica, wait for replication lag = 0 and promote it to own DB
- external MySQL
    - Percona XtraBackup to create file backup in S3, then create Aurora
    - MySQLDump to migrate directly (slower)
- external PostgreSQL
    - create backup, put in S3, use aws_s3 Aurora extension
- use DMS for continuous replication

## Other Migration Strategy

- Amazon Linux AMI iso file is available
- VM import/export to migrate to/from cloud
    - may cause some downtime
- AWS Application Discovery Service, gather info on on-premise server for planning
- Server Migration Service(SMS), incremental replication of live servers

## AWS Backup

- fully managed, centrally manage backup EC2, EBS, S3, RDS, database, fil system, storage gateway...

- support cross region and cross account

- support point-in-time recovery (if supported originally)

- support on demand, scheduled backup, tag-based policy (e.g. backup env=prod only)

- backup plan = frequency, window, move to cold storage, retention rule

- support WORM, cannot edit or delete

- automated backup max retention = 35 days, AWS backup is 100 years

## AWS Application Discovery Service

- gather info on on-premise server, e.g. server utilization, dependency mapping

- Agent-less discovery(Discovery Connector), configuration, performance

- Agent based (Discovery Agent), configuration, performance, detailed network

- result data in AWS migration hub

## Application Migration Service(MGN)

- lift-and-shift (re-host) solution to AWS

- need to install AWS Replication Agent on source servers

- use AWS replication agent to do continuous replication to build a staging, then promote to prod

- VMWare Cloud on AWS, for migration from using VMWare Cloud service on-premise to AWS

# Section 29 P.803 Architecture

## Events

- SQS(FIFO) + DLQ in SQS side -> Lambda Poll SQS x 5 retry

- SNS -> Lambda(blocking) x 5 retry + DLQ on lambda side

- SDK -> SNS -> many subscribers (fan-out)

- S3 object creation/ deletion -> SNS, SQS, Lambda

- S3 event -> EventBridge

    - advanced filter with json rules
    - multiple destination
    - EventBridge: archive, replay

- API call -> CloudTrail -> EventBridge -> SNS alert

- API Gateway -> Data Stream -> Data Firehose -> store in S3

## Cache

- CloudFront(edge) -> S3
- CloudFront + TTL -> API Gateway + TTL -> EC2/Lambda + (DB/RDS + Redis/Memcached/DAX)

## Block IP

- public EC2 in VPC

    - set up NACL deny rule
    - SG is not helpful(allow rules only)
    - firewall software on EC2 to not process the request, but have CPU cost

- ALB + private EC2

    - in EC2 side, the destination IP is ALB
    - NACL deny rule

- NLB + private EC2

    - NACL deny

- ALB + WAF

    - WAF for IP address filtering, but has some cost

- CloudFront + ALB

    - all traffic comes to ALB is from CloudFront address, so NACL no helpful
    - use WAF on CloudFront or CloudFront Geo Restriction

# High Performance Computing (HPC)

- high computation power X short period = low cost

**Data Transfer**

- Direct Connect: private network
- Snowball, Snowmobile: move PB of data to cloud
- DataSync: move between on-premise <> S3, EFS, FSx

**EC2**

- compute optimized, gpu optimized

- spot instance + auto scaling for money saving

- use EC2 Cluster Placement Group (machine in same rack) for better network performance, 10Gbps inside

- EC2 Enhanced Networking (SR-IOV)

    - Elastic Network Adapter(100Gbps)
    - Legacy, Intel 82599 VF
    - give high bandwidth, high packer per second, low latency

- Elastic Fabric Adapter

    - more improvement from ENA, Linux only
    - Message Passing Interface(MFI) standard
    - message bypass Linux OS for even lower latency

**Storage**

- EBS, 256000 IOPS for io2

- Instance store, but may lost data

- S3, object storage

- EFS

- FSx for Lustre

## EFS lifecycle

- transit to IA type only after 90 days
- max EFS lifecycle policy is 365 days

**Automation**

- AWS batch, run jobs in parallel and span multi EC2 instance
  - usually for ML task
- AWS ParallelCluster
  - config in text file, auto create resources
  - can enable EFA

# Miscellaneous

## Cost Explorer

- Resource Optimization
  - identify under-utilized EC2 instance

## Compute Optimizer

- suggest optimal instance type based on historical usage

## Trusted Advisor

- guidance on provision AWS resource in best practice
- check reserved instance will be expired in 30 days

## ECS launch template vs launch configuration

- launch template is more flexible, support mixture of instance (on-demand + spot, different type)

## VPC sharing vs VPC peering

- sharing is share the VPC to other account within same organization
- peering is joining two VPC/subnet with disjoint CIDR across different account and region

## ASG

- EC2 termination strategy:
  - if multi AZ, choose the one with most instance not protected
  - pick the one with oldest launch template
  - then pick the one with closet billing hour
  - then random
- Predictive scaling
  - assume all instance has same capacity
- instance warm-up time condition
  - do not route traffic to new instance until warm-up time (for instance having long start time)

## Lambda

- lambda env var are encrypted with default kms key, better to use new KMS key

API Gateway

- with throttling limit, the client SDk will retry automatically
- throttling can bursts APIs.

# RDS

## RDS read replica

- async

## RDS event

- operation event only, eg db instance event, snapshot event, use native function or a stored procedure -> lambda for data changes

## RDS monitoring

- default does not have memory
- use enhanced monitoring(with agent) to track on different process
- CloudWatch get CPU utilization from hypervisor of DB instance

## EKS

- use secret encryption with new AWS KMS key on EKS cluster to store sensitive data within etcd key-value store

## Fargate Storage

- default Fargate task has 20GB ephemeral storage (container image counts into it), up to 200GB

## SSL cert in ALB

- when having multiple domains
- use one cert with Subject Alternative Names
  - add/ remove new domain need to provision new cert
- use multiple cert
  - more expensive but flexible
  - can add/ remove individual cert

## Route 53

- active-active: all endpoints will receive traffic
- active-passive: route to standby endpoint only if primary is not available
- A record: domain name -> IPV4 address
- AAAA record: domain name -> IPV6 address
- Cname record: domain name -> domain name
- MX record: to email service
- alias record(53 specific): point to AWS resources, e.g. CloudFront, S3, ELB, without using Cname
  - e.g. alias with type A: domain map -> AWS resource without IP

## AWS Outpost

- on-premise AWS infrastructure server

## AWS License Manager

- control software license usage, block creation of new resource if exceed limit

## AWS Health

- check for how service and resource changes(AWS issues) affecting yours

## K8S Autoscaling

- Karpenter
  - more dynamic and flexible
- Cluster Autoscaler
  - manage node counts based on utilization

## RDS enhanced monitoring metrics

- RDS child processes
- OS processes

## Load Balancing to on-premise + AWS

- use ALB with weighted target group routing to divert some traffic to on-premise, some AWS
- NLB cannot

## Database Choices

- OLAP: Redshift, Athena, EMR, QuickSight
- OLTP: RDS, Aurora, DynamoDB, DocumentDB

## AppSync

- support GraphQL(mostly used) and restful API
- good for aggregating data from multiple database table

## EC2 instance type

- end with 'a' is AMD cpu, not support 82599 VF

## System Manager Run Command

- manage and automate tasks on multiple EC2 instance, e.g. update, patch, install software

## DataSync vs Storage Gateway

- DataSync is for data transfer between on-premise and cloud
- Storage Gateway is for backup, archive, disaster recovery

## IAM database authentication:

- Network traffic to/from database is encrypted with SSL
- use IAM to centrally manage access to your database resources
- in EC2, use profile credentials specific to your EC2 instance to access your database instead of a password, for greater security

## AWS Elastic Disaster Recovery (AWS DRS)

- crash-consistent recovery point objective (RPO) of seconds, and a recovery time objective (RTO) typically ranging within an hour

## Amazon Pinpoint

- Scalable 2-way (outbound/inbound) marketing communications service
- Supports email, SMS, push, voice, and in-app messaging

## Prometheus/ Grafana

- Prometheus: metrics monitoring system
- Grafana: visualization tool
- create workspace on AWS Managed Prometheus to collect metrics, and set it as data source to AWS Managed Grafana

## Resource Access Manager

- prefix list: set of CIDR/ IP address
- RAM allows sharing of resources by creating a resource share using customer-managed/ AWS managed prefix lists

## AWS Proton

- streamlined way to manage microservices and containerized applications

## AWS Local Zones

- AWS infrastructure in some large populations area, for low latency

## AWS Wavelength Zone

- ultra low latency to mobile device by extending AWS to telecom network