

AWS Certified Cloud Practitioner (CLF-C02)

AWS Certified Cloud Practitioner(CLF-C02) Notebook

NOTEBOOK 2024

www.akcoding.com

Table Of Contents

Cloud Computing.....	9
Problems with traditional IT approach.....	9
What is Cloud Computing?.....	10
The Five Characteristics of Cloud Computing.....	10
Six Advantages of Cloud Computing.....	10
Problems solved by the Cloud.....	11
Types of Cloud Computing.....	11
Example of Cloud ComputingTypes.....	11
Pricing of the Cloud – Quick Overview.....	11
AWS Regions.....	12
AWS Availability Zones.....	12
AWS Points of Presence (Edge Locations).....	12
Tour of the AWS Console.....	12
Shared Responsibility Model diagram.....	13
IAM Section.....	13
IAM: Users & Groups.....	13
IAM: Permissions.....	14
IAM Policies inheritance.....	14
IAM – Password Policy.....	14
Multi Factor Authentication - MFA.....	14
MFA devices options in AWS.....	15
How can users access AWS ?.....	15
What's the AWS CLI?.....	16
What's the AWS SDK?.....	16
IAM Roles for Services.....	16
IAM Security Tools.....	16
IAM Guidelines & Best Practices.....	17
Shared Responsibility Model for IAM.....	17
IAM Section – Summary.....	17
EC2 Section.....	18
Amazon EC2.....	18
EC2 sizing & configuration options.....	18
EC2 User Data.....	18
EC2 InstanceTypes - Overview.....	18
EC2 Instance Types – General Purpose.....	19
EC2 Instance Types – Compute Optimized.....	19

EC2 InstanceTypes – Memory Optimized.....	19
EC2 InstanceTypes – Storage Optimized.....	19
Introduction to Security Groups.....	20
Security Groups Deep Dive.....	20
Security Groups Good to know.....	20
Classic Ports to know.....	21
EC2 Instance Connect.....	21
EC2 Instances Purchasing Options.....	21
EC2 On Demand.....	21
EC2 Reserved Instances.....	22
EC2 Savings Plans.....	22
EC2 Spot Instances.....	22
EC2 Dedicated Hosts.....	22
EC2 Dedicated Instances.....	23
EC2 Capacity Reservations.....	23
Which purchasing option is right for me?.....	23
EC2 Section – Summary.....	23
EC2 Instance Storage Section.....	24
What's an EBS Volume?.....	24
EBS Volume.....	24
EBS – Delete on Termination attribute.....	25
EBS Snapshots.....	25
EBS Snapshots Features.....	26
AMI Overview.....	26
AMI Process (from an EC2 instance).....	26
EC2 Image Builder.....	26
EC2 Instance Store.....	27
EFS – Elastic File System.....	27
EBS vs EFS.....	28
EFS Infrequent Access (EFS-IA).....	28
Shared Responsibility Model for EC2 Storage.....	28
Amazon FSx – Overview.....	29
Amazon FSx for Windows File Server.....	29
Amazon FSx for Lustre.....	29
EC2 Instance Storage - Summary.....	29
Elastic Load Balancing & Auto Scaling Groups Section.....	30
Scalability & High Availability.....	30
Vertical Scalability.....	30
Horizontal Scalability.....	30
High Availability.....	30
High Availability & Scalability For EC2.....	30

Scalability vs Elasticity (vs Agility).....	31
What is load balancing?.....	31
Why use a load balancer?.....	31
Why use an Elastic Load Balancer?.....	32
What's an Auto Scaling Group?.....	32
Auto Scaling Groups – Scaling Strategies.....	32
ELB & ASG – Summary.....	32
Amazon S3 Section.....	33
Section introduction.....	33
Amazon S3 Use cases.....	33
Amazon S3 - Buckets.....	33
Amazon S3 – Objects (cont.).....	34
Amazon S3 – Security.....	34
S3 Storage Classes.....	34
S3 Durability and Availability.....	34
S3 Standard – General Purpose.....	35
S3 Storage Classes – Infrequent Access.....	35
Amazon S3 Glacier Storage Classes.....	35
S3 Intelligent-Tiering.....	35
S3 Encryption.....	36
IAM Access Analyzer for S3.....	36
Shared Responsibility Model for S3.....	36
S3 Storage Classes Comparison Summary.....	36
AWS Snow Family.....	37
AWS Snowcone & Snowcone SSD.....	38
Snowball Edge (for data transfers).....	38
AWS Snowmobile.....	38
AWS Snow Family for Data Migrations.....	39
What is Edge Computing?.....	39
Snow Family – Edge Computing.....	39
AWS OpsHub.....	40
Snowball Edge Pricing.....	40
Hybrid Cloud for Storage.....	40
AWS Storage Cloud Native Options.....	41
AWS Storage Gateway.....	41
Amazon S3 – Summary.....	41
Databases Section.....	42
Databases Intro.....	42
NoSQL Databases.....	42
Databases & Shared Responsibility on AWS.....	42
Amazon RDS Overview.....	42

Advantage over using RDS versus deploying DB on EC2.....	43
Amazon Aurora.....	43
Amazon Aurora Serverless.....	43
Amazon ElastiCache Overview.....	43
DynamoDB.....	44
DynamoDB Accelerator - DAX.....	44
DynamoDB – Global Tables.....	44
Redshift Overview.....	44
Redshift Serverless.....	44
Amazon EMR.....	45
Amazon Athena.....	45
Amazon QuickSight.....	45
DocumentDB.....	46
DynamoDB vs DocumentDB.....	46
Amazon Neptune.....	46
Amazon Timestream.....	47
Amazon QLDB.....	47
Amazon Managed Blockchain.....	47
AWS Glue.....	48
DMS – Database Migration Service.....	48
Databases & Analytics Summary in AWS.....	48
Other Compute Section.....	49
What is Docker?.....	49
Where Docker images are stored?.....	49
Docker versus Virtual Machines.....	49
ECS.....	50
Fargate.....	50
ECR.....	51
What's serverless?.....	51
Why AWS Lambda.....	51
Benefits of AWS Lambda.....	51
Amazon API Gateway.....	52
AWS Batch.....	52
Batch vs Lambda.....	52
Amazon Lightsail.....	53
Other Compute - Summary.....	53
Lambda Summary.....	53
Deploying and Managing Infrastructure at Scale Section.....	53
What is CloudFormation.....	54
Benefits of AWS CloudFormation.....	54
Benefits of AWS CloudFormation.....	54

AWS Cloud Development Kit (CDK).....	54
AWS Elastic Beanstalk Overview.....	55
Elastic Beanstalk.....	55
AWS CodeDeploy.....	55
AWS CodeCommit.....	56
AWS CodeBuild.....	56
AWS CodePipeline.....	56
AWS CodeArtifact.....	57
AWS CodeStar.....	57
AWS Cloud9.....	57
AWS Systems Manager (SSM).....	58
Deployment - Summary.....	58
Developer Services - Summary.....	58
Global Infrastructure Section.....	59
Global AWS Infrastructure.....	59
Global Applications in AWS.....	59
Amazon Route 53 Overview.....	59
Route 53 Routing Policies.....	59
Route 53 Routing Policies.....	60
Amazon CloudFront.....	60
CloudFront – Origins.....	60
S3 Transfer Acceleration.....	60
AWS Global Accelerator.....	61
AWS Global Accelerator vs CloudFront.....	61
AWS Outposts.....	61
AWS WaveLength.....	62
AWS Local Zones.....	62
Global Applications in AWS - Summary.....	62
Cloud Integration Section.....	63
Section Introduction.....	63
Amazon SQS – Standard Queue.....	63
SQS to decouple between application tiers.....	64
Amazon SNS.....	64
Amazon MQ.....	65
Integration Section – Summary.....	65
Cloud Monitoring Section.....	66
Amazon CloudWatch Metrics.....	66
Important Metrics.....	66
Amazon CloudWatch Alarms.....	66
Amazon CloudWatch Logs.....	66
CloudWatch Logs for EC2.....	67

Amazon EventBridge (formerly CloudWatch Events).....	67
AWS CloudTrail.....	67
AWS X-Ray.....	68
Amazon CodeGuru.....	68
Amazon CodeGuru Reviewer.....	68
Amazon CodeGuru Profiler.....	69
AWS Health Dashboard - Service History.....	69
AWS Health Dashboard – Your Account.....	69
Monitoring Summary.....	69
VPC Section.....	70
VPC & Subnets Primer.....	70
Internet Gateway & NAT Gateways.....	71
Network ACL & Security Groups.....	72
VPC Flow Logs.....	72
VPC Peering.....	72
VPC Endpoints.....	73
AWS PrivateLink (VPC Endpoint Services).....	73
Site to Site VPN & Direct Connect.....	73
AWS Client VPN.....	74
Transit Gateway.....	74
VPC Closing Comments.....	75
Security & Compliance Section.....	75
Example, for RDS.....	75
Example, for S3.....	76
DDOS Protection on AWS.....	76
AWS Shield.....	76
AWS WAF – Web Application Firewall.....	76
AWS Network Firewall.....	77
AWS Firewall Manager.....	77
Penetration Testing on AWS Cloud.....	77
AWS KMS (Key Management Service).....	78
CloudHSM.....	78
Types of KMS Keys.....	78
AWS Certificate Manager (ACM).....	78
AWS Secrets Manager.....	79
AWS Artifact (not really a service).....	79
Amazon GuardDuty.....	79
Amazon Inspector.....	80
What does Amazon Inspector evaluate?.....	80
AWS Config.....	81
AWS Macie.....	81

AWS Security Hub.....	81
Amazon Detective.....	82
AWS Abuse.....	82
Root user privileges.....	82
IAM Access Analyzer.....	83
Section Summary: Security & Compliance.....	83
Machine Learning Section.....	84
Amazon Rekognition.....	84
Amazon Transcribe.....	84
Amazon Polly.....	85
Amazon Translate.....	85
Amazon Lex & Connect.....	85
Amazon Comprehend.....	86
Amazon SageMaker.....	86
Amazon Forecast.....	86
Amazon Kendra.....	86
Amazon Personalize.....	87
Amazon Textract.....	87
AWS Machine Learning - Summary.....	87
Account Management, Billing & Support Section.....	88
AWS Organizations.....	88
Multi Account Strategies.....	88
Service Control Policies (SCP).....	88
AWS Organization – Consolidated Billing.....	89
AWS Control Tower.....	89
AWS Resource Access Manager (AWS RAM).....	90
Pricing Models in AWS.....	90
Free services & free tier in AWS.....	90
Compute Pricing – EC2.....	91
Compute Pricing – EC2.....	91
Compute Pricing – Lambda & ECS.....	91
Storage Pricing – S3.....	91
Storage Pricing - EBS.....	92
Database Pricing - RDS.....	92
Content Delivery – CloudFront.....	93
Networking Costs in AWS per GB - Simplified.....	93
Savings Plan.....	93
AWS Compute Optimizer.....	94
Billing and Costing Tools.....	94
AWS Pricing Calculator.....	94
AWS Billing Dashboard.....	94

Cost Allocation Tags.....	95
Tagging and Resource Groups.....	95
Cost and Usage Reports.....	95
Cost Explorer.....	96
Cost Explorer – Savings Plan.....	96
Cost Explorer – Forecast Usage.....	96
Billing Alarms in CloudWatch.....	97
AWS Budgets.....	97
AWS Cost Anomaly Detection.....	98
AWS Service Quotas.....	98
Trusted Advisor.....	99
AWS Basic Support Plan.....	99
AWS Developer Support Plan.....	99
AWS Business Support Plan (24/7).....	99
AWS Enterprise On-Ramp Support Plan (24/7).....	100
AWS Enterprise Support Plan (24/7).....	100
Account Best Practices – Summary.....	100
Billing and CostingTools – Summary.....	100
Advanced Identity Section.....	101
AWS STS (Security Token Service).....	101
Amazon Cognito (simplified).....	101
AWS IAM Identity Center (successor to AWS Single Sign-On).....	102
Advanced Identity - Summary.....	102
Other AWS Services.....	102
AWS Migration Hub.....	102
AWS Step Functions.....	102
AWS Architecting & Ecosystem Section.....	103
Well Architected Framework General Guiding Principles.....	103
AWS Cloud Best Practices – Design Principles.....	103
Well Architected Framework 6 Pillars.....	103
1. Operational Excellence.....	104
2. Security.....	104
3. Reliability.....	104
4. Performance Efficiency.....	105
5. Cost Optimization.....	105
6. Sustainability.....	105
AWS Cloud Adoption Framework (AWS CAF).....	106
CAF Perspectives and Foundational Capabilities Business Capabilities.....	107
AWS CAF – Transformation Phases.....	108
AWS Marketplace.....	108
AWSTraining.....	108

AWS Professional Services & Partner Network.....	108
AWS IQ.....	109
AWS re:Post.....	109
AWS re:Post – Knowledge Center.....	109
AWS Managed Services (AMS).....	109

Cloud Computing

Problems with traditional IT approach

- Pay for the rent for the data center
- Pay for power supply, cooling, maintenance
- Adding and replacing hardware takes time
- Scaling is limited
- Hire 24/7 team to monitor the infrastructure
- How to deal with disasters? (earthquake, power shutdown, fire...)
- Can we externalize all this?

What is Cloud Computing?

- Cloud computing is the on-demand delivery of compute power, database storage, applications, and other IT resources
- Through a cloud services platform with pay-as-you-go pricing
- You can provision exactly the right type and size of computing resources you need
- You can access as many resources as you need, almost instantly
- Simple way to access servers, storage, databases and a set of application services
- Amazon Web Services owns and maintains the network-connected hardware required for these application services, while you provision and use what you need via a web application.

The Five Characteristics of Cloud Computing

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

6 Advantages of Cloud Computing

1. Trade capital expense (CAPEX) for operational expense (OPEX)
 - a. Pay On-Demand: don't own hardware
 - b. Reduced Total Cost of Ownership (TCO) & Operational Expense (OPEX)
2. Benefit from massive economies of scale
 - a. Prices are reduced as AWS is more efficient due to large scale
3. Stop guessing capacity
 - a. Scale based on actual measured usage
4. Increase speed and agility
5. Stop spending money running and maintaining data centers
6. Go global in minutes: leverage the AWS global infrastructure

Problems solved by the Cloud

- Flexibility: change resource types when needed
- Cost-Effectiveness: pay as you go, for what you use
- Scalability: accommodate larger loads by making hardware stronger or adding additional nodes
- Elasticity: ability to scale out and scale-in when needed
- High-availability and fault-tolerance: build across data centers
- Agility: rapidly develop, test and launch software applications

Types of Cloud Computing

- Infrastructure as a Service (**IaaS**)
 - Provide building blocks for cloud IT

- Provides networking, computers, data storage space
- Highest level of flexibility
- Easy parallel with traditional on-premises IT
- Platform as a Service (**PaaS**)
 - Removes the need for your organization to manage the underlying infrastructure
 - Focus on the deployment and management of your applications
- Software as a Service (**SaaS**)
 - Completed product that is run and managed by the service provider

Example of Cloud Computing Types

- Infrastructure as a Service:
 - Amazon EC2 (on AWS)
 - GCP, Azure, Rackspace, Digital Ocean, Linode
- Platform as a Service:
 - Elastic Beanstalk (on AWS)
 - Heroku, Google App Engine (GCP), Windows Azure (Microsoft)
- Software as a Service:
 - Many AWS services (ex: Rekognition for Machine Learning)
 - Google Apps (Gmail), Dropbox, Zoom

Pricing of the Cloud – Quick Overview

- AWS has 3 pricing fundamentals, following the pay-as-you-go pricing model
- Compute:
 - Pay for compute time
- Storage:
 - Pay for data stored in the Cloud
- Data transfer OUT of the Cloud:
 - Data transfer IN is free
- Solves the expensive issue of traditional IT

AWS Regions

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

Each AWS **Region** consists of a minimum of **3 Availability Zones (AZ)**

AWS Availability Zones

- Each region has many availability zones (usually 3, min is 3, max is 6).
- Example:
 - ap-southeast-2a
 - ap-southeast-2b
 - ap-southeast-2c
- **Each availability zone (AZ) is one or more discrete data centers with redundant power, networking, and connectivity**
- They're separate from each other, so that they're isolated from disasters
- They're connected with high bandwidth, ultra-low latency networking

AWS Points of Presence (Edge Locations)

- Content is delivered to end users with lower latency

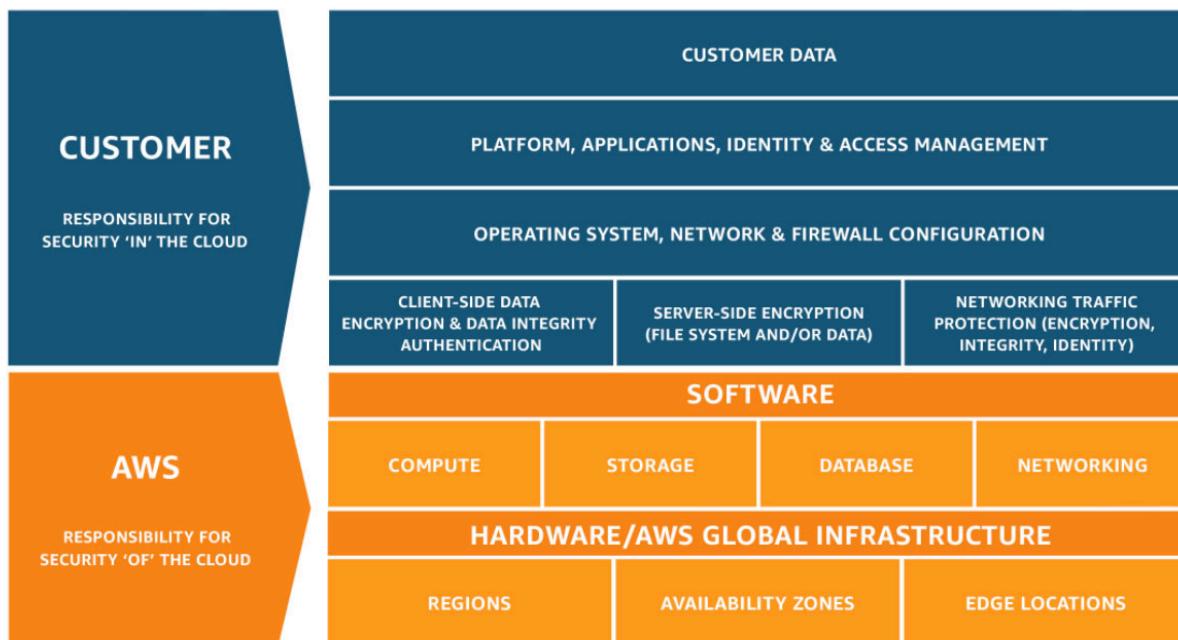
Tour of the AWS Console

- AWS has Global Services:
 - Identity and Access Management (IAM)
 - Route 53 (DNS service)
 - CloudFront (Content Delivery Network)
 - WAF (Web Application Firewall)
- Most AWS services are Region-scoped:
 - Amazon EC2 (Infrastructure as a Service)
 - Elastic Beanstalk (Platform as a Service)
 - Lambda (Function as a Service)
 - Rekognition (Software as a Service)

Shared Responsibility Model diagram

CUSTOMER = RESPONSIBILITY FOR THE SECURITY IN THE CLOUD

AWS = RESPONSIBILITY FOR THE SECURITY OF THE CLOUD



IAM Section

IAM: Users & Groups

- IAM = Identity and Access Management, **Global service**
- Root account created by default, shouldn't be used or shared
- Users are people within your organization, and can be grouped
- **Groups only contain users, not other groups**
- Users don't have to belong to a group, and user can belong to multiple groups

IAM: Permissions

- Users or Groups can be assigned JSON documents called policies
- These policies define the permissions of the users
- **In AWS you apply the least privilege principle:** don't give more permissions than a user needs

IAM Policies inheritance

IAM – Password Policy

- Strong passwords = higher security for your account
- In AWS, you can setup a password policy:
- Set a minimum password length
- Require specific character types:
 - including uppercase letters
 - lowercase letters
 - numbers
 - non-alphanumeric characters
- Allow all IAM users to change their own passwords
- Require users to change their password after some time (password expiration)
- Prevent password reuse

Multi Factor Authentication - MFA

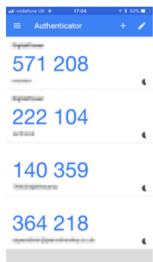
- Users have access to your account and can possibly change configurations or delete resources in your AWS account
- You want to protect your Root Accounts and IAM users
- MFA = password you know + security device you own



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

MFA devices options in AWS

Virtual MFA device



Google Authenticator
(phone only)

Support for multiple tokens on a single device.



Authy
(phone only)

Universal 2nd Factor (U2F) Security Key



YubiKey by Yubico (3rd party)

Support for multiple root and IAM users using a single security key

Hardware Key Fob MFA Device



Provided by Gemalto (3rd party)

Hardware Key Fob MFA Device for AWS GovCloud (US)



Provided by SurePassID (3rd party)

How can users access AWS ?

- To access AWS, you have three options:
 - AWS Management Console (protected by password + MFA)
 - AWS Command Line Interface (CLI): protected by access keys
 - AWS Software Developer Kit (SDK) - for code: protected by access keys
- Access Keys are generated through the AWS Console
- Users manage their own access keys
- Access Keys are secret, just like a password. Don't share them • Access Key ID ~ username
- Secret Access Key ~= password

What's the AWS CLI?

- A tool that enables you to interact with AWS services using commands in your command-line shell
- Direct access to the public APIs of AWS services
- You can develop scripts to manage your resources
- It's open-source <https://github.com/aws/aws-cli>
- Alternative to using AWS Management Console

What's the AWS SDK?

- AWS Software Development Kit (AWS SDK)
- Language-specific APIs (set of libraries)
- **Enables you to access and manage AWS services programmatically**
- Embedded within your application
- Supports
 - SDKs(JavaScript,Python,PHP,.NET,Ruby,Java,Go,Node.js, C++)
 - Mobile SDKs (Android, iOS, ...)
 - IoT Device SDKs (Embedded C, Arduino, ...)
- Example: AWS CLI is built on AWS SDK for Python

IAM Roles for Services

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

IAM Security Tools

- **IAM Credentials Report** (account-level)
 - a report that lists all your account's users and the status of their various credentials
- **IAM Access Advisor** (user-level)
 - Access advisor shows the service permissions granted to a user and when those services were last accessed.
 - You can use this information to revise your policies.

IAM Guidelines & Best Practices

- Don't use the root account except for AWS account setup

- One physical user = One AWS user
- Assign users to groups and assign permissions to groups
- Create a strong password policy
- Use and enforce the use of Multi Factor Authentication (MFA)
- Create and use Roles for giving permissions to AWS services
- Use Access Keys for Programmatic Access (CLI / SDK)
- Audit permissions of your account using IAM Credentials Report & IAM Access Advisor
- Never share IAM users & Access Keys

Shared Responsibility Model for IAM

AWS:

- Infrastructure (global network security)
- Configuration and vulnerability analysis
- Compliance validation

User:

- Users, Groups, Roles, Policies management and monitoring
- Enable MFA on all accounts
- Rotate all your keys often
- Use IAM tools to apply appropriate permissions
- Analyze access patterns & review permissions

IAM Section – Summary

- **Users:** mapped to a physical user, has a password for AWS Console
- **Groups:** contains users only
- **Policies:** JSON document that outlines permissions for users or groups
- **Roles:** for EC2 instances or AWS services
- **Security:** MFA + Password Policy
- **AWS CLI:** manage your AWS services using the command-line
- **AWS SDK:** manage your AWS services using a programming language
- **Access Keys:** access AWS using the CLI or SDK
- **Audit:** IAM Credential Reports & IAM Access Advisor

EC2 Section

Amazon EC2

- **EC2** is one of the most popular of AWS' offering

- **EC2 = Elastic Compute Cloud = Infrastructure as a Service**
- It mainly consists in the capability of :
 - Renting virtual machines (EC2)
 - Storing data on virtual drives (EBS)
 - Distributing load across machines (ELB)
 - Scaling the services using an auto-scaling group (ASG)
- Knowing EC2 is fundamental to understand how the Cloud works

EC2 sizing & configuration options

- Operating System (OS): Linux, Windows or Mac OS
- How much compute power & cores (CPU)
- How much random-access memory (RAM)
- How much storage space:
 - Network-attached (EBS & EFS)
 - Hardware (EC2 Instance Store)
- Network card: speed of the card, Public IP address
- Firewall rules: security group
- Bootstrap script (configure at first launch): EC2 User Data

EC2 User Data

- It is possible to bootstrap our instances using an **EC2 User data** script.
- **bootstrapping** means launching commands when a machine starts
- That script is only run once at the instance first start
- EC2 user data is used to automate boot tasks such as:
 - Installing updates
 - Installing software
 - Downloading common files from the internet
 - Anything you can think of
- The EC2 User Data Script runs with the root user

EC2 InstanceTypes - Overview

m5.2xlarge

- 5: generation (AWS improves them over time)
- m: instance class
- 2xlarge: size within the instance class

EC2 Instance Types – General Purpose

- Great for a diversity of workloads such as web servers or code repositories
- Balance between:

- **Compute**
- **Memory**
- **Networking**
- In the course, we will be using the t2.micro which is a General Purpose EC2 instance

EC2 Instance Types – Compute Optimized

- Great for compute-intensive tasks that require high performance processors:
 - Batch processing workloads
 - Media transcoding
 - High performance web servers
 - High performance computing (HPC)
 - Scientific modeling & machine learning
 - Dedicated gaming servers

EC2 InstanceTypes – Memory Optimized

- Fast performance for workloads that process large data sets in memory
- Use cases:
 - High performance, relational/non-relational databases
 - Distributed web scale cache stores
 - In-memory databases optimized for BI (business intelligence)
 - Applications performing real-time processing of big unstructured data

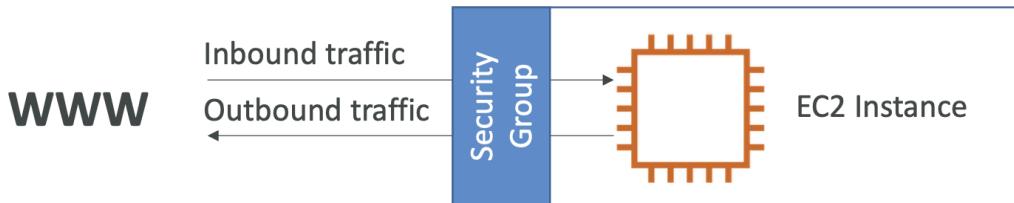
EC2 InstanceTypes – Storage Optimized

- Great for storage-intensive tasks that require high, sequential read and write access to large data sets on local storage
- Use cases:
 - High frequency online transaction processing (OLTP) systems
 - Relational & NoSQL databases
 - Cache for in-memory databases (for example, Redis)
 - Data warehousing applications
 - Distributed file systems

t2.micro is part of the AWS free tier (up to 750 hours per month)

Introduction to Security Groups

- Security Groups are the fundamental of network security in AWS
- They control how traffic is allowed into or out of our EC2 Instances.
- Security groups only contain rules
- Security groups rules can reference by IP or by security group



Security Groups Deep Dive

- Security groups are acting as a “firewall” on EC2 instances
- They regulate:
 - Access to Ports
 - Authorized IP ranges – IPv4 and IPv6
 - Control of inbound network (from other to the instance)
 - Control of outbound network (from the instance to other)

Security Groups Good to know

- Can be attached to multiple instances
- Locked down to a region / VPC combination
- Does live “outside” the EC2 – if traffic is blocked the EC2 instance won’t see it
- **It’s good to maintain one separate security group for SSH access**
- If your application is not accessible (time out), then it’s a security group issue
- If your application gives a “connection refused” error, then it’s an application error or it’s not launched
- All inbound traffic is blocked by default
- All outbound traffic is authorized by default

A Security Group can have **allow** rules only

Classic Ports to know

- 22 = SSH (Secure Shell) - log into a Linux instance
- 21 = FTP (File Transfer Protocol) – upload files into a file share
- 22 = SFTP (Secure File Transfer Protocol) – upload files using SSH
- 80 = HTTP – access unsecured websites
- 443 = HTTPS – access secured websites
- 3389 = RDP (Remote Desktop Protocol) – log into a Windows instance

EC2 Instance Connect

- **Connect to your EC2 instance within your browser**
- No need to use your key file that was downloaded
- The “magic” is that a temporary key is uploaded onto EC2 by AWS
- Works only out-of-the-box with Amazon Linux 2
- Need to make sure port 22 is still opened!

EC2 Instances Purchasing Options

- **On-Demand Instances** – short workload, predictable pricing, pay by second
- **Reserved (1 & 3 years)**
 - Reserved Instances – long workloads
 - 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

- Pay for what you use:
 - Linux or Windows - billing per second, after the first minute
 - All other operating systems - billing per hour
- Has the highest cost but no upfront payment
- No long-term commitment
- **Recommended for short-term and un-interrupted workloads**, where you can't predict how the application will behave
- Use Case:
 - **Development and Testing:**
 - Ideal for software development and testing environments where workloads are unpredictable and require flexibility.

EC2 Reserved Instances

- Up to 72% discount compared to On-demand
- You reserve a specific instance attributes (Instance Type, Region, Tenancy, OS)
- **Reservation Period – 1 year (+discount) or 3 years (+++discount)**
- Payment Options – No Upfront (+), Partial Upfront (++) All Upfront (+++)
- Reserved Instance's Scope – Regional or Zonal (reserve capacity in an AZ)
- Recommended for steady-state usage applications (think database)
- **You can buy and sell in the Reserved Instance Marketplace**

- **Convertible Reserved Instance**
 - Can change the EC2 instance type, instance family, OS, scope and tenancy
 - Up to 66% discount

EC2 Savings Plans

- Get a discount based on long-term usage (up to 72% - same as RIs)
- Commit to a certain type of usage (\$10/hour for 1 or 3 years)
- Usage beyond EC2 Savings Plans is billed at the On-Demand price
- Locked to a specific instance family & AWS region (e.g., M5 in us-east-1)
- Flexible across:
 - Instance Size (e.g., m5.xlarge, m5.2xlarge)
 - OS (e.g., Linux, Windows)
 - Tenancy (Host, Dedicated, Default)

EC2 Spot Instances

- **Can get a discount of up to 90% compared to On-demand**
- **Instances that you can “lose” at any point of time if your max price is less than the current spot price**
- The MOST cost-efficient instances in AWS
- **Useful for workloads that are resilient to failure**
- Batch jobs
- Data analysis
- Image processing
- Any distributed workloads
- Workloads with a flexible start and end time
- Not suitable for critical jobs or databases

EC2 Dedicated Hosts

- **A physical server with EC2 instance capacity fully dedicated to your use**
- Allows you address compliance requirements and use your existing server-bound software licenses (per-socket, per-core, per-VM software licenses)

- Purchasing Options:
 - On-demand – pay per second for active Dedicated Host
 - Reserved - 1 or 3 years (No Upfront, Partial Upfront, All Upfront)
- The most expensive option
- Useful for software that have complicated licensing model (BYOL – Bring Your Own License)
- Or for companies that have strong regulatory or compliance needs

EC2 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 (can move hardware after Stop / Start)

EC2 Capacity Reservations

- Reserve On-Demand instances capacity in a specific AZ for any duration
- You always have access to EC2 capacity when you need it
- No time commitment (create/cancel anytime), no billing discounts
- Combine with Regional Reserved Instances and Savings Plans to benefit from billing discounts
- You're charged at On-Demand rate whether you run instances or not
- Suitable for short-term, uninterrupted workloads that needs to be in a specific AZ

Which purchasing option is right for me?

- **On demand:** coming and staying in resort whenever we like, we pay the full price
- **Reserved:** like planning ahead and if we plan to stay for a long time, we may get a good discount.
- **Savings Plans:** pay a certain amount per hour for a certain period and stay in any room type (e.g., King, Suite, Sea View, ...)
- **Spot instances:** the hotel allows people to bid for the empty rooms and the highest bidder keeps the rooms. You can get kicked out at any time
- **Dedicated Hosts:** We book an entire building of the resort
- **Capacity Reservations:** you book a room for a period with full price even you don't stay in it

EC2 Section – Summary

- **EC2 Instance:** AMI (OS) + Instance Size (CPU + RAM) + Storage + security groups + EC2 User Data
- **Security Groups:** Firewall attached to the EC2 instance
- **EC2 User Data:** Script launched at the first start of an instance

- **SSH:** start a terminal into our EC2 Instances (port 22)
- **EC2 Instance Role:** link to IAM roles
- **Purchasing Options:** On-Demand, Spot, Reserved (Standard + Convertible), Dedicated Host, Dedicated Instance

EC2 Instance Storage Section

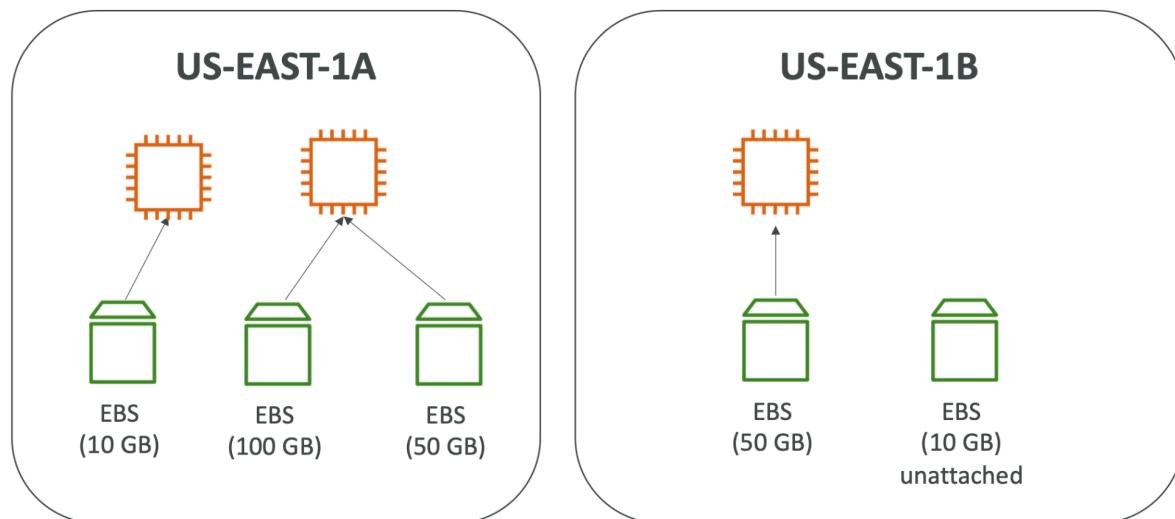
What's an EBS Volume?

- An EBS (Elastic Block Store) Volume is a network drive you can attach to your instances while they run
- It allows your instances to persist data, even after their termination
- They can only be mounted to one instance at a time (at the CCP level)
- They are bound to a specific availability zone
- Analogy: Think of them as a “network USB stick”
- Free tier: 30 GB of free EBS storage of type General Purpose (SSD) or Magnetic per month

EBS Volume

- It's a network drive (i.e. not a physical drive)
 - It uses the network to communicate the instance, which means there might be a bit of latency
 - It can be detached from an EC2 instance and attached to another one quickly
- It's locked to an Availability Zone (AZ)
 - An EBS Volume in us-east-1a cannot be attached to us-east-1b
 - To move a volume across, you first need to snapshot it
- Have a provisioned capacity (size in GBs, and IOPS)
 - You get billed for all the provisioned capacity
 - You can increase the capacity of the drive over time

EBS Volume - Example



EBS – Delete on Termination attribute

Volume Type	Device	Snapshot	Size (GiB)	Volume Type	IOPS	Throughput (MB/s)	Delete on Termination	Encryption
Root	/dev/xvda	snap-09f18f682fd23a1b1	8	General Purpose SSD (gp2)	100 / 3000	N/A	<input checked="" type="checkbox"/>	Not Encrypted
EBS	/dev/sdb	Search (case-insensit	8	General Purpose SSD (gp2)	100 / 3000	N/A	<input type="checkbox"/>	Not Encrypted

- Controls the EBS behavior when an EC2 instance terminates
 - By default, the root EBS volume is deleted (attribute enabled)
 - By default, any other attached EBS volume is not deleted (attribute disabled)
- This can be controlled by the AWS console / AWS CLI
- Use case:** preserve root volume when instance is terminated

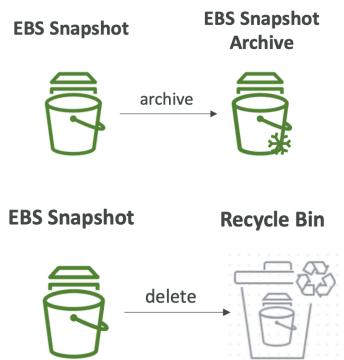
EBS Snapshots

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



EBS Snapshots Features

- **EBS Snapshot Archive**
 - Move a Snapshot to an "archive tier" that is 75% cheaper
 - Takes within 24 to 72 hours for restoring the archive
- **Recycle Bin for EBS Snapshots**
 - Setup rules to retain deleted snapshots so you can recover them after an accidental deletion
 - Specify retention (from 1 day to 1 year)

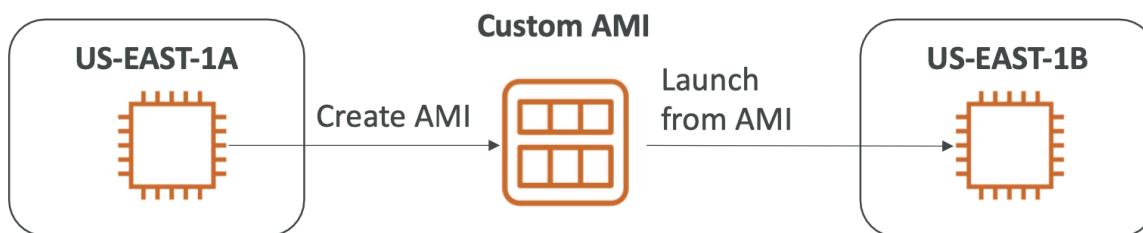


AMI Overview

- AMI = Amazon Machine Image
- AMI are a customization of an EC2 instance
 - You add your own software, configuration, operating system, monitoring
 - Faster boot / configuration time because all your software is pre-packaged
- AMI are built for a specific region (and can be copied across regions)
- You can launch EC2 instances from:
 - A Public AMI: AWS provided
 - Your own AMI: you make and maintain them yourself
 - An AWS Marketplace AMI: an AMI someone else made (and potentially sells)

AMI Process (from an EC2 instance)

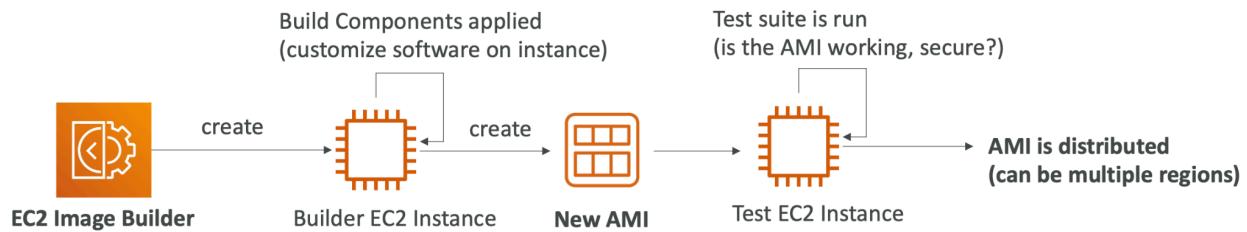
1. Start an EC2 instance and customize it
2. Stop the instance (for data integrity)
3. Build an AMI – this will also create EBS snapshots
4. Launch instances from other AMIs



EC2 Image Builder

- Used to automate the creation of Virtual Machines or container images
- => Automate the creation, maintain, validate and test EC2 AMIs
- Can be run on a schedule (weekly, whenever packages are updated, etc...)

- Free service (only pay for the underlying resources)

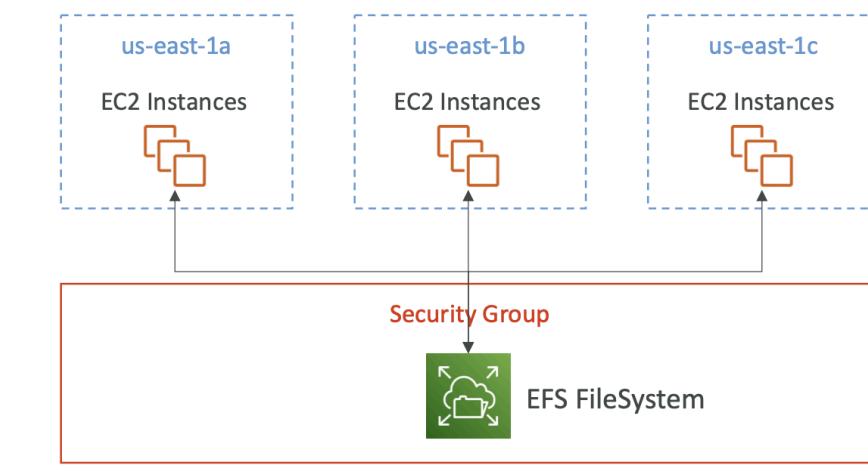


EC2 Instance Store

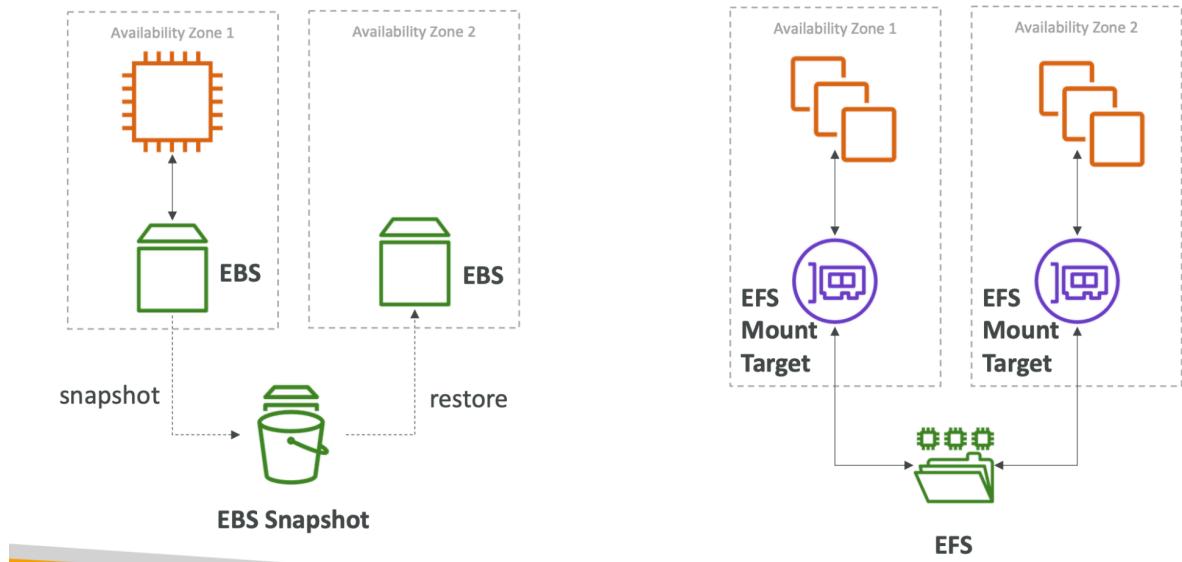
- EBS volumes are network drives with good but “limited” performance
- If you need a high-performance hardware disk, use EC2 Instance Store
- Better I/O performance
- EC2 Instance Store lose their storage if they’re stopped (ephemeral)
- Good for buffer / cache / scratch data / temporary content
- Risk of data loss if hardware fails
- Backups and Replication are your responsibility

EFS – Elastic File System

- Managed NFS (network file system) that can be mounted on 100s of EC2
- EFS works with Linux EC2 instances in multi-AZ
- Highly available, scalable, **expensive** (3x gp2), pay per use, no capacity planning

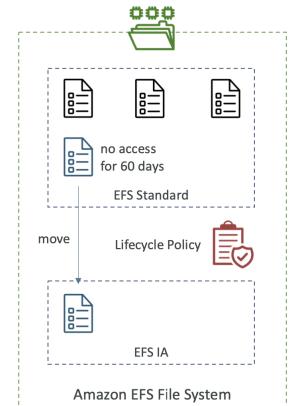


EBS vs EFS



EFS Infrequent Access (EFS-IA)

- Storage class that is cost-optimized for files not accessed every day
- Up to 92% lower cost compared to EFS Standard
- EFS will automatically move your files to EFS-IA based on the last time they were accessed
- Enable EFS-IA with a Lifecycle Policy
- **Example:** move files that are not accessed for 60 days to EFS-IA
- Transparent to the applications accessing EFS



Shared Responsibility Model for EC2 Storage

AWS:

- Infrastructure
- Replication for data for EBS volumes & EFS drives
- Replacing faulty hardware
- Ensuring their employees cannot access your data

User:

- Setting up backup / snapshot procedures
- Setting up data encryption
- Responsibility of any data on the drives
- Understanding the risk of using EC2 Instance Store

Amazon FSx – Overview

- Launch 3rd party high-performance file systems on AWS
- Fully managed service



Amazon FSx for Windows File Server

- A fully managed, highly reliable, and scalable Windows native shared file system
- Built on Windows File Server
- Supports SMB protocol & Windows NTFS
- Integrated with Microsoft Active Directory
- Can be accessed from AWS or your on-premise infrastructure

Amazon FSx for Lustre

- A fully managed, high-performance, scalable file storage for High Performance Computing (HPC)
- The name Lustre is derived from “Linux” and “cluster”
- Machine Learning, Analytics, Video Processing, Financial Modeling,
- Scales up to 100s GB/s, millions of IOPS, sub-ms latencies

EC2 Instance Storage - Summary

- **EBS volumes:**
 - Network drives attached to one EC2 instance at a time
 - Mapped to an Availability Zones
 - Can use EBS Snapshots for backups / transferring EBS volumes across AZ
- **AMI:** create ready-to-use EC2 instances with our customizations
- **EC2 Image Builder :** automatically build, test and distribute AMIs
- **EC2 Instance Store:**
 - High performance hardware disk attached to our EC2 instance
 - Lost if our instance is stopped / terminated
- **EFS:** network file system, can be attached to 100s of instances in a region
- **EFS-IA:** cost-optimized storage class for infrequent accessed files
- **FSx for Windows:** Network File System for Windows servers
- **FSx for Lustre:** High Performance Computing Linux file system

Elastic Load Balancing & Auto Scaling Groups Section

Scalability & High Availability

- Scalability means that an application / system can handle greater loads by adapting.
- There are two kinds of scalability:
 - Vertical Scalability
 - Horizontal Scalability (= elasticity)
- Scalability is linked but different to High Availability

Vertical Scalability

- Vertical Scalability means increasing the size of the instance
- For example, your application runs on a t2.micro
- Scaling that application vertically means running it on a t2.large
- Vertical scalability is very common for non distributed systems, such as a database.
- There's usually a limit to how much you can vertically scale (hardware limit)

Horizontal Scalability

- Horizontal Scalability means increasing the number of instances / systems for your application
- Horizontal scaling implies distributed systems.
- This is very common for web applications / modern applications
- It's easy to horizontally scale thanks the cloud offerings such as Amazon EC2

High Availability

- High Availability usually goes hand in hand with horizontal scaling
- High availability means running your application / system in at least 2 Availability Zones
- The goal of high availability is to survive a data center loss (disaster)

High Availability & Scalability For EC2

- **Vertical Scaling:** Increase instance size (= scale up / down)
 - From: t2.nano - 0.5G of RAM, 1 vCPU
 - To: u-12tb1.metal – 12.3 TB of RAM, 448 vCPUs
- **Horizontal Scaling:** Increase number of instances (= scale out / in)

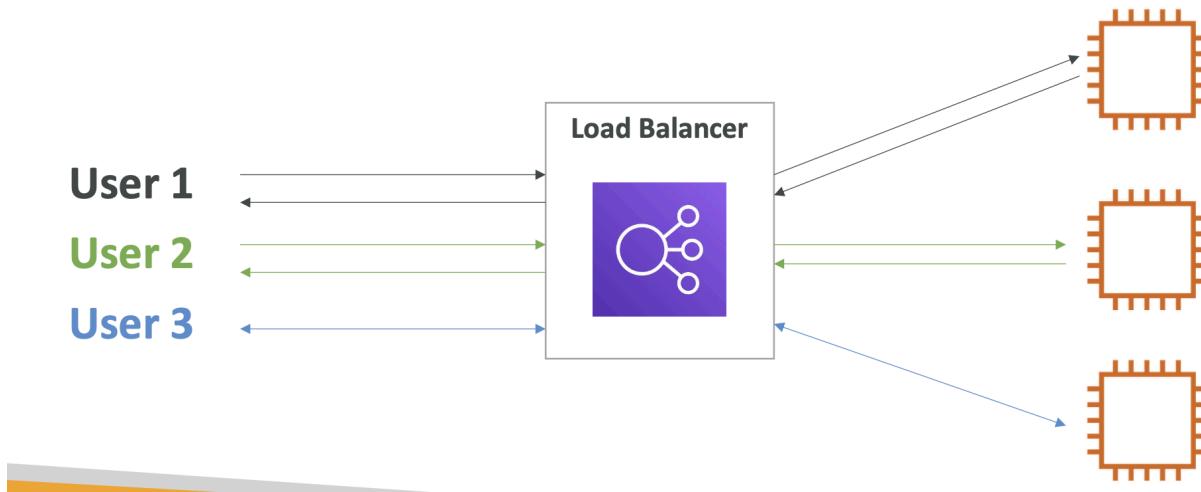
- Auto Scaling Group
- Load Balancer
- **High Availability:** Run instances for the same application across multi AZ
 - Auto Scaling Group multi AZ
 - Load Balancer multi AZ

Scalability vs Elasticity (vs Agility)

- **Scalability:** ability to accommodate a larger load by making the hardware stronger (scale up), or by adding nodes (scale out)
- **Elasticity:** once a system is scalable, elasticity means that there will be some “auto-scaling” so that the system can scale based on the load. This is “cloud-friendly”: pay-per-use, match demand, optimize costs
- **Agility:** (not related to scalability - distractor) new IT resources are only a click away, which means that you reduce the time to make those resources available to your developers from weeks to just minutes.

What is load balancing?

- Load balancers are servers that forward internet traffic to multiple servers (EC2 Instances) downstream.



Why use a load balancer?

- Spread load across multiple downstream instances
- Expose a single point of access (DNS) to your application
- Seamlessly handle failures of downstream instances
- Do regular health checks to your instances
- Provide SSL termination (HTTPS) for your websites
- High availability across zones

Why use an Elastic Load Balancer?

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

What's an Auto Scaling Group?

- In real-life, the load on your websites and application can change
- In the cloud, you can create and get rid of servers very quickly
- The goal of an Auto Scaling Group (ASG) is to:
 - Scale out (add EC2 instances) to match an increased load
 - Scale in (remove EC2 instances) to match a decreased load
 - Ensure we have a minimum and a maximum number of machines running
 - Automatically register new instances to a load balancer
 - Replace unhealthy instances
- Cost Savings: only run at an optimal capacity (principle of the cloud)

Auto Scaling Groups – Scaling Strategies

- **Manual Scaling:** Update the size of an ASG manually
- **Dynamic Scaling:** Respond to changing demand
 - Simple / Step Scaling
 - When a CloudWatch alarm is triggered (example CPU > 70%), then add 2 units
 - When a CloudWatch alarm is triggered (example CPU < 30%), then remove 1
 - Target Tracking Scaling
 - Example: I want the average ASG CPU to stay at around 40%
 - Scheduled Scaling
 - Anticipate a scaling based on known usage patterns
 - Example: increase the min. capacity to 10 at 5 pm on Fridays

ELB & ASG – Summary

- **High Availability vs Scalability** (vertical and horizontal) vs Elasticity vs Agility in the Cloud
- **Elastic Load Balancers (ELB)**

- Distribute traffic across backend EC2 instances, can be Multi-AZ
- Supports health checks
- 4 types: Classic (old), Application (HTTP – L7), Network (TCP – L4), Gateway (L3)
- Auto Scaling Groups (ASG)
 - Implement Elasticity for your application, across multiple AZ
 - Scale EC2 instances based on the demand on your system, replace unhealthy
 - Integrated with the ELB

Amazon S3 Section

Section introduction

- Amazon S3 is one of the main building blocks of AWS
- It's advertised as "infinitely scaling" storage
- Many websites use Amazon S3 as a backbone
- Many AWS services use Amazon S3 as an integration as well

Amazon S3 Use cases

- Backup and storage
- Disaster Recovery
- Archive
- Hybrid Cloud storage
- Application hosting
- Media hosting
- Data lakes & big data analytics
- Software delivery
- Static website

Amazon S3 - Buckets

- Amazon S3 allows people to store objects (files) in "buckets" (directories)
- Buckets must have a globally unique name (across all regions all accounts)
- Buckets are defined at the region level
- S3 looks like a global service but buckets are created in a region
- Naming convention
 - No uppercase, No underscore
 - 3-63 characters long
 - Not an IP

- Must start with lowercase letter or number
- Must NOT start with the prefix xn--
- Must NOT end with the suffix -s3alias

Amazon S3 – Objects (cont.)

- Object values are the content of the body:
 - Max. Object Size is 5TB (5000GB)
 - If uploading more than 5GB, must use “multi-part upload”
- Metadata (list of text key / value pairs – system or user metadata)
- Tags (Unicode key / value pair – up to 10) – useful for security / lifecycle
- Version ID (if versioning is enabled)

Amazon S3 – Security

- User-Based
 - IAM Policies – which API calls should be allowed for a specific user from IAM
- Resource-Based
 - Bucket Policies – bucket wide rules from the S3 console - allows cross account
 - Object Access Control List (ACL) – finer grain (can be disabled)
 - Bucket Access Control List (ACL) – less common (can be disabled)
- Note: an IAM principal can access an S3 object if
 - The user IAM permissions ALLOW it OR the resource policy ALLOWS it
 - AND there's no explicit DENY
- Encryption: encrypt objects in Amazon S3 using encryption keys

S3 Storage Classes

- Amazon S3 Standard - General Purpose
- Amazon S3 Standard-Infrequent Access (IA)
- Amazon S3 One Zone-Infrequent Access
- Amazon S3 Glacier Instant Retrieval
- Amazon S3 Glacier Flexible Retrieval
- Amazon S3 Glacier Deep Archive
- Amazon S3 Intelligent Tiering
- Can move between classes manually or using S3 Lifecycle configurations

No charge any data retrieval **fee**

- Amazon S3 Standard
- Amazon S3 Intelligent-Tiering

S3 Durability and Availability

- **Durability:**

- High durability (99.99999999%, 11 9's) of objects across multiple AZ
- If you store 10,000,000 objects with Amazon S3, you can on average expect to incur a loss of a single object once every 10,000 years
- Same for all storage classes
- **Availability:**
 - Measures how readily available a service is
 - Varies depending on storage class
 - Example: S3 standard has 99.99% availability = not available 53 minutes a year

S3 Standard – General Purpose

- 99.99% Availability
- Used for frequently accessed data
- Low latency and high throughput
- Sustain 2 concurrent facility failures
- Use Cases: Big Data analytics, mobile & gaming applications, content distribution

S3 Storage Classes – Infrequent Access

- For data that is less frequently accessed, but requires rapid access when needed
- Lower cost than S3 Standard
- **Amazon S3 Standard-Infrequent Access (S3 Standard-IA)**
 - 99.9% Availability
 - Use cases: Disaster Recovery, backups
- **Amazon S3 One Zone-Infrequent Access (S3 One Zone-IA)**
 - High durability (99.99999999%) in a single AZ; data lost when AZ is destroyed
 - 99.5% Availability
 - Use Cases: Storing secondary backup copies of on-premise data, or data you can recreate

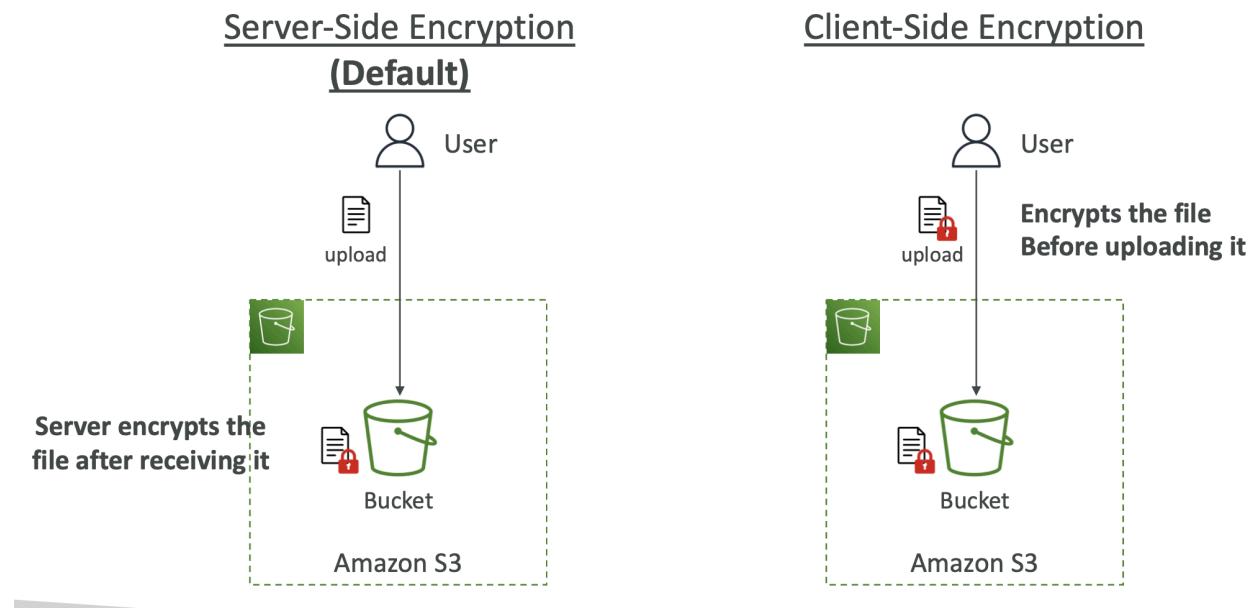
Amazon S3 Glacier Storage Classes

- **Low-cost object storage meant for archiving / backup**
- Pricing: price for storage + object retrieval cost
- **Amazon S3 Glacier Instant Retrieval**
 - Millisecond retrieval, great for data accessed once a quarter
 - Minimum storage duration of 90 days
- **Amazon S3 Glacier Flexible Retrieval (formerly Amazon S3 Glacier):**
 - Expedited (1 to 5 minutes), Standard (3 to 5 hours), Bulk (5 to 12 hours) – free
 - Minimum storage duration of 90 days
- **Amazon S3 Glacier Deep Archive** – for long term storage:
 - Standard (12 hours), Bulk (48 hours)
 - Minimum storage duration of 180 days

S3 Intelligent-Tiering

- Small monthly monitoring and auto-tiering fee
- Moves objects automatically between Access Tiers based on usage
- There are no retrieval charges in S3 Intelligent-Tiering
- Frequent Access tier (automatic): default tier
Infrequent Access tier (automatic): objects not accessed for 30 days
- Archive Instant Access tier (automatic): objects not accessed for 90 days
- Archive Access tier (optional): configurable from 90 days to 700+ days
- Deep Archive Access tier (optional): config. from 180 days to 700+ days

S3 Encryption



IAM Access Analyzer for S3

- Ensures that only intended people have access to your S3 buckets
- Example: publicly accessible bucket, bucket shared with other AWS account
- Evaluates S3 Bucket Policies, S3 ACLs, S3 Access Point Policies
- Powered by IAM Access Analyzer

Shared Responsibility Model for S3

AWS:

- Infrastructure (global security, durability, availability, sustain concurrent loss of data in two facilities)
- Configuration and vulnerability analysis
- Compliance validation

User:

- S3 Versioning
- S3 Bucket Policies
- S3 Replication Setup
- Logging and Monitoring
- S3 Storage Classes
- Data encryption at rest and in transit

S3 Storage Classes Comparison Summary

Storage Class	Durability	Availability	Retrieval Times	Use Cases
S3 Standard	99.999999999 %	99.99%	Milliseconds	Frequently accessed data, dynamic websites, big data analytics
S3 Intelligent-Tiering	99.999999999 %	99.9% (frequent), 99% (infrequent)	Milliseconds	Data with unpredictable or changing access patterns
S3 Standard-IA	99.999999999 %	99.9%	Milliseconds	Infrequently accessed data, backup and disaster recovery, compliance storage
S3 One Zone-IA	99.999999999 %	99.5%	Milliseconds	Infrequently accessed data that can be recreated, secondary backups
S3 Glacier	99.999999999 %	N/A	Minutes to hours	Long-term archival, compliance, backups
S3 Glacier Deep Archive	99.999999999 %	N/A	12 hours or more	Long-term archival, regulatory and compliance archiving, digital preservation
S3 Outposts	99.999999999 %	Varies by Outposts	Milliseconds	On-premises data processing, workloads with data residency requirements, hybrid cloud storage

AWS Snow Family

- Highly-secure, portable devices to collect and process data at the edge, and migrate data into and out of AWS



Snowcone



Snowball Edge



Snowmobile

- Data migration:



Snowcone



Snowball Edge

- Edge computing:

AWS Snow Family: offline devices to perform data migrations

If it takes more than a week to transfer over the network, use Snowball devices!

AWS Snowcone & Snowcone SSD

- Small, portable computing, anywhere, rugged & secure, withstands harsh environments
- Light (4.5 pounds, 2.1 kg)
- Device used for edge computing, storage, and data transfer
- Snowcone – 8 TB of HDD Storage
- Snowcone SSD – 14 TB of SSD Storage
- Use Snowcone where Snowball does not fit (space- constrained environment)
- Must provide your own battery / cables
- Can be sent back to AWS offline, or connect it to internet and use AWS DataSync to send data

Snowball Edge (for data transfers)

- Physical data transport solution: move TBs or PBs of data in or out of AWS
- Alternative to moving data over the network (and paying network fees)
- Pay per data transfer job
- Provide block storage and Amazon S3-compatible object storage
- Snowball Edge Storage Optimized
 - 80 TB of HDD or 210TB NVMe capacity for block volume and S3 compatible object storage
- Snowball Edge Compute Optimized
 - 42 TB of HDD or 28TB NVMe capacity for block volume and S3 compatible object storage
- Use cases: large data cloud migrations, DC decommission, disaster recovery

AWS Snowmobile

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

AWS Snow Family for Data Migrations



	Snowcone & Snowcone SSD	Snowball Edge Storage Optimized	Snowmobile
Storage Capacity	8 TB HDD 14 TB SSD	80 TB - 210 TB	< 100 PB
Migration Size	Up to 24 TB, online and offline	Up to petabytes, offline	Up to exabytes, offline
DataSync agent	Pre-installed		

Snow Family – Usage Process

1. Request Snowball devices from the AWS console for delivery
2. Install the snowball client / AWS OpsHub on your servers
3. Connect the snowball to your servers and copy files using the client
4. Ship back the device when you're done (goes to the right AWS facility)
5. Data will be loaded into an S3 bucket
6. Snowball is completely wiped

What is Edge Computing?

- Process data while it's being created on an edge location
 - A truck on the road,a ship on the sea,a mining station underground...



- These locations may have
 - Limited / no internet access
 - Limited / no easy access to computing power
- We set up a Snowball Edge / Snowcone device to do edge computing
- Use cases of Edge Computing:

- Preprocess data
- Machine learning at the edge
- Transcoding media streams
- Eventually (if need be) we can ship back the device to AWS (for transferring data for example)

Snow Family – Edge Computing

- Snowcone & Snowcone SSD (smaller)
 - 2 CPUs, 4 GB of memory, wired or wireless access
 - USB-C power using a cord or the optional battery
- Snowball Edge – Compute Optimized •
 - 104 vCPUs, 416 GiB of RAM
 - Optional GPU (useful for video processing or machine learning)
 - 28 TB NVMe or 42TB HDD usable storage
 - Storage Clustering available (up to 16 nodes)
- Snowball Edge – Storage Optimized
 - Up to 40 vCPUs, 80 GiB of RAM, 80 TB storage
 - Up to 104 vCPUs, 416 GiB of RAM, 210 TB NVMe storage
- All: Can run EC2 Instances & AWS Lambda functions (using AWS IoT Greengrass)
- Long-term deployment options: 1 and 3 years discounted pricing

AWS OpsHub

- Historically, to use Snow Family devices, you needed a CLI (Command Line Interface tool)
- Today, you can use AWS OpsHub (a software you install on your computer / laptop) to manage your Snow Family Device
 - Unlocking and configuring single or clustered devices
 - Transferring files
 - Launching and managing instances running on Snow Family Devices
 - Monitor device metrics (storage capacity, active instances on your device)
 - Launch compatible AWS services on your devices (ex: Amazon EC2 instances, AWS DataSync, Network File System (NFS))

Snowball Edge Pricing

- You pay for device usage and data transfer out of AWS
- Data transfer IN to Amazon S3 is \$0.00 per GB
- On-Demand
 - Includes a one-time service fee per job, which includes:
 - 10 days of usage for Snowball Edge Storage Optimized 80TB
 - 15 days of usage for Snowball Edge Storage Optimized 210TB
 - Shipping days are NOT counted towards the included 10 or 15 days
 - Pay per day for any additional days

- Committed Upfront
 - Pay in advance for monthly, 1-year, and 3-years of usage (Edge Computing)
 - Up to 62% discounted pricing

Hybrid Cloud for Storage

- AWS is pushing for "hybrid cloud"
 - Part of your infrastructure is on-premises
 - Part of your infrastructure is on the cloud
- This can be due to
 - Long cloud migrations
 - Security requirements
 - Compliance requirements
 - IT strategy
- S3 is a proprietary storage technology (unlike EFS / NFS), so how do you expose the S3 data on-premise?
- AWS Storage Gateway!

AWS Storage Cloud Native Options



AWS Storage Gateway

- Bridge between on-premise data and cloud data in S3
- Hybrid storage service to allow on-premises to seamlessly use the AWS Cloud
- Use cases: disaster recovery, backup & restore, tiered storage
- Types of Storage Gateway:
 - File Gateway
 - Volume Gateway
 - Tape Gateway
- No need to know the types at the exam

Amazon S3 – Summary

- **Buckets vs Objects:** global unique name, tied to a region

- **S3 security:** IAM policy, S3 Bucket Policy (public access), S3 Encryption
- **S3 Websites:** host a static website on Amazon S3
- **S3 Versioning:** multiple versions for files, prevent accidental deletes
- **S3 Replication:** same-region or cross-region, must enable versioning
- **S3 Storage Classes:** Standard, IA, 1Z-IA, Intelligent, Glacier (Instant, Flexible, Deep)
- **Snow Family:** import data onto S3 through a physical device, edge computing
- **OpsHub:** desktop application to manage Snow Family devices
- **Storage Gateway:** hybrid solution to extend on-premises storage to S3

Databases Section

Databases Intro

- Storing data on disk (EFS, EBS, EC2 Instance Store, S3) can have its limits
- Sometimes, you want to store data in a database...
- You can structure the data
- You build indexes to efficiently query / search through the data
- You define relationships between your datasets
- Databases are optimized for a purpose and come with different features, shapes and constraints

NoSQL Databases

- NoSQL = non-SQL = non relational databases
- NoSQL databases are purpose built for specific data models and have flexible schemas for building modern applications.
- **Benefits:**
 - Flexibility: easy to evolve data model
 - Scalability: designed to scale-out by using distributed clusters
 - High-performance: optimized for a specific data model
 - Highly functional: types optimized for the data model
- Examples: Key-value, document, graph, in-memory, search databases

Databases & Shared Responsibility on AWS

- AWS offers use to manage different databases
- Benefits include:
 - Quick Provisioning, High Availability, Vertical and Horizontal Scaling
 - Automated Backup & Restore, Operations, Upgrades
 - Operating System Patching is handled by AWS
 - Monitoring, alerting

- Note: many database technologies could be run on EC2, but you must handle yourself the resiliency, backup, patching, high availability, fault tolerance, scaling...

Amazon RDS Overview

- RDS stands for Relational Database Service
- It's a managed DB service for DB use SQL as a query language.
- It allows you to create databases in the cloud that are managed by AWS
 - Postgres
 - MySQL
 - MariaDB
 - Oracle
 - Microsoft SQL Server
 - IBM DB2
 - Aurora (AWS Proprietary database)

AWS automatically patches RDS instances as it is an AWS managed service.

Advantage over using RDS versus deploying DB on EC2

- RDS is a managed service:
 - Automated provisioning, OS patching
 - Continuous backups and restore to specific timestamp (Point in Time Restore)!
 - Monitoring dashboards
 - Read replicas for improved read performance
 - Multi AZ setup for DR (Disaster Recovery)
 - Maintenance windows for upgrades
 - Scaling capability (vertical and horizontal)
 - Storage backed by EBS
- BUT you can't SSH into your instances

Amazon Aurora

- Aurora is a proprietary technology from AWS (not open sourced)
- PostgreSQL and MySQL are both supported as Aurora DB
- Aurora is “AWS cloud optimized” and claims 5x performance improvement over MySQL on RDS, over 3x the performance of Postgres on RDS
- Aurora storage automatically grows in increments of 10GB, up to 128 TB
- Aurora costs more than RDS (20% more) – but is more efficient
- Not in the free tier

Amazon Aurora Serverless

- Automated database instantiation and auto-scaling based on actual usage
- PostgreSQL and MySQL are both supported as Aurora Serverless DB
- No capacity planning needed
- Least management overhead
- Pay per second, can be more cost- effective
- **Use cases:** good for infrequent, intermittent or unpredictable workloads...

Amazon ElastiCache Overview

- The same way RDS is to get managed Relational Databases...
- ElastiCache is to get managed Redis or Memcached
- **Caches are in-memory databases with high performance, low latency**
- Helps reduce load off databases for read intensive workloads
- AWS takes care of OS maintenance / patching, optimizations, setup, configuration, monitoring, failure recovery and backups

DynamoDB

- **Fully Managed Highly available with replication across 3 AZ**
- NoSQL database - not a relational database
- Scales to massive workloads, distributed “**serverless**” database
- Millions of requests per seconds, trillions of row, 100s of TB of storage
- Fast and consistent in performance
- Single-digit millisecond latency – low latency retrieval
- Integrated with IAM for security, authorization and administration
- Low cost and auto scaling capabilities
- Standard & Infrequent Access (IA) Table Class

DynamoDB Accelerator - DAX

- **Fully Managed in-memory cache for DynamoDB**
- 10x performance improvement – single- digit millisecond latency to microseconds latency – when accessing your DynamoDB tables
- Secure, highly scalable & highly available
- Difference with ElastiCache at the CCP level: DAX is only used for and is integrated with DynamoDB, while ElastiCache can be used for other databases

DynamoDB – Global Tables

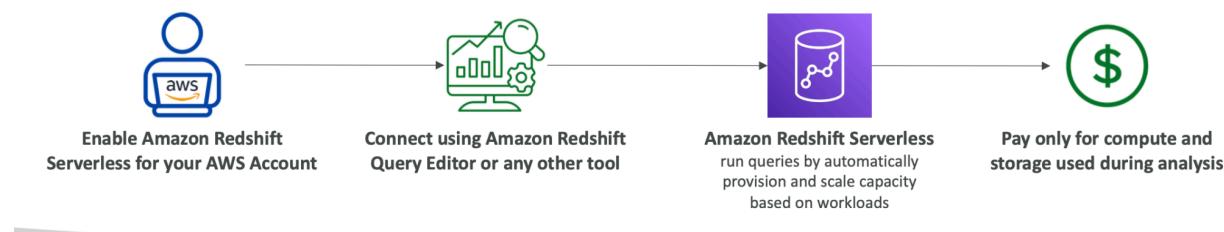
- **Make a DynamoDB table accessible with low latency in multiple-regions**
- **Active-Active replication (read/write to any AWS Region)**

Redshift Overview

- Redshift is based on PostgreSQL, but it's not used for OLTP
- It's OLAP – online analytical processing (analytics and data warehousing)
- Load data once every hour, not every second
- 10x better performance than other data warehouses, scale to PBs of data
- Columnar storage of data (instead of row based)
- Massively Parallel Query Execution (MPP), highly available
- Pay as you go based on the instances provisioned
- Has a SQL interface for performing the queries
- BI tools such as AWS Quicksight or Tableau integrate with it

Redshift Serverless

- Automatically provisions and scales data warehouse underlying capacity
- Run analytics workloads without managing data warehouse infrastructure
- Pay only for what you use (save costs)
- **Use cases:** Reporting, dashboarding applications, real-time analytics...



Amazon EMR

- EMR stands for “Elastic MapReduce”
- **EMR helps creating Hadoop clusters (Big Data) to analyze and process vast amount of data**
- The clusters can be made of hundreds of EC2 instances
- Also supports Apache Spark, HBase, Presto, Flink...
- **EMR takes care of all the provisioning and configuration**
- Auto-scaling and integrated with Spot instances
- **Use cases:** data processing, machine learning, web indexing, big data...

Amazon Athena

- Serverless query service to analyze data stored in Amazon S3
- Uses standard SQL language to query the files
- Supports CSV,JSON,ORC,Avro, and Parquet(bultonPresto)
- Pricing: \$5.00 per TB of data scanned
- Use compressed or columnar data for cost-savings (less scan)
- Use cases: Business intelligence / analytics / reporting, analyze & query VPC Flow Logs, ELB Logs, CloudTrail trails, etc...



- Exam Tip: analyze data in S3 using serverless SQL, use Athena

Amazon QuickSight

- **Serverless machine learning-powered business intelligence service to create interactive dashboards**

- Fast, automatically scalable, embeddable, with per-session pricing
- Use cases:
 - Business analytics
 - Building visualizations
 - Perform ad-hoc analysis
 - Get business insights using data
- Integrated with RDS, Aurora, Athena, Redshift, S3...



DocumentDB

- Aurora is an “AWS-implementation” of PostgreSQL / MySQL ...
- DocumentDB is the same for MongoDB (which is a NoSQL database)
- MongoDB is used to store, query, and index JSON data
- Similar “deployment concepts” as Aurora
- Fully Managed, highly available with replication across 3 AZ
- DocumentDB storage automatically grows in increments of 10GB
- Automatically scales to workloads with millions of requests per seconds

DocumentDB is serverless?

Not a serverless service. It requires you to provision instances (or nodes) and configure clusters with specified instance sizes.

DynamoDB vs DocumentDB

Feature	Amazon DynamoDB	Amazon DocumentDB
Data Model	Key-value, Document	Document (JSON-like)
Compatibility	Proprietary	MongoDB 3.6 API
Use Cases	Low-latency, high-throughput, global apps	MongoDB-based applications
Performance	Single-digit ms response times	Optimized for read and write operations
Scalability	Automatic horizontal scaling	Auto scales storage, compute separately

Transactions	ACID transactions supported	Supports MongoDB transactions
Global Deployment	Global Tables for multi-region replication	Multi-AZ replication
Management	Fully managed, serverless options	Fully managed
Security	Encryption at rest and in transit, IAM	Encryption, VPC isolation, IAM
Backups	Continuous backups, point-in-time recovery	Automated backups, point-in-time recovery

Amazon Neptune

- Fully managed graph database
- A popular graph dataset would be a social network
 - Users have friends
 - Posts have comments
 - Comments have likes from users
 - Users share and like posts...
- Highly available across 3 AZ, with up to 15 read replicas
- Build and run applications working with highly connected datasets – optimized for these complex and hard queries
- Can store up to billions of relations and query the graph with milliseconds latency
- Highly available with replications across multiple AZs
- Great for knowledge graphs (Wikipedia), fraud detection, recommendation engines, social networking

Amazon Timestream

- Fully managed, fast, scalable, **serverless** time series database
- Automatically scales up/down to adjust capacity
- Store and analyze trillions of events per day 1000s times faster & 1/10th the cost of relational databases
- Built-in time series analytics functions (helps you identify patterns in your data in near real-time)

Amazon QLDB

- QLDB stands for "Quantum Ledger Database"
- A ledger is a book recording financial transactions
- FullyManaged, Serverless, Highavailable, Replication across 3AZ
- Used to review history of all the changes made to your application data over time
- Immutable system: no entry can be removed or modified, cryptographically verifiable

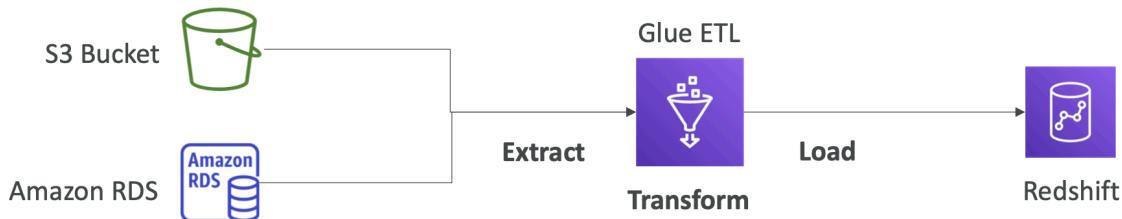
- 2-3x better performance than common ledger blockchain frameworks, manipulate data using SQL
- Difference with Amazon Managed Blockchain: no decentralization component, in accordance with financial regulation rules

Amazon Managed Blockchain

- Blockchain makes it possible to build applications where multiple parties can execute transactions without the need for a trusted, central authority.
- Amazon Managed Blockchain is a managed service to:
 - Join public blockchain networks
 - Or create your own scalable private network
- Compatible with the frameworks **Hyperledger Fabric & Ethereum**

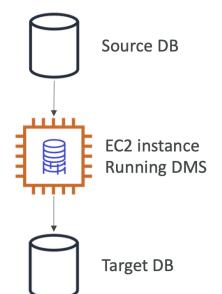
AWS Glue

- Managed extract, transform, and load (ETL) service
- Useful to prepare and transform data for analytics
- Fully serverless service
- Glue Data Catalog: catalog of datasets
 - can be used by Athena, Redshift, EMR



DMS – Database Migration Service

- Quickly and securely migrate databases EC2 instance
Running DMS to AWS, resilient, self healing
- The source database remains available during the migration
- **Supports:**
 - Homogeneous migrations: ex Oracle to Oracle
 - Heterogeneous migrations: ex Microsoft SQL Server to Aurora



Databases & Analytics Summary in AWS

- **Relational Databases** - OLTP: RDS & Aurora (SQL)
- Differences between Multi-AZ, Read Replicas, Multi-Region
- **In-memory Database**: ElastiCache
- **Key/Value Database**: DynamoDB (serverless) & DAX (cache for DynamoDB)

- **Warehouse** - OLAP: Redshift (SQL)
- **Hadoop Cluster**: EMR
- **Athena**: query data on Amazon S3 (serverless & SQL)
- **QuickSight**: dashboards on your data (serverless)
- **DocumentDB**: “Aurora for MongoDB” (JSON–NoSQL database)
- **AmazonQLDB**: Financial Transactions Ledger (immutable journal, cryptographically verifiable)
- Amazon Managed Blockchain: managed Hyperledger Fabric & Ethereum blockchains
- **Glue**: Managed ETL (Extract Transform Load) and Data Catalog service
- **Database Migration**: DMS
- **Neptune**: graph database
- **Timestream**: time-series database

Other Compute Section

What is Docker?

- Docker is a software development platform to deploy apps
- Apps are packaged in containers that can be run on any OS
- Apps run the same, regardless of where they're run
 - Any machine
 - No compatibility issues
 - Predictable behavior
 - Less work
 - Easier to maintain and deploy
 - Works with any language, any OS, any technology
- Scale containers up and down very quickly (seconds)

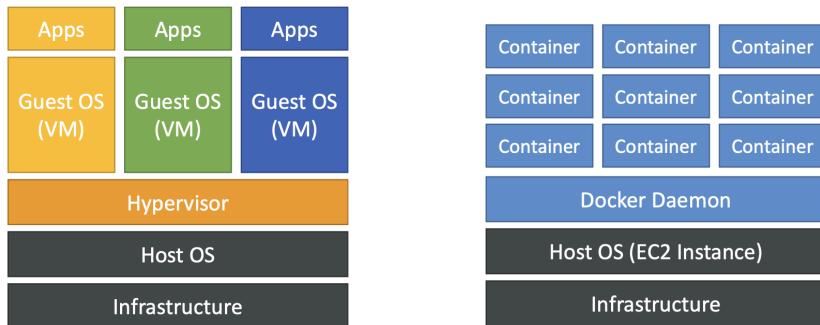
Where Docker images are stored?

- Docker images are stored in Docker Repositories
- Public: Docker Hub <https://hub.docker.com/>
- Find base images for many technologies or OS:
 - Ubuntu
 - MySQL
 - NodeJS, Java...
- Private: Amazon ECR (Elastic Container Registry)

Docker versus Virtual Machines

- Docker is “sort of” a virtualization technology, but not exactly

- Resources are shared with the host => many containers on one server



ECS

Amazon ECS is a highly scalable, high-performance container orchestration service that supports Docker containers. It allows you to run and manage Docker containers on a cluster of EC2 instances.

- **ECS = Elastic Container Service**
- Launch Docker containers on AWS
- You must provision & maintain the infrastructure (the EC2 instances)
- AWS takes care of starting / stopping containers
- Has integrations with the Application Load Balancer

Amazon ECS is best suited for scenarios where you need more control over the underlying infrastructure, such as custom configurations, specific instance types, and hybrid deployments.

Fargate

AWS Fargate is a serverless compute engine for containers that works with both Amazon ECS and Amazon EKS (Elastic Kubernetes Service). With Fargate, **you do not need to manage the underlying EC2 instances**. Fargate manages the infrastructure for you, allowing you to focus on developing and running your applications.

- Launch Docker containers on AWS
- You do not provision the infrastructure (no EC2 instances to manage) – simpler!
- Serverless offering
- AWS just runs containers for you based on the CPU / RAM you need

ECS vs. Fargate

Feature	Amazon ECS	AWS Fargate
---------	------------	-------------

Infrastructure Management	Managed by the user (EC2 instances)	Managed by AWS (serverless)
Control	Full control over EC2 instances	No control over underlying instances
Scaling	Manual or auto-scaling of instances	Automatic scaling of containers
Cost	Pay for EC2 instances	Pay for vCPU and memory resources used
Security	Managed by user with IAM and security groups	Enhanced security with AWS Nitro system
Use Cases	Custom infrastructure, hybrid deployments, legacy applications	Microservices, batch processing, event-driven applications, development/testing environments

ECR

- Elastic Container Registry
- Private Docker Registry on AWS
- This is where you store your Docker images so they can be run by ECS or Fargate

Store, manage, and deploy Docker container images

What's serverless?

- Serverless is a new paradigm in which the developers don't have to manage servers anymore...
- They just deploy code
- They just deploy... functions !
- Initially... Serverless == FaaS (Function as a Service)
- Serverless was pioneered by AWS Lambda but now also includes anything that's managed: "databases, messaging, storage, etc."
- Serverless does not mean there are no servers. It means you just don't manage / provision / see them

Why AWS Lambda

Amazon EC2

- Virtual Servers in the Cloud
- Limited by RAM and CPU

- Continuously running
- Scaling means intervention to add / remove servers

Amazon Lambda

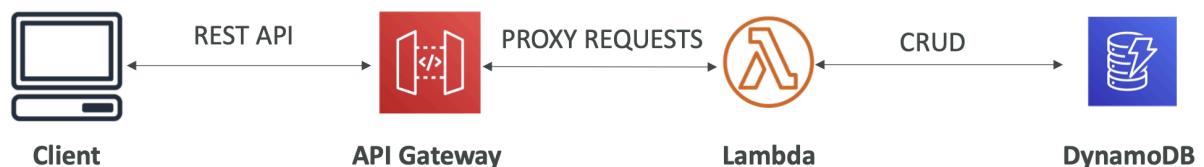
- Virtual functions – no servers to manage!
- Limited by time - short executions
- Run on-demand
- Scaling is automated!

Benefits of AWS Lambda

- Easy Pricing:
 - Pay per request and compute time
 - Free tier of 1,000,000 AWS Lambda requests and 400,000 GBs of compute time
- Integrated with the whole AWS suite of services
- Event-Driven: functions get invoked by AWS when needed
- Integrated with many programming languages
- Easy monitoring through AWS CloudWatch
- Easy to get more resources per function (up to 10GB of RAM!)
- Increasing RAM will also improve CPU and network!

Amazon API Gateway

- Example: building a serverless API



- Fully managed service for developers to easily create, publish, maintain, monitor, and secure APIs
- **Serverless and scalable**
- Supports RESTful APIs and WebSocket APIs
- Supports for security, user authentication, API throttling, API keys, monitoring...

AWS Batch

- Fully managed batch processing at any scale
- Efficiently run 100,000s of computing batch jobs on AWS
- A “batch” job is a job with a start and an end (opposed to continuous)
- Batch will dynamically launch EC2 instances or Spot Instances
- AWS Batch provisions the right amount of compute / memory
- You submit or schedule batch jobs and AWS Batch does the rest!
- Batch jobs are defined as Docker images and run on ECS

- Helpful for cost optimizations and focusing less on the infrastructure

Batch vs Lambda

- **Lambda:**
 - Time limit
 - Limited runtimes
 - Limited temporary disk space
 - Serverless
- **Batch:**
 - No time limit
 - Any runtime as long as it's packaged as a Docker image
 - Rely on EBS / instance store for disk space
 - Relies on EC2 (can be managed by AWS)

Amazon Lightsail

- Virtual servers, storage, databases, and networking
- Low & predictable pricing
- Simpler alternative to using EC2, RDS, ELB, EBS, Route 53...
- Great for people with little cloud experience!
- Can setup notifications and monitoring of your Lightsail resources
- Use cases:
 - Simple web applications (has templates for LAMP, Nginx, MEAN, Node.js...)
 - Websites (templates for WordPress, Magento, Plesk, Joomla)
 - Dev /Test environment
- Has high availability but no auto-scaling, limited AWS integrations

Other Compute - Summary

- **Docker** : container technology to run applications
- **ECS**: run Docker containers on EC2 instances
- **Fargate**:
 - Run Docker containers without provisioning the infrastructure
 - Serverless offering (no EC2 instances)
- **ECR**: Private Docker Images Repository
- **Batch**: run batch jobs on AWS across managed EC2 instances
- **Lightsail**: predictable & low pricing for simple application & DB stacks

Lambda Summary

- Lambda is Serverless, Function as a Service, seamless scaling, reactive
- **Lambda Billing**:
 - By the time run x by the RAM provisioned
 - By the number of invocations

- **Language Support:** many programming languages except (arbitrary) Docker
- **Invocation time:** up to 15 minutes
- **Use cases:**
 - Create Thumbnails for images uploaded onto S3
 - Run a Serverless cron job
- API Gateway: expose Lambda functions as HTTP API

Deploying and Managing Infrastructure at Scale Section

What is CloudFormation

- **CloudFormation** is a declarative way of outlining your AWS Infrastructure, for any resources (most of them are supported).
- For example, within a **CloudFormation template**, you say:
 - I want a security group
 - I want two EC2 instances using this security group
 - I want an S3 bucket
 - I want a load balancer (ELB) in front of these machines
- Then CloudFormation creates those for you, in the right order, with the exact configuration that you specify

Benefits of AWS CloudFormation

- **Infrastructure as code**
 - No resources are manually created, which is excellent for control
 - Changes to the infrastructure are reviewed through code
- **Cost**
 - Each resources within the stack is tagged with an identifier so you can easily see how much a stack costs you
 - You can estimate the costs of your resources using the CloudFormation template
 - Savings strategy: In Dev, you could automation deletion of templates at 5 PM and recreated at 8 AM, safely

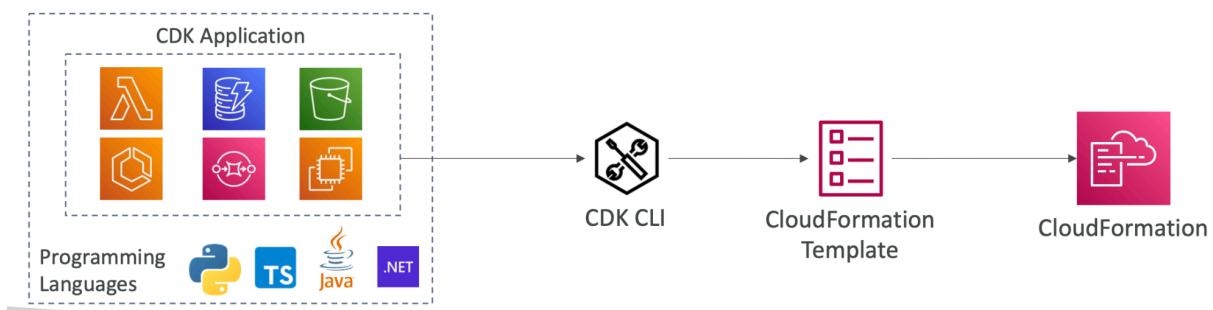
Benefits of AWS CloudFormation

- **Productivity**
 - Ability to destroy and re-create an infrastructure on the cloud on the fly
 - Automated generation of Diagrams for your templates!
 - Declarative programming (no need to figure out ordering and orchestration)
- Don't reinvent the wheel
 - Leverage existing templates on the web!
 - Leverage the documentation

- Supports (almost) all AWS resources:
 - Everything we'll see in this course is supported
 - You can use "custom resources" for resources that are not supported

AWS Cloud Development Kit (CDK)

- Define your cloud infrastructure using a familiar language:
 - JavaScript/TypeScript, Python, Java, and .NET
- The code is "compiled" into a CloudFormation template (JSON/YAML)
- You can therefore deploy infrastructure and application runtime code together
- Great for Lambda functions
- Great for Docker containers in ECS / EKS



AWS Elastic Beanstalk Overview

AWS Elastic Beanstalk is a powerful service for deploying and managing web applications and services. It simplifies the deployment process, automates infrastructure management, and provides built-in monitoring and scaling capabilities. With support for multiple programming languages and platforms, Elastic Beanstalk allows you to focus on developing your applications while AWS handles the operational details. Whether you are building a small web application or a large-scale service, Elastic Beanstalk provides the tools and infrastructure to deploy, manage, and scale your applications effectively.

- Elastic Beanstalk is a developer centric view of deploying an application on AWS
- It uses all the components we've seen before: EC2, ASG, ELB, RDS, etc...
- But it's all in one view that's easy to make sense of!
- We still have full control over the configuration
- **Beanstalk = Platform as a Service (PaaS)**
- Beanstalk is free but you pay for the underlying instances

Elastic Beanstalk

- Managed service
 - Instance configuration / OS is handled by Beanstalk
 - Deployment strategy is configurable but performed by Elastic Beanstalk

- Capacity provisioning
- Load balancing & auto-scaling
- Application health-monitoring & responsiveness
- Just the application code is the responsibility of the developer
- Three architecture models:
 - Single Instance deployment: good for dev
 - LB + ASG: great for production or pre-production web applications
 - ASG only: great for non-web apps in production (workers, etc..)

AWS CodeDeploy

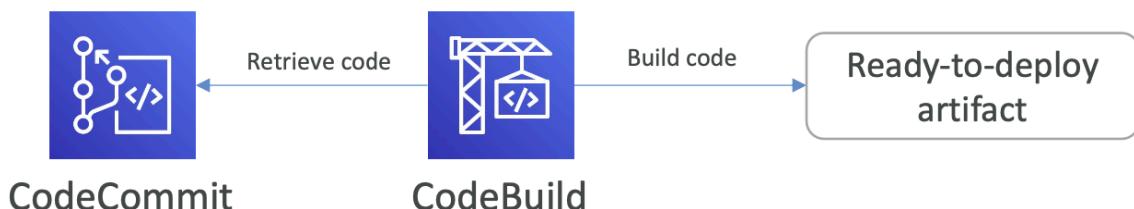
- We want to **deploy our application automatically**
- Works with EC2 Instances
- Works with On-Premises Servers
- **Hybrid service**
- Servers / Instances must be provisioned and configured ahead of time with the CodeDeploy Agent

AWS CodeCommit

- Before pushing the application code to servers, it needs to be stored somewhere
- Developers usually store code in a repository, using the Git technology
- A famous public offering is GitHub, AWS' competing product is CodeCommit
- CodeCommit:
 - Source-control service that hosts Git-based repositories
 - Makes it easy to **collaborate** with others on code
 - The code changes are automatically versioned
- Benefits:
 - Fully managed
 - Scalable & highly available
 - Private, Secured, Integrated with AWS

AWS CodeBuild

- Code building service in the cloud (name is obvious)
- Compiles source code, run tests, and produces packages that are ready to be deployed (by CodeDeploy for example)

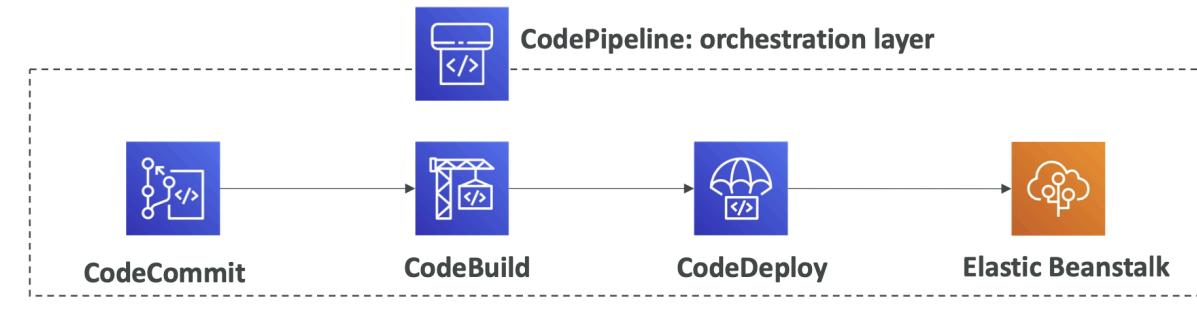


- **Benefits**

- Fully managed, serverless
- Continuously scalable & highly available
- Secure
- Pay-as-you-go pricing – only pay for the build time

AWS CodePipeline

- Orchestrate the different steps to have the code automatically pushed to production
 - **Code=>Build=>Test=>Provision=>Deploy**
 - Basis for CICD (Continuous Integration & Continuous Delivery)
- Benefits:
 - Fully managed, compatible with CodeCommit, CodeBuild, CodeDeploy, Elastic Beanstalk,
 - CloudFormation, GitHub, 3rd-party services (GitHub...) & custom plugins...
 - Fast delivery & rapid updates



AWS CodeArtifact

- Software packages depend on each other to be built (also called code dependencies), and new ones are created
- **Storing and retrieving these dependencies is called artifact management**
- Traditionally you need to setup your own artifact management system
- CodeArtifact is a secure, scalable, and cost-effective artifact management for software development
- **Works with common dependency management tools such as Maven, Gradle, npm, yarn, twine, pip, and NuGet**
- Developers and CodeBuild can then retrieve dependencies straight from CodeArtifact

AWS CodeStar

- Unified UI to easily manage software development activities in one place
- "Quick way" to get started to correctly set-up CodeCommit, CodePipeline, CodeBuild, CodeDeploy, Elastic Beanstalk, EC2, etc...
- Can edit the code "in-the-cloud" using AWS Cloud9

AWS Cloud9

- AWS Cloud9 is a cloud IDE (Integrated Development Environment) for writing, running and debugging code
- "Classic" IDE (like IntelliJ, Visual Studio Code...) are downloaded on a computer before being used
- A cloud IDE can be used within a web browser, meaning you can work on your projects from your office, home, or anywhere with internet with no setup necessary
- AWS Cloud9 also allows for code collaboration in real-time (pair programming)

AWS Systems Manager (SSM)

AWS Systems Manager is a versatile and powerful service that simplifies the management and automation of your AWS and on-premises resources. By providing tools for automation, configuration management, operational insights, and secure access, it helps organizations improve operational efficiency, maintain compliance, and enhance security.

- Helps you manage your EC2 and On-Premises systems at scale
- **Another Hybrid AWS service**
- Get operational insights about the state of your infrastructure
- Suite of 10+ products
- Most important features are:
 - Patching automation for enhanced compliance
 - Run commands across an entire fleet of servers
 - Store parameter configuration with the SSM Parameter Store
- Works for Linux, Windows, MacOS, and Raspberry Pi OS (Raspbian)

Deployment - Summary

- **CloudFormation:** (AWS only)
 - Infrastructure as Code, works with almost all of AWS resources
 - Repeat across Regions & Accounts
- **Beanstalk:** (AWS only)
 - Platform as a Service (PaaS), limited to certain programming languages or Docker
 - Deploy code consistently with a known architecture: ex, ALB + EC2 + RDS
- **CodeDeploy (hybrid):** deploy & upgrade any application onto servers
- **Systems Manager (hybrid):** patch, configure and run commands at scale

Developer Services - Summary

- **CodeCommit:** Store code in private git repository (version controlled)
- **CodeBuild:** Build & test code in AWS
- **CodeDeploy:** Deploy code onto servers
- **CodePipeline:** Orchestration of pipeline (from code to build to deploy)
- **CodeArtifact:** Store software packages / dependencies on AWS
- **CodeStar:** Unified view for allowing developers to do CICD and code
- **Cloud9:** Cloud IDE (Integrated Development Environment) with collab
- **AWS CDK:** Define your cloud infrastructure using a programming language

Global Infrastructure Section

Global AWS Infrastructure

- **Regions:** For deploying applications and infrastructure
- **Availability Zones:** Made of multiple data centers
- **Edge Locations (Points of Presence):** for content delivery as close as possible to users

Global Applications in AWS

- **Global DNS: Route 53**
 - Great to route users to the closest deployment with least latency
 - Great for disaster recovery strategies
- **Global Content Delivery Network (CDN): CloudFront**
 - Replicate part of your application to AWS Edge Locations – decrease latency
 - Cache common requests – improved user experience and decreased latency
- **S3 Transfer Acceleration**
- Accelerate global uploads & downloads into Amazon S3
- **AWS Global Accelerator:**
 - Improve global application availability and performance using the AWS global network

Amazon Route 53 Overview

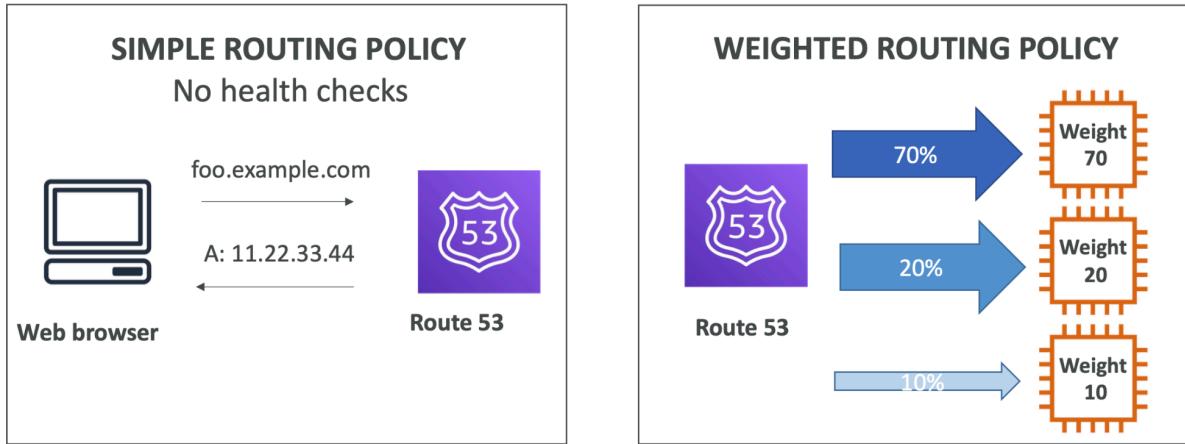
- Route 53 is a Managed DNS (Domain Name System)
- DNS is a collection of rules and records which helps clients understand how to reach a server through URLs.

Route 53 also supports hybrid cloud architectures

Amazon Route 53 can monitor the **health and performance** of your application as well as your web servers and other resources.

Route 53 Routing Policies

- Need to know them at a high-level for the Cloud Practitioner Exam



Route 53 Routing Policies

LATENCY ROUTING POLICY
FAILOVER ROUTING POLICY

Amazon CloudFront

- Content Delivery Network (CDN)
- Improves read performance, content is cached at the edge
- Improves users experience
- 216 Point of Presence globally (edge locations)
- DDoS protection (because worldwide), integration with Shield, AWS Web Application Firewall

CloudFront – Origins

- S3 bucket
 - For distributing files and caching them at the edge
 - Enhanced security with CloudFront Origin Access Control (OAC)
 - OAC is replacing Origin Access Identity (OAI)
 - CloudFront can be used as an ingress (to upload files to S3)
- Custom Origin (HTTP)

- Application Load Balancer
- EC2 instance
- S3 website (must first enable the bucket as a static S3 website)
- Any HTTP backend you want

S3 Transfer Acceleration

- Increase transfer speed by transferring file to an AWS edge location which will forward the data to the S3 bucket in the target region

AWS Global Accelerator

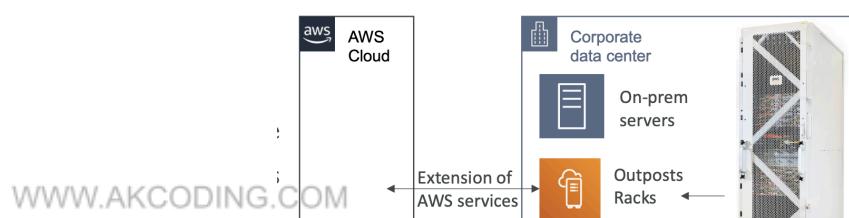
- Improve global application availability and performance using the AWS global network
- Leverage the AWS internal network to optimize the route to your application (60% improvement)
- 2 Anycast IP are created for your application and traffic is sent through Edge Locations
- The Edge locations send the traffic to your application

AWS Global Accelerator vs CloudFront

- They both use the AWS global network and its edge locations around the world
- Both services integrate with AWS Shield for **DDoS** protection.
- CloudFront – Content Delivery Network
 - Improves performance for your cacheable content (such as images and videos)
 - Content is served at the edge
- Global Accelerator
 - No caching, proxying packets at the edge to applications running in one or more AWS Regions.
 - Improves performance for a wide range of applications over TCP or UDP
 - **Good for HTTP use cases that require static IP addresses**
 - Good for HTTP use cases that required deterministic, fast regional failover

AWS Outposts

- **Hybrid Cloud:** businesses that keep an on-premises infrastructure alongside a cloud infrastructure
- Therefore, two ways of dealing with IT systems:
 - One for the AWS cloud (using the AWS console, CLI, and AWS APIs)
 - One for their on-premises infrastructure
- **AWS Outposts** are “server racks” that offers the same AWS infrastructure, services, APIs & tools to build your own applications on-premises just as in the cloud
- **AWS will setup and manage “Outposts Racks” within your on-premises infrastructure and you can start leveraging AWS services on-premises**



- You are responsible for the Outposts Rack physical security

Benefits:

- Low-latency access to on-premises systems
- Local data processing
- Data residency
- Easier migration from on-premises to the cloud
- Fully managed service

Some services that work on Outposts:



Amazon EC2



Amazon EBS



Amazon S3



Amazon EKS



Amazon ECS



Amazon RDS



Amazon EMR

AWS WaveLength

- WaveLength Zones are infrastructure deployments embedded within the telecommunications providers' datacenters at the edge of the 5G networks
- Brings AWS services to the edge of the 5G networks
- Example:EC2,EBS,VPC...
- Ultra-low latency applications through 5G networks
- Traffic doesn't leave the Communication Service Provider's (CSP) network
- High-bandwidth and secure connection to the parent AWS Region
- No additional charges or service agreements
- Use cases: Smart Cities, ML-assisted diagnostics, Connected Vehicles, Interactive Live Video Streams, AR/VR, Real-time Gaming, ...

AWS Local Zones

- Places AWS compute, storage, database, and other selected AWS services closer to end users to run latency-sensitive applications
- Extend your VPC to more locations – “Extension of an AWS Region”
- Compatible with EC2, RDS, ECS, EBS, ElastiCache, Direct Connect ...
- Example:
 - AWS Region: N.Virginia (us-east-1)
 - AWS Local Zones: Boston, Chicago, Dallas, Houston, Miami, ...

Global Applications in AWS - Summary

- **Global DNS: Route 53**
 - Great to route users to the closest deployment with least latency

- Great for disaster recovery strategies
- **Global Content Delivery Network (CDN): CloudFront**
 - Replicate part of your application to AWS Edge Locations – decrease latency
 - Cache common requests – improved user experience and decreased latency
- **S3 Transfer Acceleration**
 - Accelerate global uploads & downloads into Amazon S3
- **AWS Global Accelerator**
 - Improve global application availability and performance using the AWS global network
- **AWS Outposts**
 - Deploy Outposts Racks in your own Data Centers to extend AWS services
- **AWS WaveLength**
 - Brings AWS services to the edge of the 5G networks
 - Ultra-low latency applications
- **AWS Local Zones**
 - Bring AWS resources (compute, database, storage, ...) closer to your users
 - Good for latency-sensitive applications

Cloud Integration Section

Section Introduction

- When we start deploying multiple applications, they will inevitably need to communicate with one another
- There are two patterns of application communication

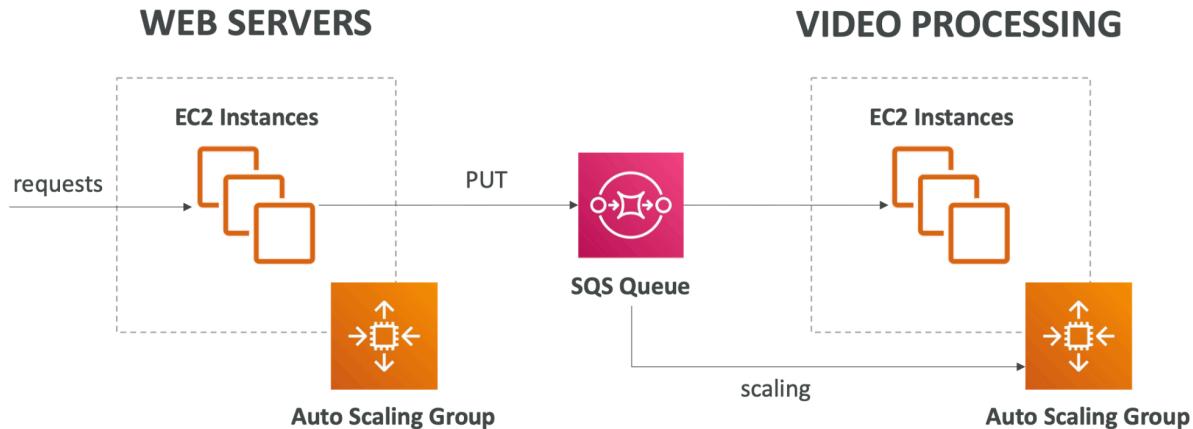


- Synchronous between applications can be problematic if there are sudden spikes of traffic
- What if you need to suddenly encode 1000 videos but usually it's 10?
- In that case, it's better to **decouple** your applications:
 - using SQS: queue model
 - using SNS: pub/sub model
 - using Kinesis: real-time data streaming model
- These services can scale independently from our application!

Amazon SQS – Standard Queue

- Oldest AWS offering (over 10 years old)
- Fully managed service (~serverless), use to decouple applications
- Scales from 1 message per second to 10,000s per second
- Default retention of messages: 4 days, maximum of 14 days
- No limit to how many messages can be in the queue
- Messages are deleted after they're read by consumers
- Low latency (<10 ms on publish and receive)
- Consumers share the work to read messages & scale horizontally

SQS to decouple between application tiers



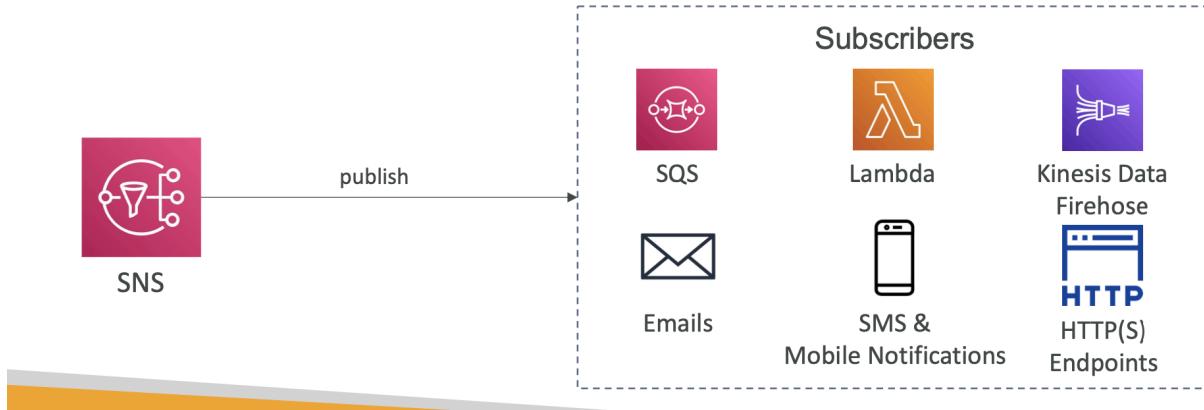
Amazon Kinesis

- For the exam: Kinesis = real-time big data streaming
- Managed service to collect, process, and analyze real-time streaming data at any scale
- Too detailed for the Cloud Practitioner exam but good to know:
 - **Kinesis Data Streams:** low latency streaming to ingest data at scale from hundreds of thousands of sources
 - **Kinesis Data Firehose:** load streams into S3, Redshift, ElasticSearch, etc...
 - **Kinesis Data Analytics:** perform real-time analytics on streams using SQL
 - **Kinesis Video Streams:** monitor real-time video streams for analytics or ML

Amazon SNS

- What if you want to send one message to many receivers?
- The “event publishers” only sends message to one SNS topic

- As many “event subscribers” as we want to listen to the SNS topic notifications • Each subscriber to the topic will get all the messages
- Up to 12,500,000 subscriptions per topic, 100,000 topics limit



Amazon MQ

- SQS, SNS are “cloud-native” services: proprietary protocols from AWS
- Traditional applications running from on-premises may use open protocols such as: MQTT, AMQP, STOMP, Openwire, WSS
- When migrating to the cloud, instead of re-engineering the application to use SQS and SNS, we can use Amazon MQ
- Amazon MQ is a managed message broker service for



- Amazon MQ doesn’t “scale” as much as SQS / SNS
- Amazon MQ runs on servers, can run in Multi-AZ with failover
- Amazon MQ has both queue feature (~SQS) and topic features (~SNS)

Integration Section – Summary

- **SQS:**
 - Queue service in AWS
 - Multiple Producers, messages are kept up to 14 days
 - Multiple Consumers share the read and delete messages when done • Used to decouple applications in AWS
- **SNS:**
 - Notification service in AWS
 - Subscribers: Email, Lambda, SQS, HTTP, Mobile...
 - Multiple Subscribers, send all messages to all of them
 - No message retention
- **Kinesis:** real-time data streaming, persistence and analysis

- **Amazon MQ:** managed message broker for ActiveMQ and RabbitMQ in the cloud (MQTT, AMQP.. protocols)

Cloud Monitoring Section

Amazon CloudWatch Metrics

- CloudWatch provides metrics for every services in AWS
- Metric is a variable to monitor (CPU Utilization, NetworkIn...)
- Metrics have timestamps
- Can create CloudWatch dashboards of metrics

Important Metrics

- **EC2 instances:** CPU Utilization, Status Checks, Network (not RAM)
 - Default metrics every 5 minutes
 - Option for Detailed Monitoring (\$\$\$): metrics every 1 minute
- **EBS volumes:** Disk Read/Writes
- **S3 buckets:** Bucket Size Bytes, Number Of Objects, All Requests
- **Billing:** Total Estimated Charge (only in us-east-1)
- **Service Limits:** how much you've been using a service API
- **Custom metrics:** push your own metrics

Amazon CloudWatch Alarms

- Alarms are used to trigger notifications for any metric
- Alarms actions
 - **Auto Scaling:** increase or decrease EC2 instances “desired” count
 - **EC2 Actions:** stop, terminate, reboot or recover an EC2 instance
 - **SNS notifications:** send a notification into an SNS topic
- Various options (sampling, %, max, min, etc...)
- Can choose the period on which to evaluate an alarm
- Example: create a billing alarm on the CloudWatch Billing metric
- **Alarm States:** OK, INSUFFICIENT_DATA, ALARM

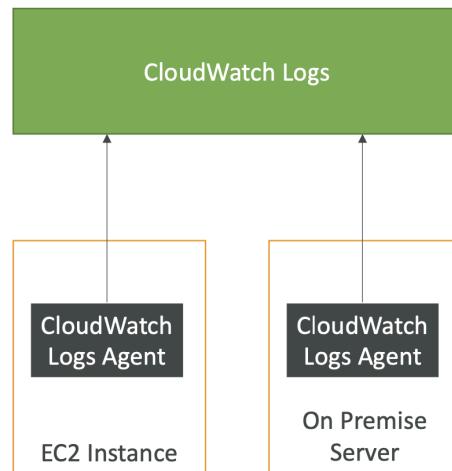
Amazon CloudWatch Logs

- CloudWatch Logs can collect log from:
 - Elastic Beanstalk: collection of logs from application
 - ECS: collection from containers
 - AWS Lambda: collection from function logs
 - CloudTrail based on filter

- **CloudWatch log agents: on EC2 machines or on-premises servers**
- Route 53: Log DNS queries
- Enables real-time monitoring of logs
- Adjustable CloudWatch Logs retention

CloudWatch Logs for EC2

- By default, no logs from your EC2 instance will go to CloudWatch
- You need to run a CloudWatch agent on EC2 to push the log files you want
- Make sure IAM permissions are correct
- The CloudWatch log agent can be setup on-premises too



You can also use **CloudWatch** in hybrid cloud architectures by using the **CloudWatch Agent** or API to monitor your on-premises resources

Amazon EventBridge (formerly CloudWatch Events)

- Schedule: Cron jobs (scheduled scripts)
- Event Pattern: Event rules to react to a service doing something
- Trigger Lambda functions, send SQS/SNS messages...

AWS CloudTrail

- Provides governance, compliance and audit for your AWS Account
- **CloudTrail is enabled by default!**
- Get an history of events / API calls made within your AWS Account by:
 - Console
 - SDK
 - CLI
 - AWS Services

- Can put logs from CloudTrail into CloudWatch Logs or S3
- A trail can be applied to All Regions (default) or a single Region.
- If a resource is deleted in AWS, investigate CloudTrail first!

CloudTrail Diagram



AWS X-Ray

- Debugging in the good old way:
 - Test locally
 - Add log statements everywhere
 - Re-deploy in production
 - Log formats differ across applications and log analysis is hard.
 - Debugging: one big monolith “easy”, distributed services “hard”
 - No common views of your entire architecture
 - Enter... AWS X-Ray!
- Production,

AWS X-Ray advantages

- Troubleshooting performance (bottlenecks)
- Understand dependencies in a microservice architecture
- Pinpoint service issues
- Review request behavior
- Find errors and exceptions
- Are we meeting time SLA?
- Where am I throttled?
- Identify users that are impacted

Amazon CodeGuru

- An ML-powered service for automated code reviews and application performance recommendations
- Provides two functionalities
 - CodeGuru Reviewer: automated code reviews for static code analysis (development)
 - CodeGuru Profiler: visibility/recommendations about application performance during runtime (production)

Amazon CodeGuru Reviewer

- Identify critical issues, security vulnerabilities, and hard-to-find bugs

- Example: common coding best practices, resource leaks, security detection, input validation
- Uses Machine Learning and automated reasoning
- Hard-learned lessons across millions of code reviews on 1000s of open-source and Amazon repositories
- Supports Java and Python
- Integrates with GitHub, Bitbucket, and AWS CodeCommit

Amazon CodeGuru Profiler

- Helps understand the runtime behavior of your application
- Example: identify if your application is consuming excessive CPU capacity on a logging routine
- Features:
 - Identify and remove code inefficiencies
 - Improve Application Performance(e.g.,reduceCPU utilization)
 - Decrease compute costs
 - Provides heap summary (identify which objects using up memory)
 - Anomaly Detection
- Support applications running on AWS or on-premise
- Minimal overhead on application

AWS Health Dashboard - Service History

- Shows all regions, all services health
- Shows historical information for each day
- Has an RSS feed you can subscribe to
- Previously called AWS Service Health Dashboard

AWS Health Dashboard – Your Account

- Previously called AWS Personal Health Dashboard (PHD)
- AWS Account Health Dashboard provides alerts and remediation guidance when AWS is experiencing events that may impact you.
- While the Service Health Dashboard displays the general status of AWS services, Account Health Dashboard gives you a personalized view into the performance and availability of the AWS services underlying your AWS resources.
- The dashboard displays relevant and timely information to help you manage events in progress and provides proactive notification to help you plan for scheduled activities.
- Can aggregate data from an entire AWS Organization

Monitoring Summary

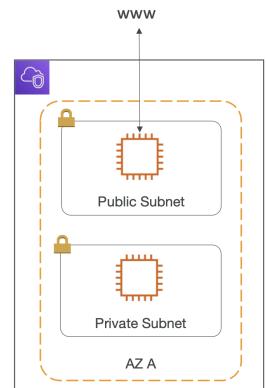
- **CloudWatch:**
 - **Metrics:** monitor the performance of AWS services and billing metrics

- **Alarms:** automate notification, perform EC2 action, notify to SNS based on metric
- **Logs:** collect log files from EC2 instances, servers, Lambda Functions...
- **Events** (or EventBridge): react to events in AWS, or trigger a rule on a schedule
- **CloudTrail:** audit API calls made within your AWS account
- **CloudTrail Insights:** automated analysis of your CloudTrail Events
- **X-Ray:** trace requests made through your distributed applications
- **AWS Health Dashboard:** status of all AWS services across all regions
- **AWS Account Health Dashboard:** AWS events that impact your infrastructure
- **Amazon CodeGuru:** automated code reviews and application performance recommendations

VPC Section

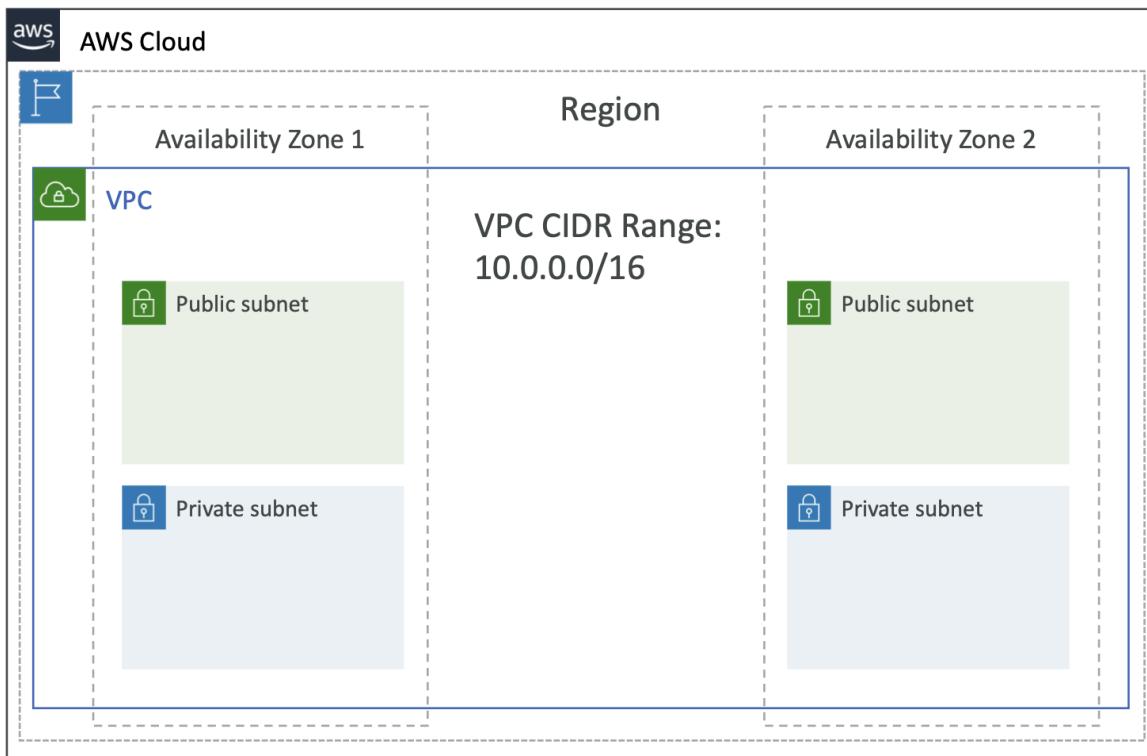
VPC & Subnets Primer

- **VPC -Virtual Private Cloud:** private network to deploy your resources (regional resource)
- **Subnets** allow you to partition your network inside your VPC (Availability Zone resource)
- A **public subnet** is a subnet that is accessible from the internet
- A **private subnet** is a subnet that is not accessible from the internet
- To define access to the internet and between subnets, we use **Route Tables**.



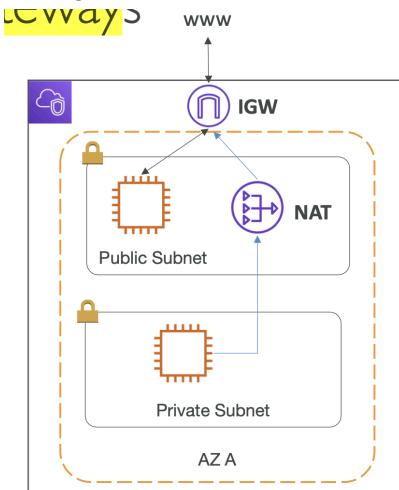
An Amazon Virtual Private Cloud (Amazon VPC) spans all of the **Availability Zones (AZ)** in the Region whereas a **subnet** spans only **one Availability Zone (AZ)** in the Region

VPC Diagram



Internet Gateway & NAT Gateways

- Internet Gateways helps our VPC instances connect with the internet
- Public Subnets have a route to the internet gateway.
- **NAT Gateways (AWS-managed) & NAT Instances (self-managed) allow your instances in your Private Subnets to access the internet while remaining private**

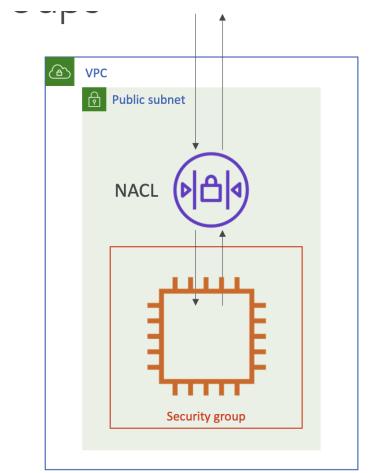


Network ACL & Security Groups

- NACL (Network ACL)
 - A firewall which controls traffic from and to subnet
 - Can have ALLOW and DENY rules
 - Are attached at the Subnet level
 - Rules only include IP addresses

The NACL that is created by default includes an outbound and inbound rule to allow all traffic to and from the subnet.

- Security Groups
 - A firewall that controls traffic to and from an EC2 Instance
 - Can have only ALLOW rules
 - Rules include IP addresses and other security groups

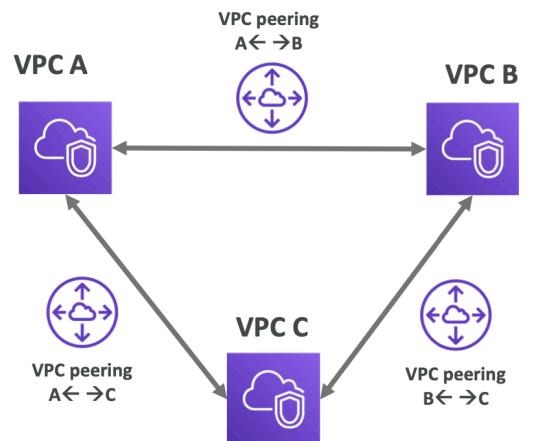


VPC Flow Logs

- Capture information about IP traffic going into your interfaces:
 - VPC Flow Logs
 - Subnet Flow Logs
 - Elastic Network Interface Flow Logs
- Helps to monitor & troubleshoot connectivity issues. Example:
 - Subnets to internet
 - Subnets to subnets
 - Internet to subnets
- Captures network information from AWS managed interfaces too: Elastic Load Balancers, ElastiCache, RDS, Aurora, etc...
- VPC Flow logs data can go to S3, CloudWatch Logs, and Kinesis Data Firehose

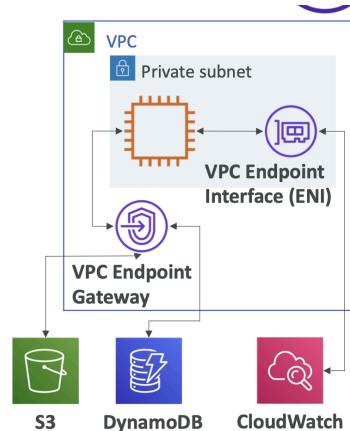
VPC Peering

- Connect two VPC, privately using AWS' network
- Make them behave as if they were in the same network
- Must not have overlapping CIDR (IP address range)
- VPC Peering connection is not transitive (must be established for each VPC that need to communicate with one another)



VPC Endpoints

- Endpoints allow you to connect to AWS Services using a private network instead of the public www network
- This gives you enhanced security and lower latency to access AWS services
- VPC Endpoint Gateway: S3 & DynamoDB**
- VPC Endpoint Interface: the rest

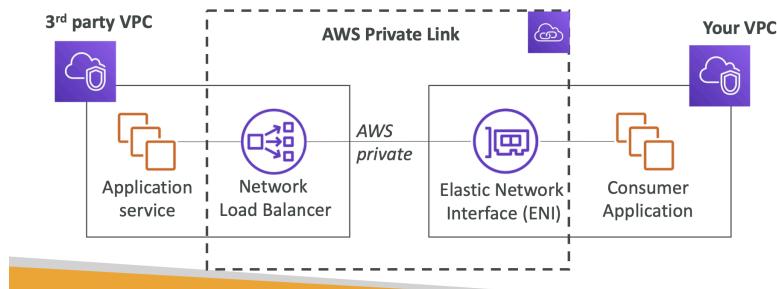


Just remember that only **Amazon S3** and **Amazon DynamoDB** support VPC gateway endpoint. All other services that support VPC Endpoints use a **VPC interface endpoint** (note that Amazon S3 supports the VPC interface endpoint as well).

AWS PrivateLink (VPC Endpoint Services)

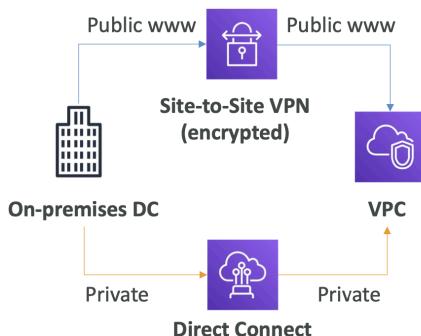
AWS PrivateLink is a service that enables you to securely access services hosted on AWS in a highly available and scalable manner, without using public IPs and without exposing traffic to the public internet. PrivateLink simplifies the security of data shared with cloud-based applications by eliminating the exposure of data to the public internet.

- Most secure & scalable way to expose a service to 1000s of VPCs
- Does not require VPC peering, internet gateway, NAT, route tables...
- Requires a network load balancer (Service VPC) and ENI (Customer VPC)



Site to Site VPN & Direct Connect

- Site to Site VPN
 - Connect an on-premises VPN to AWS
 - The connection is automatically encrypted
 - Goes over the public internet
- Direct Connect (DX)
 - Establish a physical connection between on-premises and AWS
 - The connection is private, secure and fast
 - Goes over a private network
 - Takes at least a month to establish



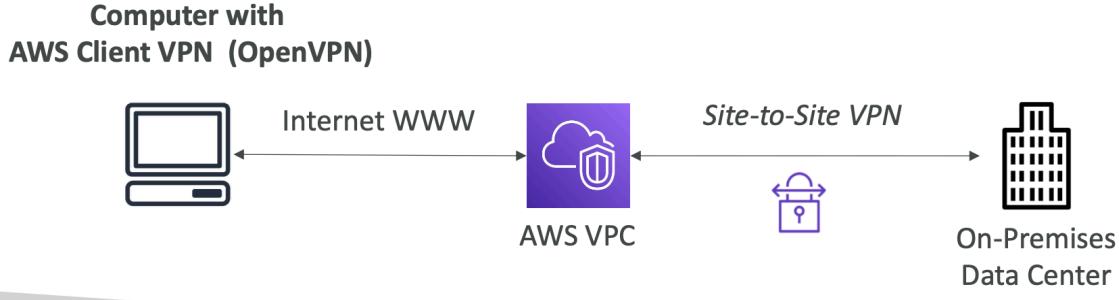
Site-to-Site VPN

- **On-premises:** must use a Customer Gateway (CGW)
- **AWS:** must use a Virtual Private Gateway (VGW)



AWS Client VPN

- Connect from your computer using OpenVPN to your private network in AWS and on-premises
- Allow you to connect to your EC2 instances over a private IP (just as if you were in the private VPC network)
- **Goes over public Internet**



Transit Gateway

- For having transitive peering between thousands of VPC and on-premises, hub-and-spoke (star) connection
- One single Gateway to provide this functionality
- Works with Direct Connect Gateway, VPN connections

VPC Closing Comments

- **VPC** – Virtual Private Cloud
- **Subnets** – Tied to an AZ, network partition of the VPC
- **Internet Gateway** – at the VPC level, provide Internet Access
- **NAT Gateway / Instances** – give internet access to private subnets
- **NAACL** – Stateless, subnet rules for inbound and outbound
- **Security Groups** – Stateful, operate at the EC2 instance level or ENI
- **VPC Peering** – Connect two VPC with non overlapping IP ranges, nontransitive
- **Elastic IP** – fixed public IPv4, ongoing cost if not in-use
- **VPC Endpoints** – Provide private access to AWS Services within VPC
- **PrivateLink** – Privately connect to a service in a 3rd party VPC
- **VPC Flow Logs** – network traffic logs
- **Site to Site VPN** – VPN over public internet between on-premises DC and AWS
- **Client VPN** – OpenVPN connection from your computer into your VPC
- **Direct Connect** – direct private connection to AWS
- **Transit Gateway** – Connect thousands of VPC and on-premises networks together

Security & Compliance Section

- AWS Shared Responsibility Model

- AWS responsibility - Security of the Cloud
- Protecting infrastructure (hardware, software, facilities, and networking) that runs all the AWS services
- Managed services like S3, DynamoDB, RDS, etc.
- Customer responsibility - Security in the Cloud
 - For EC2 instance, customer is responsible for management of the guest OS (including security patches and updates), firewall & network configuration, IAM
 - Encrypting application data
- Shared controls:
 - Patch Management, Configuration Management, Awareness & Training

Example, for RDS

- AWS responsibility:
 - Manage the underlying EC2 instance, disable SSH access
 - Automated DB patching
 - Automated OS patching
 - Audit the underlying instance and disks & guarantee it functions
- Your responsibility:
 - Check the ports / IP / security group inbound rules in DB's SG
 - In-database user creation and permissions
 - Creating a database with or without public access
 - Ensure parameter groups or DB is configured to only allow SSL connections
 - Database encryption setting

Example, for S3

- AWS responsibility:
 - Guarantee you get unlimited storage
 - Guarantee you get encryption
 - Ensure separation of the data between different customers
 - Ensure AWS employees can't access your data
- Your responsibility:
 - Bucket configuration
 - Bucket policy / public setting
 - IAM user and roles
 - Enabling encryption

DDOS Protection on AWS

- **AWS Shield Standard:** protects against DDOS attack for your website and applications, for all customers at no additional costs
- **AWS Shield Advanced:** 24/7 premium DDoS protection
- **AWS WAF:** Filter specific requests based on rules

- **CloudFront and Route 53:**
 - Availability protection using global edge network
 - Combined with AWS Shield, provides attack mitigation at the edge
 - Be ready to scale – leverage AWS Auto Scaling

AWS Shield

- **AWS Shield Standard:**
 - Free service that is activated for every AWS customer
 - Provides protection from attacks such as SYN/UDP Floods, Reflection attacks and other layer 3/layer 4 attacks
- **AWS Shield Advanced:**
 - Optional DDoS mitigation service (\$3,000 per month per organization)
 - Protect against more sophisticated attack on Amazon EC2, Elastic Load Balancing (ELB), Amazon CloudFront, AWS Global Accelerator, and Route 53
 - **24/7 access to AWS DDoS response team (DRP)**
 - Protect against higher fees during usage spikes due to DDoS

AWS WAF – Web Application Firewall

- **Protects your web applications from common web exploits (Layer 7)**
- Layer 7 is HTTP (vs Layer 4 is TCP)
- Deploy on Application Load Balancer, API Gateway, CloudFront
- Define Web ACL (Web Access Control List):
 - Rules can include IP addresses, HTTP headers, HTTP body, or URI strings
 - Protects from common attack - SQL injection and Cross-Site Scripting (XSS)
 - **Size constraints, geo-match (block countries)**
 - Rate-based rules (to count occurrences of events) – for DDoS protection

AWS Network Firewall

- Protect your entire Amazon VPC
- From Layer 3 to Layer 7 protection
- Any direction, you can inspect
 - VPC to VPC traffic
 - Outbound to internet
 - Inbound from internet
 - To / from Direct Connect & Site- to-Site VPN

AWS Firewall Manager

- **Manage security rules in all accounts of an AWS Organization**

- Security policy: common set of security rules
 - VPC Security Groups for EC2, Application Load Balancer, etc...
 - WAF rules
 - AWS Shield Advanced
 - AWS Network Firewall
- Rules are applied to new resources as they are created (good for compliance) across all and future accounts in your Organization

Penetration Testing on AWS Cloud

- AWS customers are welcome to carry out security assessments or penetration tests against their AWS infrastructure without prior approval for 8 services:
 - Amazon EC2 instances, NAT Gateways, and Elastic Load Balancers
 - Amazon RDS
 - Amazon CloudFront
 - Amazon Aurora
 - Amazon API Gateways
 - AWS Lambda and Lambda Edge functions
 - Amazon Lightsail resources
 - Amazon Elastic Beanstalk environments
- List can increase over time (you won't be tested on that at the exam)

AWS KMS (Key Management Service)

- Anytime you hear “encryption” for an AWS service, it’s most likely KMS
- KMS = AWS manages the encryption keys for us
- Encryption Opt-in:
 - EBS volumes: encrypt volumes
 - S3 buckets: Server-side encryption of objects
 - Redshift database: encryption of data
 - RDS database: encryption of data
 - EFS drives: encryption of data
- Encryption Automatically enabled:
 - CloudTrail Logs
 - S3 Glacier
 - Storage Gateway

CloudHSM

- KMS => AWS manages the software for encryption
- **CloudHSM => AWS provisions encryption hardware**
- Dedicated Hardware (HSM = Hardware Security Module)
- You manage your own encryption keys entirely (not AWS)
- HSM device is tamper resistant, FIPS 140-2 Level 3 compliance



Types of KMS Keys

- **Customer Managed Key:**
 - Create, manage and used by the customer, can enable or disable
 - Possibility of rotation policy (new key generated every year, old key preserved)
 - Possibility to bring-your-own-key
- **AWS Managed Key:**
 - Created, managed and used on the customer's behalf by AWS
 - Used by AWS services(aws/s3,aws/ebs,aws/redshift)
- **AWS Owned Key:**
 - Collection of CMKs that an AWS service owns and manages to use in multiple accounts
 - AWS can use those to protect resources in your account (but you can't view the keys)
- **CloudHSM Keys (custom keystore):**
 - Keys generated from your own CloudHSM hardware device
 - Cryptographic operations are performed within the CloudHSM cluster

AWS Certificate Manager (ACM)

- Let's you easily provision, manage, and deploy SSL/TLS Certificates
- Used to provide in-flight encryption for websites (HTTPS)
- Supports both public and private TLS certificates
- Free of charge for public TLS certificates
- Automatic TLS certificate renewal
- Integrations with (load TLS certificates on)
 - Elastic Load Balancers
 - CloudFront Distributions
 - APIs on API Gateway

AWS Secrets Manager

AWS Secrets Manager is a service that helps you protect access to your applications, services, and IT resources without the upfront cost and complexity of managing your own hardware security module (HSM) or infrastructure. Secrets Manager enables you to rotate, manage, and retrieve database **credentials**, **API keys**, and other secrets throughout their lifecycle.

- Newer service, meant for storing secrets
- Capability to force rotation of secrets every X days
- Automate generation of secrets on rotation (uses Lambda)
- Integration with Amazon RDS (MySQL, PostgreSQL, Aurora)
- Secrets are encrypted using KMS
- Mostly meant for RDS integration

AWS Artifact (not really a service)

- Portal that provides customers with on-demand access to **AWS compliance documentation and AWS agreements**
- **Artifact Reports** - Allows you to download AWS security and compliance documents from third-party auditors, like AWS ISO certifications, Payment Card Industry (PCI), and System and Organization Control (SOC) reports
- **Artifact Agreements** - Allows you to review, accept, and track the status of AWS agreements such as the Business Associate Addendum (BAA) or the Health Insurance Portability and Accountability Act (HIPAA) for an individual account or in your organization
- Can be used to support internal audit or compliance

Amazon GuardDuty

- **GuardDuty is a threat detection service that continuously monitors for malicious activity and anomalous behavior in AWS accounts.**
- Uses Machine Learning algorithms, anomaly detection, 3rd party data
- One click to enable (30 days trial), no need to install software
- Input data includes:
 - **CloudTrail Events Logs** – unusual API calls, unauthorized deployments
 - **CloudTrailManagementEvents**–createVPCsubnet, createtrail,...
 - CloudTrailS3DataEvents–getobject,listobjects,deleteobject,...
 - **VPC Flow Logs** – unusual internal traffic, unusual IP address
 - **DNS Logs** – compromised EC2 instances sending encoded data within DNS queries
 - Optional Features – EKS Audit Logs, RDS & Aurora, EBS, Lambda, S3 Data Events...
- Can set up EventBridge rules to be notified in case of findings
- EventBridge rules can target AWS Lambda or SNS
- Can protect against CryptoCurrency attacks (has a dedicated “finding” for it)

Amazon Inspector

Amazon Inspector is a powerful tool for automating security assessments of your AWS resources. By identifying vulnerabilities and deviations from best practices, it helps improve your security posture and ensures compliance with industry standards. Whether you need to conduct regular security audits, monitor your environment continuously, or prepare for compliance assessments, Amazon Inspector provides the tools and insights necessary to secure your AWS infrastructure effectively.

- **Automated Security Assessments**
- For EC2 instances
 - Leveraging the AWS System Manager (SSM) agent
 - **Analyze against unintended network accessibility**

- Analyze the running OS against known vulnerabilities
- For Container Images push to Amazon ECR
 - Assessment of Container Images as they are pushed
- For Lambda Functions
 - Identifies software vulnerabilities in function code and package dependencies
 - Assessment of functions as they are deployed
- Reporting & integration with AWS Security Hub
- Send findings to Amazon Event Bridge

What does Amazon Inspector evaluate?

- Remember : only for EC2 instances, Container Images & Lambda functions
- Continuous scanning of the infrastructure, only when needed
- Package vulnerabilities (EC2, ECR & Lambda) – database of CVE
- Network reachability (EC2)
- A risk score is associated with all vulnerabilities for prioritization

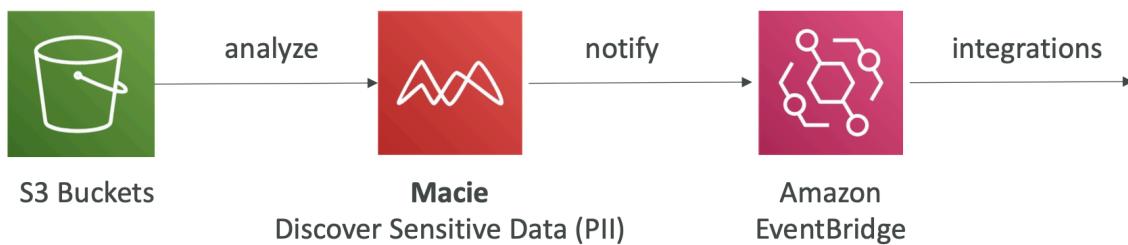
AWS Config

AWS Config is a fully managed service that provides you with an AWS resource inventory, configuration history, and configuration change notifications to enable security and governance. With AWS Config, you can discover AWS resources, record their configurations, and assess configuration compliance against desired configurations.

- Helps with auditing and recording compliance of your AWS resources
- Helps record configurations and changes over time
- Possibility of storing the configuration data into S3 (analyzed by Athena)
- Questions that can be solved by AWS Config:
 - Is there unrestricted SSH access to my security groups?
 - Do my buckets have any public access?
 - How has my ALB configuration changed over time?
- You can receive alerts (SNS notifications) for any changes
- AWS Config is a per-region service
- Can be aggregated across regions and accounts

AWS Macie

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



AWS Security Hub

- Central security tool to manage security across several AWS accounts and automate security checks
- Integrated dashboards showing current security and compliance status to quickly take actions
- Automatically aggregates alerts in predefined or personal findings formats from various AWS services & AWS partner tools:
 - Config
 - GuardDuty
 - Inspector
 - Macie
 - IAMAccessAnalyzer
 - AWS Systems Manager
 - AWS Firewall Manager
 - AWS Health
 - AWS Partner Network Solutions
- Must first enable the AWS Config Service

Amazon Detective

- GuardDuty, Macie, and Security Hub are used to identify potential security issues, or findings
- Sometimes security findings require deeper analysis to isolate the root cause and take action – it's a complex process
- Amazon Detective analyzes, investigates, and quickly identifies the root cause of security issues or suspicious activities (using ML and graphs)
- Automatically collects and processes events from VPC Flow Logs, CloudTrail, GuardDuty and create a unified view
- Produces visualizations with details and context to get to the root cause

AWS Abuse

AWS Abuse refers to any activity conducted on the Amazon Web Services platform that violates the **AWS Acceptable Use Policy** (AUP) or is considered malicious or harmful. AWS takes such activities seriously and has mechanisms in place to detect, report, and address abuse.

- Report suspected AWS resources used for abusive or illegal purposes
- Abusive & prohibited behaviors are:
 - Spam – receiving undesired emails from AWS-owned IP address, websites & forums spammed by AWS resources
 - Port scanning – sending packets to your ports to discover the insecure ones
 - DoS or DDoS attacks – AWS-owned IP addresses attempting to overwhelm or crash your servers/softwares
 - Intrusion attempts – logging in on your resources
 - Hosting objectionable or copyrighted content – distributing illegal or copyrighted content without consent
 - Distributing malware – AWS resources distributing softwares to harm computers or machines
- Contact the AWS Abuse team: AWS abuse form, or abuse@amazonaws.com

Root user privileges

- Root user = Account Owner (created when the account is created)
- Has complete access to all AWS services and resources
- Lock away your AWS account root user access keys!
- Do not use the root account for everyday tasks, even administrative tasks
- ROOT
- • Actions that can be performed only by the root user:
 - Change account settings (account name, email address, root user password, root user access keys)
 - View certain tax invoices
 - Close your AWS account
 - Restore IAM user permissions
 - Change or cancel your AWS Support plan
 - Register as a seller in the Reserved Instance Marketplace
 - Configure an Amazon S3 bucket to enable MFA
 - Edit or delete an Amazon S3 bucket policy that includes an invalid VPC ID or VPC endpoint ID
 - Sign up for GovCloud

Note:

1. It is highly recommended to enable Multi-Factor Authentication (MFA) for root user account
2. Root user access credentials are the email address and password used to create the AWS account

IAM Access Analyzer

- Find out which resources are shared externally
 - S3 Buckets
 - IAM Roles
 - KMS Keys
 - Lambda Functions and Layers • SQS queues
 - Secrets Manager Secrets
- Define Zone of Trust = AWS Account or AWS Organization
- Access outside zone of trusts => findings

Section Summary: Security & Compliance

- **Shared Responsibility on AWS**
- **Shield**: Automatic DDoS Protection + 24/7 support for advanced
- **WAF**: Firewall to filter incoming requests based on rules
- **KMS**: Encryption keys managed by AWS
- **CloudHSM**: Hardware encryption, we manage encryption keys
- **AWS Certificate Manager**: provision, manage, and deploy SSL/TLS Certificates
- **Artifact**: Get access to compliance reports such as PCI, ISO, etc...
- **GuardDuty**: Find malicious behavior with VPC, DNS & CloudTrail Logs
- **Inspector** : find software vulnerabilities in EC2, ECR Images, and Lambda functions
- **Network Firewall**: Protect VPC against network attacks
- **Config**: Track config changes and compliance against rules
- **Macie**: Find sensitive data (ex: PII data) in Amazon S3 buckets
- **CloudTrail**:Track API calls made by users within account
- **AWS Security Hub**: gather security findings from multiple AWS accounts
- **Amazon Detective**: find the root cause of security issues or suspicious activities
- **AWS Abuse**: Report AWS resources used for abusive or illegal purposes
- **Root user privileges**:
 - Change account settings
 - Close your AWS account
 - Change or cancel your AWS Support plan
 - Register as a seller in the Reserved Instance Marketplace
- **IAM Access Analyzer** : identify which resources are shared externally
- **Firewall Manager** : manage security rules across an Organization (WAF, Shield...)

Machine Learning Section

Amazon Rekognition

- Find objects, people, text, scenes in images and videos using ML

- Facial analysis and facial search to do user verification, people counting
- Create a database of “familiar faces” or compare against celebrities
- Use cases:
 - Labeling
 - Content Moderation
 - Text Detection
 - Face Detection and Analysis (gender, age range, emotions...)
 - Face Search and Verification
 - Celebrity Recognition
 - Pathing (ex: for sports game analysis)

Amazon Transcribe

- Automatically convert speech to text
- Uses a deep learning process called automatic speech recognition (ASR) to convert speech to text quickly and accurately
- Automatically remove Personally Identifiable Information (PII) using Redaction
- Supports Automatic Language Identification for multilingual audio
- Use cases:
 - transcribe customer service calls
 - automate closed captioning and subtitling
 - generate metadata for media assets to create a fully searchable archive

Amazon Polly

- Turn text into lifelike speech using deep learning(text to speech)
- Allowing you to create applications that talk

Amazon Translate

- Natural and accurate language translation
- Amazon Translate allows you to localize content - such as websites and applications - for international users, and to easily translate large volumes of text efficiently.

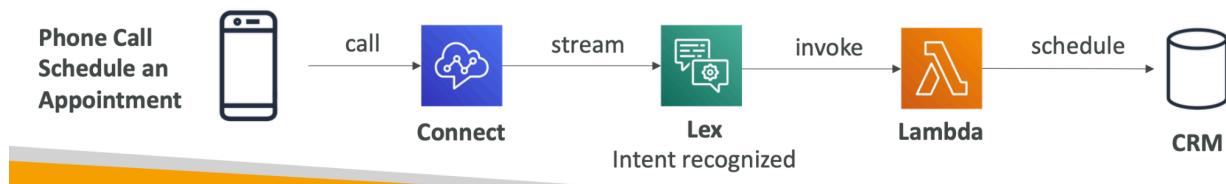
Amazon Lex & Connect

Amazon Lex is a fully managed artificial intelligence (AI) service by AWS that enables developers to build conversational interfaces into applications using voice and text. It is the technology behind Amazon Alexa, providing automatic speech recognition (ASR) to convert speech to text and natural language understanding (NLU) to recognize the intent of the text. With Amazon Lex, you can create sophisticated, **natural language chatbots**, virtual agents, and interactive voice response (IVR) systems.

- Amazon Lex: (same technology that powers Alexa)

- Automatic Speech Recognition (ASR) to convert speech to text
- Natural Language Understanding to recognize the intent of text, callers
- Helps build chatbots, call center bots
- Amazon Connect:
 - Receive calls, create contact flows, cloud-based virtual contact center
 - Can integrate with other CRM systems or AWS
 - No upfront payments, 80% cheaper than traditional contact center solutions

Amazon Connect is a cloud-based contact center service offered by AWS. It allows businesses to set up and manage a customer contact center with ease, providing tools for voice and chat interactions. Amazon Connect aims to simplify the contact center setup process, reduce costs, and enhance customer experiences through features such as intelligent routing, real-time and historical analytics, and seamless integration with other AWS services.



Amazon Comprehend

- For Natural Language Processing – NLP
- Fully managed and serverless service
- Uses machine learning to find insights and relationships in text
 - Language of the text
 - Extracts key phrases, places, people, brands, or events
 - Understands how positive or negative the text is
 - Analyzes text using tokenization and parts of speech
 - Automatically organizes a collection of text files by topic
- Sample use cases:
 - analyze customer interactions (emails) to find what leads to a positive or negative experience
 - Create and group articles by topics that Comprehend will uncover

Amazon SageMaker

- Fully managed service for developers / data scientists to build ML models
- Typically, difficult to do all the processes in one place + provision servers
- Machine learning process (simplified): predicting your exam score

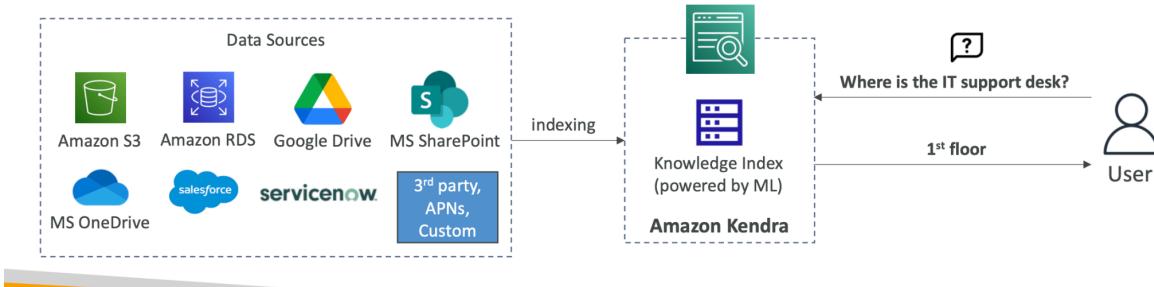
Amazon Forecast

- Fully managed service that uses ML to deliver highly accurate forecasts

- Example: predict the future sales of a raincoat
- 50% more accurate than looking at the data itself
- Reduce forecasting time from months to hours
- Use cases: Product Demand Planning, Financial Planning, Resource Planning, ...

Amazon Kendra

- Fully managed document search service powered by Machine Learning
- Extract answers from within a document (text, pdf, HTML, PowerPoint, MS Word, FAQs...)
- Natural language search capabilities
- Learn from user interactions/feedback to promote preferred results (Incremental Learning)
- Ability to manually fine-tune search results (importance of data, freshness, custom, ...)



Amazon Personalize

- Fully managed ML-service to build apps with real-time personalized recommendations
- Example: personalized product recommendations/re-ranking, customized direct marketing
 - Example: User bought gardening tools, provide recommendations on the next to buy
- Same technology used by Amazon.com
- Integrates into existing websites, applications, SMS, email marketing systems,
- Implement in days, not months (you don't need to build, train, and deploy ML solutions)
- Use cases: retail stores, media and entertainment

Amazon Textract

- Automatically extracts text, handwriting, and data from any scanned documents using AI and ML



- Extract data from forms and tables
- Read and process any type of document (PDFs, images, ...)
- Use cases:
 - Financial Services (e.g., invoices, financial reports)
 - Healthcare (e.g., medical records, insurance claims)
 - Public Sector (e.g., tax forms, ID documents, passports)

AWS Machine Learning - Summary

- **Rekognition**: face detection, labeling, celebrity recognition
- **Transcribe**: audio to text (ex: subtitles)
- **Polly**: text to audio
- **Translate**: translations
- **Lex**: build conversational bots – chatbots
- **Connect**: cloud contact center
- **Comprehend**: natural language processing
- **SageMaker**: machine learning for every developer and data scientist
- **Forecast**: build highly accurate forecasts
- **Kendra**: ML-powered search engine
- **Personalize**: real-time personalized recommendations
- **Textract**: detect text and data in documents

Account Management, Billing & Support Section

AWS Organizations

- **Global service**
- Allows to manage multiple AWS accounts
- The main account is the master account
- Cost Benefits:
 - Consolidated Billing across all accounts - single payment method
 - Pricing benefits from aggregated usage (**volume discount** for EC2, S3...)
 - Pooling of Reserved EC2 instances for optimal savings
- API is available to automate AWS account creation

- Restrict account privileges using **Service Control Policies (SCP)**

Multi Account Strategies

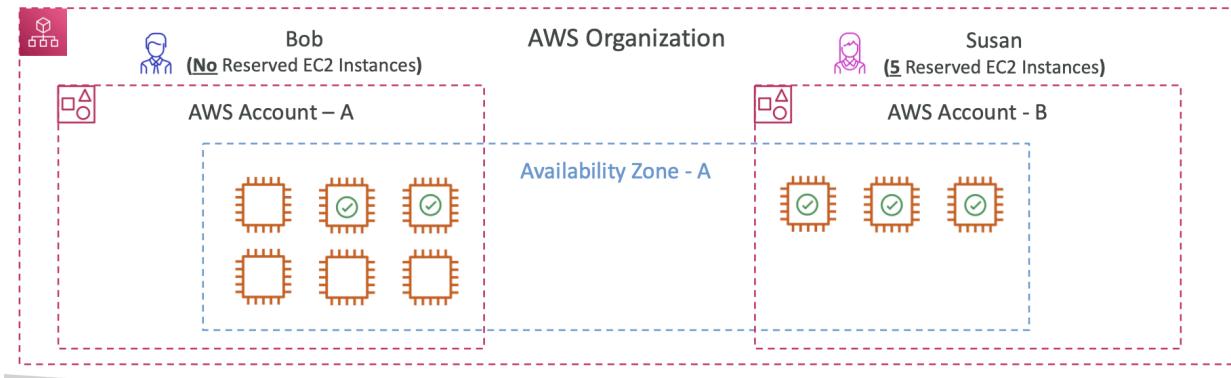
- Create accounts per department, per cost center, per dev / test / prod, based on regulatory restrictions (using SCP), for better resource isolation (ex:VPC), to have separate per-account
- service limits, isolated account for logging
- Multi Account vs One Account Multi VPC
- Use tagging standards for billing purposes
- Enable CloudTrail on all accounts, send logs to central S3 account
- Send CloudWatch Logs to central logging account

Service Control Policies (SCP)

- Whitelist or blacklist IAM actions
- Applied at the OU or Account level
- Does not apply to the Master Account
- SCP is applied to all the **Users and Roles** of the Account, including Root user
- The SCP does not affect service-linked roles
 - Service-linked roles enable other AWS services to integrate with AWS Organizations and can't be restricted by SCPs.
- SCP must have an explicit Allow (does not allow anything by default)
- Use cases:
 - Restrict access to certain services (for example: can't use EMR)
 - Enforce PCI compliance by explicitly disabling services

AWS Organization – Consolidated Billing

- When enabled, provides you with:
 - Combined Usage – combine the usage across all AWS accounts in the AWS Organization to share the volume pricing, Reserved Instances and Savings Plans discounts
 - One Bill – get one bill for all AWS Accounts in the AWS Organization
- The management account can turn off Reserved Instances discount sharing for any account in the AWS Organization, including itself

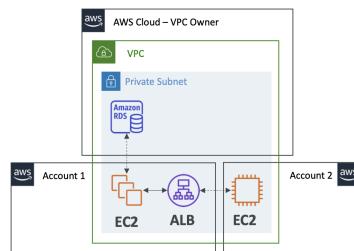


AWS Control Tower

- Easy way to set up and govern a secure and compliant multi-account AWS environment based on best practices
- Benefits:
 - Automate the setup of your environment in a few clicks
 - Automate ongoing policy management using guardrails
 - Detect policy violations and remediate them
 - Monitor compliance through an interactive dashboard
- AWS Control Tower runs on top of AWS Organizations:
 - It automatically sets up AWS Organizations to organize accounts and implement SCPs (Service Control Policies)

AWS Resource Access Manager (AWS RAM)

- Share AWS resources that you own with other AWS accounts
- Share with any account or within your Organization
- Avoid resource duplication!
- Supported resources include Aurora, VPC Subnets, Transit Gateway, Route 53, EC2 Dedicated Hosts, License Manager Configurations...



Pricing Models in AWS

- AWS has 4 pricing models:
 - **Pay as you go**: pay for what you use, remain agile, responsive, meet scale demands
 - **Save when you reserve**: minimize risks, predictably manage budgets, comply with long-term requirements
 - Reservations are available for EC2 Reserved Instances, DynamoDB Reserved Capacity, ElastiCache Reserved Nodes, RDS Reserved Instance, Redshift Reserved Nodes
 - **Pay less by using more**: volume-based discounts
 - **Pay less as AWS grows**

Free services & free tier in AWS

- IAM
 - VPC
 - Consolidated Billing
 - Elastic Beanstalk
 - CloudFormation
 - Auto Scaling Groups
 - Free Tier : <https://aws.amazon.com/free/>
 - EC2 t2.micro instance for a year
 - S3, EBS, ELB, AWS Data transfer
- }  You do pay for the resources created

Compute Pricing – EC2

- Only charged for what you use
- Number of instances
- Instance configuration:
 - Physical capacity
 - Region
 - OS and software
 - Instance type
 - Instance size
- ELB running time and amount of data processed
- Detailed monitoring

Compute Pricing – EC2

- **On-demand instances:**
 - Minimum of 60s
 - Pay per second (Linux/Windows) or per hour (other)

- **Reserved instances:**
 - Up to 75% discount compared to On-demand on hourly rate
 - 1- or 3-years commitment
 - All upfront, partial upfront, no upfront
- **Spot instances:**
 - Up to 90% discount compared to On-demand on hourly rate
 - Bid for unused capacity
- **Dedicated Host:**
 - On-demand
 - Reservation for 1 year or 3 years commitment
- Savings plans as an alternative to save on sustained usage

Compute Pricing – Lambda & ECS

- Lambda:
 - Pay per call
 - Pay per duration
- ECS:
 - EC2 Launch Type Model: No additional fees, you pay for AWS resources stored and created in your application
- Fargate:
 - Fargate Launch Type Model: Pay for vCPU and memory resources allocated to your applications in your containers

Storage Pricing – S3

- **Storage class:** S3 Standard, S3 Infrequent Access, S3 One-Zone IA, S3 Intelligent Tiering, S3 Glacier and S3 Glacier Deep Archive
- Number and size of objects: Price can be tiered (based on volume)
- Number and type of requests
- Data transfer OUT of the S3 region
- S3 Transfer Acceleration
- Lifecycle transitions
- Similar service: EFS (pay per use, has infrequent access & lifecycle rules)

Storage Pricing - EBS

- Volume type (based on performance)
Storage volume in GB per month provisioned
- IOPS:
 - General Purpose SSD: Included
 - Provisioned IOPS SSD: Provisioned amount in IOPS
 - Magnetic: Number of requests
- Snapshots:
 - Added data cost per GB per month

- Data transfer:
 - Outbound data transfer are tiered for volume discounts
 - Inbound is free

Database Pricing - RDS

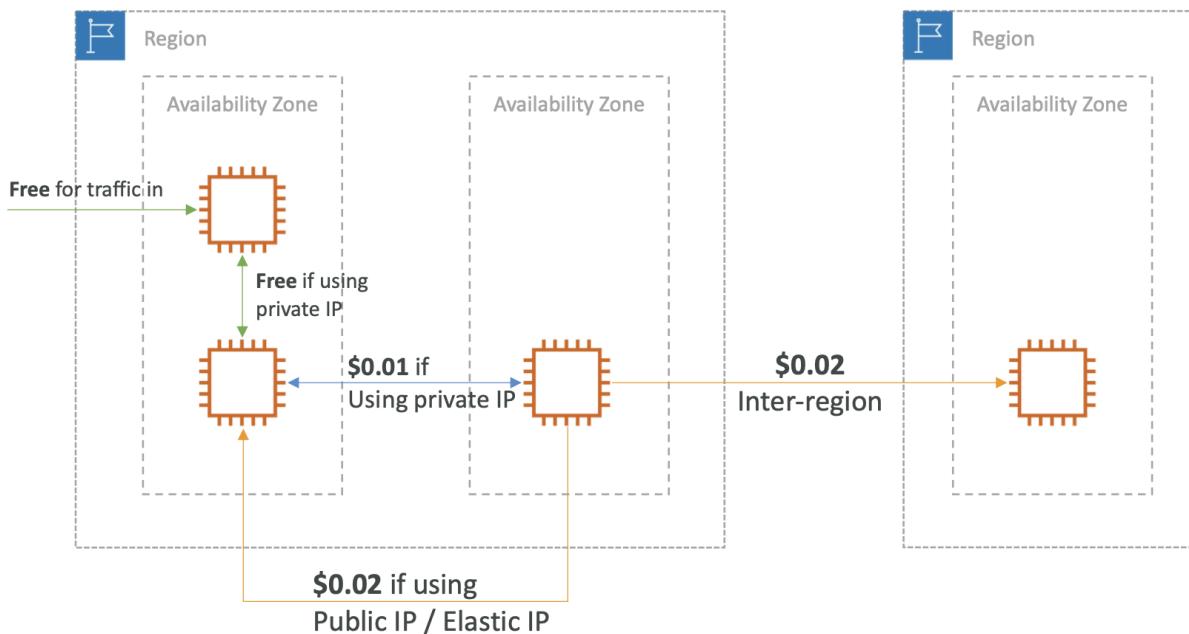
- Per hour billing
- Database characteristics:
 - Engine
 - Size
 - Memory class
- Purchase type:
 - On-demand
 - Reserved instances (1 or 3 years) with required up-front
- Backup Storage: There is no additional charge for backup storage up to 100% of your total database storage for a region.
- Additional storage (per GB per month)
- Number of input and output requests per month
- Deployment type (storage and I/O are variable):
 - Single AZ
 - Multiple AZs
- Data transfer:
- Outbound data transfer are tiered for volume discounts
- Inbound is free

Content Delivery – CloudFront

- Pricing is different across different geographic regions
- Aggregated for each edge location, then applied to your bill
- Data Transfer Out (volume discount)
- Number of HTTP/HTTPS requests

Networking Costs in AWS per GB - Simplified

- Use Private IP instead of Public IP for good savings and better network performance
- Use same AZ for maximum savings (at the cost of high availability)



Savings Plan

- Commit a certain \$ amount per hour for 1 or 3 years
- Easiest way to setup long-term commitments on AWS
- EC2 Savings Plan
 - Up to 72% discount compared to On-Demand
 - Commit to usage of individual instance families in a region (e.g. C5 or M5)
 - Regardless of AZ, size (m5.xl to m5.4xl), OS (Linux/Windows) or tenancy
 - All upfront, partial upfront, no upfront
- Compute Savings Plan
 - Up to 66% discount compared to On-Demand
 - Regardless of Family, Region, size, OS, tenancy, compute options
 - ComputeOptions:EC2,Fargate,Lambda
- Machine Learning Savings Plan: SageMaker
- Setup from the AWS Cost Explorer console

AWS Compute Optimizer

- Reduce costs and improve performance by recommending optimal AWS resources for your workloads
- Helps you choose optimal configurations and right-size your workloads (over/under provisioned)
- Uses Machine Learning to analyze your resources' configurations and their utilization CloudWatch metrics
- Supported resources
 - a. EC2 instances

- b. EC2 Auto Scaling Groups
- c. EBS volumes
- d. Lambda functions
- Lower your costs by up to 25%
- Recommendations can be exported to S3

Billing and Costing Tools

- Estimating costs in the cloud:
 - Pricing Calculator
- Tracking costs in the cloud:
 - Billing Dashboard
 - Cost Allocation Tags
 - Cost and Usage Reports
 - Cost Explorer
- Monitoring against costs plans:
 - Billing Alarms
 - Budgets

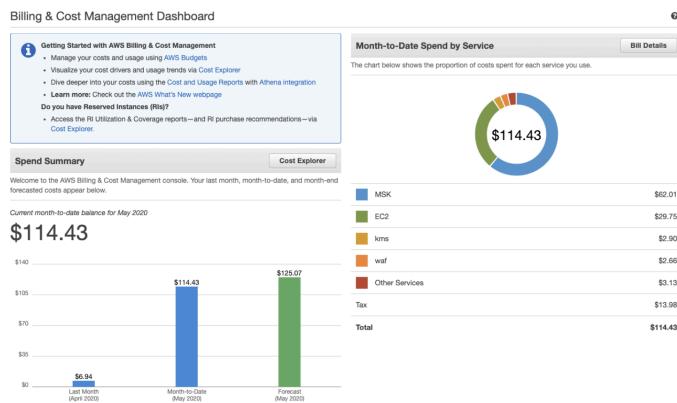
AWS Pricing Calculator

The AWS Pricing Calculator is a web-based tool provided by Amazon Web Services (AWS) that helps users estimate the cost of AWS services based on their specific usage scenarios. It allows you to model your architecture and get a detailed breakdown of costs, helping you plan and budget effectively.

- Estimate the cost for your solution architecture

AWS Billing Dashboard

The AWS Billing Dashboard is a comprehensive and user-friendly interface within the AWS Management Console that provides an overview of your AWS account's usage, charges, and billing information. It helps users manage their costs and optimize their spending on AWS services.



Cost Allocation Tags

- Use cost allocation tags to track your AWS costs on a detailed level
- AWS generated tags
 - Automatically applied to the resource you create
 - Starts with Prefix aws: (e.g. aws: createdBy)
- User-defined tags
 - Defined by the user
 - Starts with Prefix user:

Total Cost	user:Owner	user:Stack	user:Cost Center	user:Application
0.95	DbAdmin	Test	80432	Widget2
0.01	DbAdmin	Test	80432	Widget2
3.84	DbAdmin	Prod	80432	Widget2
6.00	DbAdmin	Test	78925	Widget1
234.63	SysEng	Prod	78925	Widget1
0.73	DbAdmin	Test	78925	Widget1
0.00	DbAdmin	Prod	80432	Portal
2.47	DbAdmin	Prod	78925	Portal

Tagging and Resource Groups

- Tags are used for organizing resources:
 - EC2: instances, images, load balancers, security groups
 - RDS, VPC resources, Route 53, IAM users, etc.
 - Resources created by CloudFormation are all tagged the same way
- Free naming, common tags are: Name, Environment, Team
- Tags can be used to create **Resource Groups**
 - Create, maintain, and view a collection of resources that share common tags
 - Manage these tags using the Tag Editor

Cost and Usage Reports

- Dive deeper into your AWS costs and usage
- The AWS Cost & Usage Report contains the most comprehensive set of AWS cost and usage data available, including additional metadata about AWS services, pricing, and reservations (e.g., Amazon EC2 Reserved Instances (RIs)).
- The AWS Cost & Usage Report lists AWS usage for each service category used by an account and its IAM users in hourly or daily line items, as well as any tags that you have activated for cost allocation purposes.
- Can be integrated with Athena, Redshift or QuickSight

AWS Billing Dashboard vs AWS Cost and Usage Reports

AWS Billing Dashboard: Ideal for users needing a quick and comprehensive overview of their AWS costs and billing information. It's great for managing payments, setting up budgets, and monitoring overall spending.

AWS Cost and Usage Reports (CUR): Best for users who require detailed, **granular** data about their AWS costs and usage. It's perfect for in-depth analysis, custom reporting, and integration with other analytical tools.

Cost Explorer

AWS Cost Explorer is an essential tool for any organization using AWS, providing detailed insights into costs and usage. With its powerful reporting, filtering, and visualization capabilities, Cost Explorer enables you to understand spending patterns, identify cost-saving opportunities, and make informed financial decisions. Whether you are managing a single account or a complex multi-account environment, AWS Cost Explorer helps you gain control over your cloud costs and optimize your AWS spending.

- Visualize, understand, and manage your AWS costs and usage over time
- Create custom reports that analyze cost and usage data.
- Analyze your data at a high level: total costs and usage across all accounts
- Or Monthly, hourly, resource level granularity
- Choose an optimal Savings Plan (to lower prices on your bill)
- Forecast usage up to 12 months based on previous usage

Cost Explorer – Savings Plan

AWS Cost Explorer Savings Plans are a powerful tool for reducing your AWS costs by committing to a consistent amount of usage over a specified term. By leveraging the flexibility, significant cost savings, and detailed insights provided by Cost Explorer, you can optimize your AWS spending and achieve greater financial efficiency. Whether you are new to AWS or looking to optimize existing workloads, Savings Plans offer a straightforward and effective way to manage and reduce your cloud costs.

Cost Explorer – Forecast Usage

Forecasting in AWS Cost Explorer is a valuable tool for predicting future AWS costs and usage based on historical data. By providing detailed and customizable forecasts, it helps you plan your budget, manage resources proactively, and optimize your cloud spending. Whether you are managing a single account or a complex multi-account environment, the forecasting capabilities of AWS Cost Explorer enable you to make informed financial decisions and achieve greater control over your AWS expenditures.



Billing Alarms in CloudWatch

- Billing data metric is stored in CloudWatch us-east-1
- Billing data are for overall worldwide AWS costs
- It's for actual cost, not for projected costs
- Intended a simple alarm (not as powerful as AWS Budgets)

! Alert
Threshold-based, via SNS

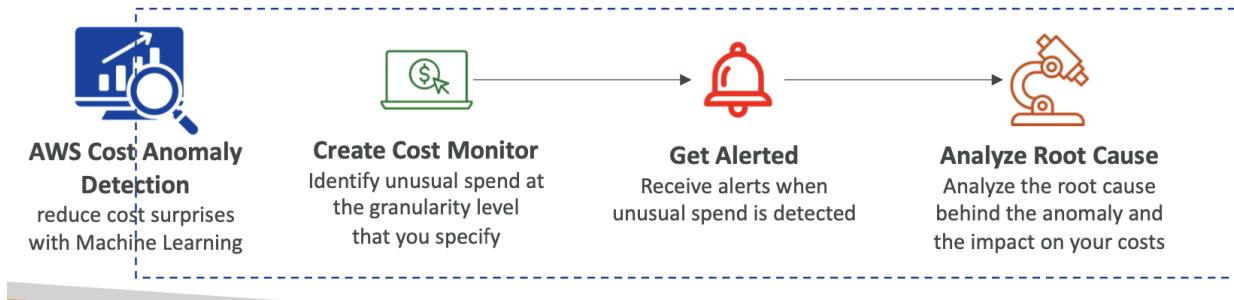
AWS Budgets

- Create budget and send alarms when costs exceeds the budget
- 4 types of budgets: Usage, Cost, Reservation, Savings Plans
- For Reserved Instances (RI)
 - Track utilization
 - Supports EC2, ElastiCache, RDS, Redshift
- Up to 5 SNS notifications per budget
- Can filter by: Service, Linked Account, Tag, Purchase Option, Instance Type, Region, Availability Zone, API Operation, etc...
- Same options as AWS Cost Explorer!
- 2 budgets are free, then \$0.02/day/budget

! Alert
Detailed, multi-threshold, via email/SNS

AWS Cost Anomaly Detection

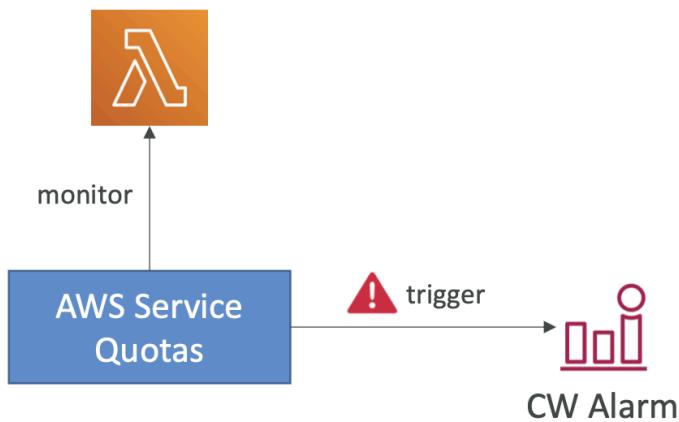
- Continuously monitor your cost and usage using ML to detect unusual spends
- It learns your unique, historic spend patterns to detect one-time cost spike and/or continuous cost increases (you don't need to define thresholds)
- Monitor AWS services, member accounts, cost allocation tags, or cost categories
- Sends you the anomaly detection report with root-cause analysis
- Get notified with individual alerts or daily/weekly summary (using SNS)



AWS Service Quotas

- Notify you when you're close to a service quota value threshold
- Create CloudWatch Alarms on the Service Quotas console
- Example: Lambda concurrent executions
- Request a quota increase from AWS Service Quotas or shutdown resources before limit is reached

AWS Lambda Quota



Trusted Advisor

- No need to install anything – high level AWS account assessment
- Analyze your AWS accounts and provides recommendation on 6 categories:

- Cost optimization
- Performance
- Security
- Fault tolerance
- Service limits
- Operational Excellence
- Business & Enterprise Support plan
- Full Set of Checks
- Programmatic Access using AWS Support API

AWS Basic Support Plan

- **Customer Service & Communities** - 24x7 access to customer service, documentation, whitepapers, and support forums.
- **AWS Trusted Advisor** - Access to the 7 core Trusted Advisor checks and guidance to provision your resources following best practices to increase performance and improve security.
- **AWS Personal Health Dashboard** - A personalized view of the health of AWS services, and alerts when your resources are impacted.

AWS Developer Support Plan

- All Basic Support Plan +
- **Business hours email access to Cloud Support Associates**
- Unlimited cases / 1 primary contact
- Case severity / response times:
 - General guidance: < 24 business hours
 - System impaired: < 12 business hours

AWS Business Support Plan (24/7)

- Intended to be used if you have production workloads
- Trusted Advisor – Full set of checks + API access
- **24x7 phone, email, and chat access to Cloud Support Engineers**
- Unlimited cases / unlimited contacts
- **Access to Infrastructure Event Management for additional fee.**
- Case severity / response times:
 - General guidance: < 24 business hours
 - System impaired: < 12 business hours
 - Production system impaired: < 4 hours
 - **Production system down: < 1 hour**

AWS Enterprise On-Ramp Support Plan (24/7)

- Intended to be used if you have production or business critical workloads
- All of Business Support Plan +

- Access to a pool of Technical Account Managers (TAM)
- Concierge Support Team (for billing and account best practices)
- Infrastructure Event Management, Well-Architected & Operations Reviews
- Case severity / response times:
 - Production system impaired: < 4 hours
 - Production system down: < 1 hour
 - Business-critical system down: < 30 minutes

AWS Enterprise Support Plan (24/7)

- Intended to be used if you have mission critical workloads
- All of Business Support Plan +
- Access to a designated Technical Account Manager (TAM)
- Concierge Support Team (for billing and account best practices)
- Infrastructure Event Management, Well-Architected & Operations Reviews
- Case severity / response times:
 - Production system impaired: < 4 hours
 - Production system down: < 1 hour
 - Business-critical system down: < 15 minutes

Account Best Practices – Summary

- Operate multiple accounts using **Organizations**
- Use **SCP** (service control policies) to restrict account power
- Easily setup multiple accounts with best-practices with **AWS Control Tower**
- **Use Tags & Cost Allocation Tags** for easy management & billing
- **IAM guidelines**: MFA, least-privilege, password policy, password rotation
- **Config** to record all resources configurations & compliance over time
- **CloudFormation** to deploy stacks across accounts and regions
- Trusted Advisor to get insights, Support Plan adapted to your needs
- Send Service Logs and Access Logs to S3 or CloudWatch Logs
- **CloudTrail** to record API calls made within your account
- If your Account is compromised: change the root password, delete and rotate all passwords / keys, contact the AWS support
- Allow users to create pre-defined stacks defined by admins using AWS Service Catalog

Billing and Costing Tools – Summary

- **Compute Optimizer** : recommends resources' configurations to reduce cost
- **Pricing Calculator**: cost of services on AWS
- **Billing Dashboard**: high level overview + free tier dashboard
- **Cost Allocation Tags**: tag resources to create detailed reports
- **Cost and Usage Reports**: most comprehensive billing dataset
- **Cost Explorer** : View current usage (detailed) and forecast usage
- **Billing Alarms**: in us-east-1 – track overall and per-service billing

- **Budgets:** more advanced – track usage, costs, RI, and get alerts
- **Savings Plans:** easy way to save based on long-term usage of AWS
- **Cost Anomaly Detection:** detect unusual spends using Machine Learning
- **Service Quotas:** notify you when you're close to service quota threshold

Advanced Identity Section

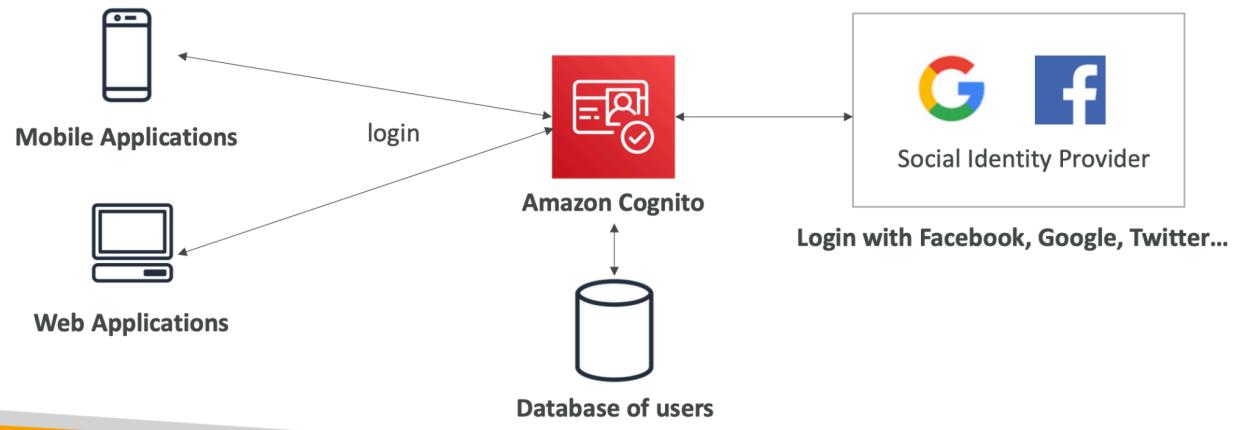
AWS STS (Security Token Service)

- Enables you to create temporary, limited- privileges credentials to access your AWS resources
- Short-term credentials: you configure expiration period
- Use cases
 - Identity federation: manage user identities in external systems, and provide them with STS tokens to access AWS resources
 - IAM Roles for cross/same account access
 - IAM Roles for Amazon EC2: provide temporary credentials for EC2 instances to access AWS resources

Amazon Cognito (simplified)

Cognito can allow users access to AWS resources using federation through third-party sources like Google, Facebook, etc.

- Identity for your Web and Mobile applications users (potentially millions)
- Instead of creating them an IAM user, you create a user in Cognito



AWS IAM Identity Center (successor to AWS Single Sign-On)

- One login (single sign-on) for all your
 - AWS accounts in AWS Organizations

- Business cloud applications (e.g., Salesforce, Box, Microsoft 365, ...)
- SAML2.0-enabled applications
- EC2 Windows Instances
- Identity providers
 - Built-in identity store in IAM Identity Center
 - 3rd party: Active Directory (AD), OneLogin, Okta...

Advanced Identity - Summary

- IAM
 - Identity and Access Management inside your AWS account
 - For users that you trust and belong to your company
- **Organizations:** Manage multiple accounts
- **Security Token Service (STS):** temporary, limited-privileges credentials to access AWS resources
- **Cognito:** create a database of users for your mobile & web applications
- **Directory Services:** integrate Microsoft Active Directory in AWS
- **IAM Identity Center :** one login for multiple AWS accounts & applications

Other AWS Services

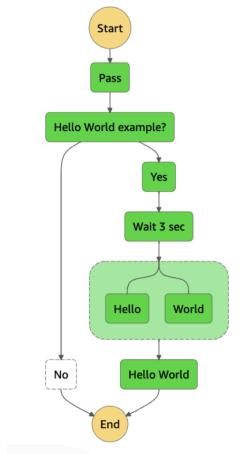
AWS Migration Hub

- Central location to collect servers and applications inventory data for the assessment, planning, and tracking of migrations to AWS
- Helps accelerate your migration to AWS, automate lift-and-shift
- AWS Migration Hub Orchestrator – provides pre-built templates to save time and effort migrating enterprise apps (e.g., SAP, Microsoft SQL Server...)
- Supports migrations status updates from Application Migration Service (MGN) and Database Migration Service (DMS)

AWS Step Functions

AWS Step Functions is a serverless workflow service that enables you to coordinate and orchestrate multiple AWS Lambda functions in a specific order. It provides visual workflows, error handling, and retries to ensure reliable execution.

- Build serverless visual workflow to orchestrate your Lambda functions
- Features: sequence, parallel, conditions, timeouts, error handling, ...
- Can integrate with EC2, ECS, On-premises servers, API Gateway, SQS queues, etc...
- Possibility of implementing human approval feature
- Use cases: order fulfillment, data processing, web applications, any workflow



AWS Architecting & Ecosystem Section

Well Architected Framework General Guiding Principles

- **Stop guessing your capacity needs**
- Test systems at production scale
- Automate to make architectural experimentation easier
- Allow for evolutionary architectures
 - Design based on changing requirements
- Drive architectures using data
- Improve through game days
 - Simulate applications for flash sale days

AWS Cloud Best Practices – Design Principles

- **Scalability:** vertical & horizontal
- **Disposable Resources:** servers should be disposable & easily configured
- **Automation:** Serverless, Infrastructure as a Service, Auto Scaling...
- **Loose Coupling:**
 - Monolith are applications that do more and more over time, become bigger
 - Break it down into smaller, loosely coupled components
 - A change or a failure in one component should not cascade to other components
- **Services, not Servers:**
 - Don't use just EC2
 - Use managed services, databases, serverless, etc !

Well Architected Framework 6 Pillars

1. Operational Excellence
2. Security
3. Reliability
4. Performance Efficiency
5. Cost Optimization
6. Sustainability

1. Operational Excellence

- Includes the ability to run and monitor systems to deliver business value and to continually improve supporting processes and procedures
- Design Principles
 - **Perform operations as code** - Infrastructure as code
 - **Make frequent, small, reversible changes** - So that in case of any failure, you can reverse it
 - **Refine operations procedures frequently** - And ensure that team members are familiar with it
 - **Anticipate failure** - Learn from all operational failures
 - **Use managed services** - to reduce operational burden
 - Implement observability for actionable insights - performance, reliability, cost

2. Security

- Includes the ability to protect information, systems, and assets while delivering business value through risk assessments and mitigation strategies
- Design Principles
 - **Implement a strong identity foundation** - Centralize privilege management and reduce (or even eliminate) reliance on long-term credentials - Principle of least privilege - IAM
 - **Enable traceability** - Integrate logs and metrics with systems to automatically respond and take action
 - **Apply security at all layers** - Like edge network, VPC, subnet, load balancer, every instance, operating system, and application
 - Automate security best practices
 - Protect data in transit and at rest - Encryption, tokenization, and access control
 - Keep people away from data - Reduce or eliminate the need for direct access or manual processing of data
 - Prepare for security events - Run incident response simulations and use tools with automation to increase your speed for detection, investigation, and recovery
 - Shared Responsibility Model

3. Reliability

- Ability of a system to recover from infrastructure or service disruptions, dynamically acquire computing resources to meet demand, and mitigate disruptions such as misconfigurations or transient network issues
- Design Principles
 - Test recovery procedures - Use automation to simulate different failures or to recreate
 - **Automatically recover from failure** - Anticipate and remediate failures before they occur

- **Scale horizontally to increase aggregate system availability** - Distribute requests across multiple, smaller resources to ensure that they don't share a common point of failure
- **Stop guessing capacity** - Maintain the optimal level to satisfy demand without over or under provisioning - Use Auto Scaling
- **Manage change in automation** - Use automation to make changes to infrastructure

4. Performance Efficiency

- Includes the ability to use computing resources efficiently to meet system requirements, and to maintain that efficiency as demand changes and technologies evolve
- Design Principles
 - **Democratize advanced technologies** - Advance technologies become services and hence you can focus more on product development
 - **Go global in minutes** - Easy deployment in multiple regions
 - **Use serverless architectures** - Avoid burden of managing servers
 - **Experiment more often** - Easy to carry out comparative testing
 - **Mechanical sympathy** - Be aware of all AWS services

5. Cost Optimization

- Includes the ability to run systems to deliver business value at the lowest price point
- Design Principles
 - **Adopt a consumption mode** - Pay only for what you use
 - **Measure overall efficiency** - Use CloudWatch
 - **Stop spending money on data center operations** - AWS does the infrastructure part and enables customer to focus on organization projects
 - **Analyze and attribute expenditure** - Accurate identification of system usage and costs, helps measure return on investment (ROI) - Make sure to use tags
 - **Use managed and application level services to reduce cost of ownership** - As managed services operate at cloud scale, they can offer a lower cost per transaction or service

6. Sustainability

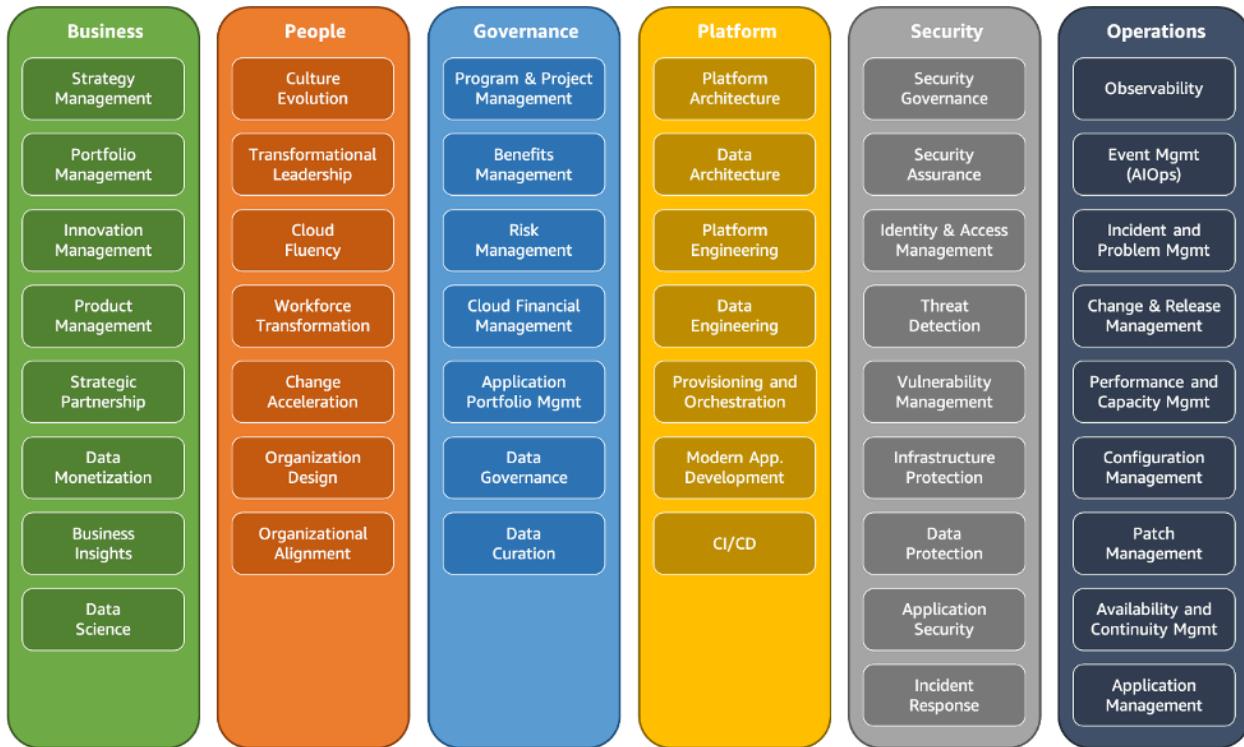
- The sustainability pillar focuses on minimizing the environmental impacts of running cloud workloads.
- Design Principles
 - **Understand your impact** – establish performance indicators, evaluate improvements
 - **Establish sustainability goals** – Set long-term goals for each workload, model return on investment (ROI)

- **Maximize utilization** – Right size each workload to maximize the energy efficiency of the underlying hardware and minimize idle resources.
- Anticipate and adopt new, more efficient hardware and software offerings – and design for flexibility to adopt new technologies over time.
- **Use managed services** – Shared services reduce the amount of infrastructure; Managed services help automate sustainability best practices as moving infrequent accessed data to cold storage and adjusting compute capacity.
- **Reduce the downstream impact of your cloud workloads** – Reduce the amount of energy or resources required to use your services and reduce the need for your customers to upgrade their devices

AWS Cloud Adoption Framework (AWS CAF)

- Helps you build and then execute a comprehensive plan for your digital transformation through innovative use of AWS
- Created by AWS Professionals by taking advantage of AWS Best Practices and lessons learned from 1000s of customers
- AWS CAF identifies specific organizational capabilities that underpin successful cloud transformations
- **AWS CAF groups its capabilities in 6 perspectives:**
Business, People, Governance, Platform, Security, and Operations

AWS Cloud Adoption Framework (AWS CAF) - Foundational capabilities:



CAF Perspectives and Foundational Capabilities Business Capabilities

- Business Perspective** helps ensure that your cloud investments accelerate your digital transformation ambitions and business outcomes.
- People Perspective** serves as a bridge between technology and business, accelerating the cloud journey to help organizations more rapidly evolve to a culture of continuous growth, learning, and where change becomes business-as-normal, with focus on culture, organizational structure, leadership, and workforce.
- Governance Perspective** helps you orchestrate your cloud initiatives while maximizing organizational benefits and minimizing transformation-related risks.
- Platform Perspective** helps you build an enterprise-grade, scalable, hybrid cloud platform; modernize existing workloads; and implement new cloud-native solutions.
- Security Perspective** helps you achieve the confidentiality, integrity, and availability of your data and cloud workloads.
- Operations Perspective** helps ensure that your cloud services are delivered at a level that meets the needs of your business.

AWS CAF – Transformation Phases

- **Envision** – demonstrate how the Cloud will accelerate business outcomes by identifying transformation opportunities and create a foundation for your digital transformation
- **Align** – identify capability gaps across the 6 AWS CAF Perspectives which results in an Action Plan
- **Launch** – build and deliver pilot initiatives in production and demonstrate incremental business value
- **Scale** – expand pilot initiatives to the desired scale while realizing the desired business benefits

AWS Marketplace

- Digital catalog with thousands of software listings from independent software vendors (3rd party)
- **Example:**
 - Custom AMI (custom OS, firewalls, technical solutions...)
 - CloudFormation templates
 - Software as a Service
 - Containers
- If you buy through the AWS Marketplace, it goes into your AWS bill
- You can sell your own solutions on the AWS Marketplace

AWS Training

- AWS Digital (online) and **ClassroomTraining** (in-person or virtual)
- AWS Private Training (for your organization)
- Training and Certification for the U.S Government
- Training and Certification for the Enterprise
- AWS Academy: helps universities teach AWS
- And your favorite online teacher teaching you all about AWS Certifications and more!

AWS Professional Services & Partner Network

- The AWS Professional Services organization is a global team of experts
- They work alongside your team and a chosen member of the APN
- APN = AWS Partner Network
- **APN Technology Partners:** providing hardware, connectivity, and software
- **APN Consulting Partners:** professional services firm to help build on AWS
- APN Training Partners: find who can help you learn AWS
- AWS Competency Program: AWS Competencies are granted to APN Partners who have demonstrated technical proficiency and proven customer success in specialized solution areas.
- AWS Navigate Program: help Partners become better Partners

AWS IQ

AWS IQ is a service provided by Amazon Web Services that connects AWS customers with AWS Certified experts for on-demand project work. This platform allows users to find, hire, and collaborate with AWS experts to solve technical challenges, build solutions, and optimize their use of AWS services.

- Quickly find professional help for your AWS projects
- Engage and pay AWS Certified 3rd party experts for on-demand project work
- Video-conferencing, contract management, secure collaboration, integrated billing

AWS re:Post

- AWS-managed Q&A service offering crowd-sourced, expert-reviewed answers to your technical questions about AWS that replaces the original AWS Forums
- Part of the AWS FreeTier
- Community members can earn reputation points to build up their community expert status by providing accepted answers and reviewing answers from other users
- Questions from AWS Premium Support customers that do not receive a response from the community are passed on to AWS Support engineers
- AWS re:Post is not intended to be used for questions that are time-sensitive or involve any proprietary information

AWS re:Post – Knowledge Center

- Contains the most frequent & common questions and requests

AWS Managed Services (AMS)

- Provides infrastructure and application support on AWS.
- AMS offers a team of AWS experts who manage and operate your infrastructure for security, reliability, and availability
- Helps organizations offload routine management tasks and focus on their business objectives.
- Fully managed service, so AWS handles common activities such as change requests, monitoring, patch management, security, and backup services
- Implements best practices and maintains your AWS infrastructure to reduce your operational overhead and risk
- AMS business hours are 24/365

