AWS Solutions Architect Associate

Complete AWS Solutions Architect Associate (SAA-CO3) Certification Study Guide

Table of Contents

- 1. Introduction and Exam Overview
- 2. AWS Fundamentals
- 3. Identity and Access Management (IAM)
- 4. Amazon EC2 Elastic Compute Cloud
- 5. Storage Services
- 6. Networking and Content Delivery
- 7. Database Services
- 8. Security and Compliance
- 9. Monitoring and Management
- 10. Application Integration
- 11. Analytics and Machine Learning
- 12. <u>Disaster Recovery and Migration</u>
- 13. Cost Optimization
- 14. Well-Architected Framework
- 15. Practice Questions and Exam Tips

Introduction and Exam Overview {#introduction}

AWS Solutions Architect Associate

Exam Details

• Exam Code: SAA-C03

• **Duration**: 130 minutes

• Questions: 65 questions (multiple choice and multiple response)

• Passing Score: 720 out of 1000

• Cost: \$150 USD

• Validity: 3 years

Exam Domains

1. Design Resilient Architectures (26%)

- 2. Design High-Performing Architectures (24%)
- 3. Design Secure Architectures (30%)
- 4. Design Cost-Optimized Architectures (20%)

Key Success Factors

- Hands-on Experience: Minimum 1 year recommended
- Scenario-based Questions: Focus on real-world implementations
- Best Practices: Always choose the most cost-effective, secure, and scalable solution

AWS Fundamentals {#aws-fundamentals}

Global Infrastructure

AWS Regions

- **Definition**: Geographic areas containing multiple data centers
- Current Count: 31+ regions globally (as of 2024)
- Key Considerations:
 - Latency: Choose regions close to users
 - Compliance: Data governance requirements
 - Service Availability: Not all services available in all regions

Pricing: Varies by region

Availability Zones (AZs)

- Definition: Isolated data centers within a region
- Characteristics:
 - Minimum 3 AZs per region (typically 3-6)
 - Physically separated (different buildings)
 - Connected via high-bandwidth, low-latency links
 - Independent power, cooling, and networking

Edge Locations

- Count: 400+ globally
- Purpose: Content caching for CloudFront CDN
- · Benefits: Reduced latency for global users

AWS Services Overview

- 200+ Services across multiple categories
- Global Services: IAM, Route 53, CloudFront, WAF
- Regional Services: EC2, RDS, Lambda, S3
- AZ-Specific Services: Subnets, EBS volumes

Identity and Access Management (IAM) {#iam}

Core Components

Users

- **Definition**: Individual identity for human or application
- Characteristics:
 - Permanent credentials (username/password)
 - Access keys for programmatic access
 - Can belong to multiple groups

Direct policy attachment possible (not recommended)

Groups

- Purpose: Collection of users with similar permissions
- Best Practice: Assign permissions to groups, not individuals
- Limitations: Cannot contain other groups

Roles

- **Definition**: Temporary credentials for AWS services or external identities
- Use Cases:
 - EC2 instances accessing other AWS services
 - Cross-account access
 - Federation with external identity providers
- Key Benefits:
 - No permanent credentials
 - Automatic credential rotation
 - More secure than access keys

Policies

- Types:
 - AWS Managed: Created and maintained by AWS
 - Customer Managed: Custom policies you create
 - **Inline**: Embedded directly in user, group, or role

Policy Structure (JSON)

```
{
"Version": "2012-10-17",
"Statement": [
    {
      "Sid": "AllowS3ReadAccess",
      "Effect": "Allow",
      "Action": [
```

```
"s3:GetObjectVersion"
],
   "Resource": "arn:aws:s3:::my-bucket/*",
   "Condition": {
       "StringEquals": {
       "s3:prefix": "home/${aws:username}/"
       }
    }
}
```

Security Best Practices

Multi-Factor Authentication (MFA)

- Types:
 - Virtual MFA: Google Authenticator, Authy
 - Hardware MFA: YubiKey, Gemalto
 - SMS MFA: Text message (least secure)

Password Policy

- Requirements:
 - Minimum length
 - Character complexity
 - Password expiration
 - Password reuse prevention

Access Keys Management

- Best Practices:
 - Rotate regularly
 - Use IAM roles instead when possible
 - Never embed in code

Use AWS Secrets Manager or Parameter Store

Advanced IAM Features

Permission Boundaries

- Purpose: Maximum permissions an entity can have
- Use Cases:
 - Delegate permissions without escalation risk
 - Compliance requirements

Service-Linked Roles

- **Definition**: Pre-defined roles for AWS services
- Characteristics:
 - Created automatically by services
 - Cannot be modified
 - Specific to service requirements

Amazon EC2 - Elastic Compute Cloud {#ec2}

Instance Types and Families

General Purpose

- T3/T3a/T4g: Burstable performance
 - Use Cases: Web servers, small databases, development environments
 - **CPU Credits**: Baseline performance with burst capability
- M5/M5a/M6i: Balanced compute, memory, networking
 - Use Cases: Web applications, microservices, enterprise applications

Compute Optimized

- C5/C5n/C6i: High-performance processors
 - Use Cases: Scientific modeling, batch processing, gaming servers
 - Features: Enhanced networking, high memory bandwidth

Memory Optimized

- R5/R6i: High memory-to-vCPU ratio
 - Use Cases: In-memory databases, real-time analytics
- X1e: Extremely high memory (up to 3.8 TB RAM)
 - Use Cases: SAP HANA, Apache Spark

Storage Optimized

- I3: High sequential read/write with NVMe SSD
 - Use Cases: Distributed file systems, data warehousing
- **D2**: Dense HDD storage
 - Use Cases: MapReduce, distributed computing

Accelerated Computing

- P3/P4: GPU instances for ML/AI
- **G4**: GPU instances for graphics workstations
- F1: FPGA instances for hardware acceleration

Purchasing Options

On-Demand Instances

- Pricing: Pay per second (Linux) or per hour (Windows)
- Use Cases: Unpredictable workloads, testing, short-term projects
- Benefits: No commitment, full control

Reserved Instances

- Discount: Up to 72% compared to On-Demand
- Terms: 1 or 3 years
- Payment Options:
 - No Upfront: Pay monthly
 - Partial Upfront: Some upfront payment
 - All Upfront: Maximum discount

- Types:
 - Standard RI: Cannot change instance attributes
 - Convertible RI: Can change instance family, OS, tenancy

Savings Plans

- Discount: Up to 72% for consistent usage
- Flexibility: Across instance families, regions, and services
- Types:
 - o Compute Savings Plans: EC2, Lambda, Fargate
 - EC2 Instance Savings Plans: Specific instance family and region

Spot Instances

- **Discount**: Up to 90% off On-Demand pricing
- Risk: Can be terminated with 2-minute notice
- Use Cases: Fault-tolerant workloads, batch jobs, data analysis
- Spot Fleet: Collection of Spot and On-Demand instances

Dedicated Hosts

- **Use Cases**: Compliance requirements, existing server licenses
- Features: Physical server isolation, socket/core visibility
- Pricing: Most expensive option

Dedicated Instances

- Isolation: No shared hardware with other accounts
- Difference from Dedicated Hosts: No visibility into underlying hardware

Storage Options

Amazon EBS (Elastic Block Store)

Volume Types

1. gp3 (General Purpose SSD)

- Performance: 3,000 IOPS baseline, up to 16,000 IOPS
- Throughput: 125 MiB/s baseline, up to 1,000 MiB/s
- **Size**: 1 GiB 16 TiB
- Use Cases: Boot volumes, general workloads

2. gp2 (Previous Generation)

- Performance: 3 IOPS per GB, burstable to 3,000 IOPS
- Size: 1 GiB 16 TiB

3. io2 Block Express (Provisioned IOPS SSD)

- Performance: Up to 256,000 IOPS, 4,000 MiB/s throughput
- Size: 4 GiB 64 TiB
- Use Cases: Critical applications, large databases

4. st1 (Throughput Optimized HDD)

- Performance: Up to 500 IOPS, 500 MiB/s throughput
- Size: 125 GiB 16 TiB
- Use Cases: Big data, data warehouses

5. sc1 (Cold HDD)

- Performance: Up to 250 IOPS, 250 MiB/s throughput
- Use Cases: Infrequently accessed data

EBS Features

- Multi-Attach: io1/io2 volumes can attach to multiple instances
- Encryption: At-rest and in-transit encryption
- Snapshots: Point-in-time backups stored in S3
- Fast Snapshot Restore: Eliminate initialization latency

Instance Store

- Characteristics: Physically attached to host
- Performance: Very high IOPS and throughput
- Durability: Data lost when instance stops/terminates

• Use Cases: Temporary storage, caching, high-performance computing

Networking:

Elastic Network Interfaces (ENI)

- Components: Private IP, public IP, Elastic IP, MAC address, security groups
- Use Cases: Failover scenarios, network appliances
- Mobility: Can move between instances in same AZ

Enhanced Networking

- SR-IOV: Single Root I/O Virtualization
- Benefits: Higher bandwidth, lower latency, lower CPU utilization
- Types:
 - ENA (Elastic Network Adapter): Up to 100 Gbps
 - Intel 82599 VF: Up to 10 Gbps (legacy)

Placement Groups

- Cluster: Low latency, high throughput within single AZ
- **Spread**: Minimize correlated failures (max 7 instances per AZ)
- Partition: Large distributed workloads (Hadoop, Kafka)

Auto Scaling:

Auto Scaling Groups (ASG)

- Benefits: Automatic scaling, health checks, load balancing integration
- Configuration:
 - Launch Template: Instance configuration blueprint
 - Min/Max/Desired Capacity: Scaling boundaries
 - Health Check Type: EC2 or ELB

Scaling Policies

1. Target Tracking: Maintain target metric value

- 2. Step Scaling: Scale based on CloudWatch alarms
- 3. Simple Scaling: Basic scaling with cooldown periods
- 4. Predictive Scaling: Machine learning-based forecasting

1. Target Tracking Scaling

- Think of it like a thermostat in your AC.
- You set a "target" metric, for example: keep CPU usage at 50%.
- The Auto Scaling system automatically adds or removes servers to keep the metric close to that target.
- Example: If CPU goes above 50%, it launches more instances. If CPU goes below 50%, it reduces instances.
- This is the most common and easiest scaling policy.

2. Step Scaling

- Scaling happens in steps depending on how much the metric crosses a threshold.
- It works with CloudWatch alarms.
- Example:
 - If CPU > 60% \rightarrow add 1 instance.
 - If CPU > 80% → add 2 instances.
 - If CPU < 30% → remove 1 instance.
- This gives you fine-grained control. Instead of reacting the same way
 every time, you scale more when load is very high and less when load is
 just a little high.

3. Simple Scaling

- The basic form of scaling.
- You set a single rule: "When X happens, add/remove Y instances."
- Example: If CPU > 70%, add 1 instance. Then it waits for a cooldown period (say 5 minutes) before making another change.

• It is less flexible than step scaling, because it does not react differently based on how big the problem is—it always takes the same action.

4. Predictive Scaling

- Uses machine learning to look at past usage patterns (like daily or weekly traffic spikes).
- It then predicts future demand and adds/removes instances before the load actually happens.
- Example: If your e-commerce site always gets a traffic spike at 9 PM, predictive scaling will add servers just before 9 PM instead of waiting for CPU to rise.
- This helps prevent delays and ensures smooth performance during predictable high-traffic times.

What is a Cooldown Period?

A **cooldown period** is a waiting time after an Auto Scaling action (like adding or removing an instance) before another scaling action can happen.

- It's like saying: "Wait for a while and see the effect of the last change before making another one."
- Default in AWS is usually 300 seconds (5 minutes), but you can adjust it.

Why Do We Need a Cooldown?

Without a cooldown:

- The system might react too fast to small fluctuations.
- Example: If CPU usage goes slightly above 70%, it adds a new server. Then traffic suddenly drops, and CPU goes low, so it removes a server. This can cause a **ping-pong effect** (constantly adding and removing servers).

With a cooldown:

- After adding a server, Auto Scaling waits a few minutes to see if CPU usage stabilizes.
- This avoids unnecessary, repeated scaling actions.

Storage Services {#storage}

Amazon S3 (Simple Storage Service)

Core Concepts

• Buckets: Containers for objects (globally unique names)

Objects: Files stored in buckets (up to 5TB)

• Keys: Unique identifier within bucket

• Regions: Bucket location affects latency and compliance

Amazon S3 Limits (Key Facts)

Feature	Limit / Value	
Total Storage	Virtually unlimited	
Max Object Size	5 TB	
Upload via Console (single PUT)	5 GB	
Multipart Upload	Required for objects > 5 GB, supports up to 5 TB	
Min Object Size for Infrequent Access / Glacier	128 KB (smaller files are charged as 128 KB)	
Objects per Bucket	Unlimited	
Number of Buckets per Account	100 (soft limit, can be increased to 1000 via AWS Support)	
Max Bucket Size	Unlimited	
Object Key Length (name)	Up to 1024 bytes	
Max Metadata per Object	2 KB	
Max File Uploads per Second	Thousands (virtually unlimited with prefix optimization.	

What is Multipart Upload in S3?

- Multipart Upload is a way of uploading a large file to S3 in smaller parts (chunks) instead of uploading it all at once.
- Each part is uploaded independently and in parallel.

• Once all parts are uploaded, S3 assembles them into a single object.

Why Use Multipart Upload?

- 1. Upload Huge Files: Required for files larger than 5 GB (can go up to 5 TB).
- 2. Faster Uploads: Parts can be uploaded in parallel, making it faster.
- 3. **Resilience:** If one part fails, you only re-upload that part (not the whole file).
- 4. Pause/Resume: Uploads can be paused and resumed later.

How Multipart Upload Works (Steps)

- 1. **Initiate Upload:** You start a multipart upload request → S3 gives you an *Upload ID*.
- 2. **Upload Parts:** You upload file chunks (parts). Each part:
 - Must be between **5 MB and 5 GB** (except the last part).
 - Numbered from 1 to 10,000.
- 3. **Complete Upload:** After all parts are uploaded, you send a *Complete Multipart Upload* request → S3 assembles them into the final file.
- 4. **Abort Upload (optional):** If you don't finish, you can abort and S3 deletes the parts.

How to Do Multipart Upload

You can do it in different ways:

1. Using AWS CLI

```
# Start multipart upload
aws s3api create-multipart-upload --bucket my-bucket --key bigfile.zip

# Upload a part
aws s3api upload-part --bucket my-bucket --key bigfile.zip \
--part-number 1 --body part1.zip --upload-id <UploadId>

# Complete upload
```

```
aws s3api complete-multipart-upload --bucket my-bucket --key bigfile.zip
\
--upload-id <UploadId> --multipart-upload file://parts.json
```

2. Using AWS SDK (Python Example)

Here:

- multipart_threshold = minimum file size to trigger multipart upload (5 MB).
- multipart_chunksize = size of each part (5 MB in this example).

3. Using AWS Management Console

⚠ Note: The **console does not support multipart upload manually**. If a file > 5 GB is uploaded via console, it fails. You must use CLI/SDK/API.

Storage Classes

Example Use Cases

- Uploading large video files (movies, training recordings).
- Backing up huge database dumps.

Uploading machine learning datasets (hundreds of GBs).

Frequently Accessed

1. S3 Standard

• Availability: 99.99%

• **Durability**: 99.99999999% (11 9's)

• Use Cases: General-purpose storage

Infrequently Accessed

1. S3 Standard-IA (Infrequent Access)

• Cost: Lower storage cost, retrieval fee

• Use Cases: Backups, disaster recovery

2. S3 One Zone-IA

• Availability: 99.5% (single AZ)

• Use Cases: Secondary backups, recreatable data

Archive Classes

1. S3 Glacier Instant Retrieval

• Retrieval: Milliseconds

Use Cases: Archive data with immediate access needs

2. S3 Glacier Flexible Retrieval

Retrieval Options:

Expedited: 1-5 minutes

Standard: 3-5 hours

Bulk: 5-12 hours (free)

3. S3 Glacier Deep Archive

• Retrieval: 12-48 hours

• Use Cases: Long-term retention, compliance

4. S3 Intelligent-Tiering

• Automation: Moves objects between tiers based on access patterns

· Cost: Small monitoring fee

How S3 Pricing Works

Amazon S3 charges you based on 4 main factors:

1. Storage Cost (per GB per month)

- Depends on **storage class** (Standard, Infrequent Access, Glacier, etc.).
- Example (US-East region, approximate):
 - S3 Standard: \$0.023 per GB / month
 - S3 Standard-IA (Infrequent Access): \$0.0125 per GB / month
 - S3 Glacier Instant Retrieval: \$0.004 per GB / month

2. Request and Data Retrieval Costs

- Every API request (PUT, GET, LIST, DELETE) costs a tiny amount.
- Example:
 - PUT, COPY, POST, LIST → \$0.005 per 1,000 requests
 - GET, SELECT → \$0.0004 per 1,000 requests
- Infrequent Access and Glacier storage classes have extra retrieval fees when you read data.

3. Data Transfer Out (Bandwidth Costs)

- Free: Data transfer into S3.
- **Charged:** Data transfer out to the internet.
- Example (approximately):
 - First 1 GB/month → Free
 - Up to 10 TB/month → \$0.09 per GB

4. Optional Features

- Versioning: You pay for every version of an object.
- Lifecycle transitions: Moving data between classes may have a small cost.

Replication: Cross-region replication doubles your storage cost.

Advanced Features:

Versioning

- Benefits: Protection against accidental deletion/modification
- MFA Delete: Require MFA for permanent deletion
- Lifecycle Integration: Transition/expire old versions

Cross-Region Replication (CRR)

- Requirements: Source and destination bucket versioning enabled
- Use Cases: Compliance, disaster recovery, reduced latency
- Same-Region Replication (SRR): For compliance and data redundancy

Event Notifications

- Triggers: Object creation, deletion, restoration
- Destinations: SNS, SQS, Lambda, EventBridge
- Use Cases: Automated processing, alerting

Transfer Acceleration

- Mechanism: CloudFront edge locations for upload acceleration
- Benefits: 50-500% faster uploads for global users

Multipart Upload

- Threshold: Recommended for files >100MB, required for >5GB
- Benefits: Improved performance, resumable uploads

Pre-signed URLs

- Purpose: Temporary access to private objects
- Expiration: Configurable (up to 7 days with IAM user credentials)
- Use Cases: Secure file sharing, temporary access

Amazon S3 Transfer Acceleration (S3TA)

What is it?

- Normally, when you upload data to an S3 bucket, it travels directly to the region where the bucket is hosted.
- If you are far from that AWS region (e.g., in India but your bucket is in US-East), uploads can be **slow** due to long network distances.
- Transfer Acceleration speeds this up by using Amazon CloudFront's globally distributed edge locations.

Instead of sending data directly to your S3 bucket's region:

- 1. You upload data to the **nearest CloudFront edge location**.
- 2. AWS uses its **private**, **optimized AWS backbone network** to transfer the data quickly to your S3 bucket's region.

This can make uploads **50–500% faster** for users who are far away from the bucket's region.

Example

- ★ Scenario:
 - You are a video production company in India.
 - Your S3 bucket is hosted in US-East (Virginia).
 - Without acceleration: Uploading a 5 GB video may take 50 minutes due to internet routing delays.
 - With S3 Transfer Acceleration:
 - You upload to the **Mumbai CloudFront edge location**.
 - Data travels on AWS's private high-speed backbone network from Mumbai → Virginia.
 - Upload time might reduce to 15 minutes.
- The farther the distance, the bigger the benefit.

When to Use S3 Transfer Acceleration

For global teams uploading large files (media, backups, data sets).

- When customers or users are spread across multiple countries but data is centralized in one region.
- For applications that require fast, consistent uploads/downloads regardless of user location.

Related Topics Around S3 Acceleration

1. S3 Multi-Part Upload (Different Concept)

- Splits a large file into chunks, uploads in parallel.
- Improves throughput & reliability.
- Used together with acceleration for huge global performance boSQSosts.

2. CloudFront vs Transfer Acceleration

- CloudFront: Used mainly for content delivery (downloads).
- S3 Transfer Acceleration: Used mainly for uploads to S3.

3. Cross-Region Replication (CRR)

- Once data is uploaded, you can replicate it to another region automatically.
- Useful when customers from multiple regions need faster local access.

4. Direct Connect / AWS Global Accelerator

• Alternative solutions for enterprises needing low-latency data transfers.

Costs

- S3 Transfer Acceleration is not free.
- You pay per GB transferred in/out via acceleration, which is slightly more expensive than normal S3 data transfer.
- Best used when speed is more important than cost.

Pre-signed URLs in S3

What is a Pre-signed URL?

By default, S3 objects in private buckets cannot be accessed publicly.

- A Pre-signed URL is a special link generated by AWS that allows temporary access to a private S3 object.
- It is "signed" using the credentials of an IAM user or role.

Key Features

1. Purpose

 Grants temporary access to download (GET) or upload (PUT) an object without making the bucket public.

2. Expiration

- You set the lifetime of the URL when creating it.
- Max validity: 7 days if using IAM user credentials.

3. Security

- Only valid until the expiry time.
- Anyone with the link can access the object until it expires.

Use Cases

1. Secure File Sharing

- Share a private document (e.g., report, image, video) with someone outside your AWS account.
- Example: A teacher uploads exam papers to S3 and shares them with students using pre-signed URLs that expire in 24 hours.

2. Temporary Downloads

- Allow users to download software, invoices, or media files for a limited time.
- Example: E-commerce site sends customers a 48-hour link to download their bill.

3. Temporary Uploads

- Let users upload files to your bucket without needing permanent IAM credentials.
- Example: A customer uploads KYC documents via a secure pre-signed URL that expires in 30 minutes.

Things to Remember

- After expiry, the URL cannot be used again.
- Permissions depend on the IAM user/role that created the URL.
- You can generate GET, PUT, and DELETE URLs.
- Useful for fine-grained, temporary access instead of making a bucket public.

Security Features:

Access Control in Amazon S3

Amazon S3 provides multiple ways to control who can access your data and what they can do with it.

1. Bucket Policies

- Written in JSON format.
- Define permissions at the **bucket level**.
- Example: Allow only a specific IAM role to read/write objects in the bucket.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
        "Effect": "Allow",
        "Principal": {"AWS": "arn:aws:iam::123456789012:role/Developer"},
        "Action": "s3:GetObject",
        "Resource": "arn:aws:s3:::my-bucket/*"
    }
]
```

2. Access Control Lists (ACLs)

- A legacy way to set permissions.
- Can grant permissions at object or bucket level.
- Example: Make a single file readable by the public.
- Best practice → avoid ACLs unless required (use bucket policies instead).

3. Access Points

- A newer feature for large datasets.
- Each access point has its own policy and hostname.
- Example: You can give a research team an access point that only allows them to read certain prefixes (folders) in a data lake.

4. Block Public Access

- Protects buckets from accidental public exposure.
- Can be applied at account level or bucket level.
- Example: Even if a bucket policy allows public access, Block Public Access can override it and keep it private.

Encryption in Amazon S3

To keep data secure, S3 supports encryption at rest and in transit.

1. Server-Side Encryption (SSE)

Data is encrypted after upload by S3.

- SSE-S3 → Amazon manages the keys automatically. (Easiest option)
- SSE-KMS → Uses AWS Key Management Service (KMS). Customers can use AWS-managed or customer-managed keys for better control.
- SSE-C → Customer provides their own keys with every request. AWS never stores the keys.

2. Client-Side Encryption

- The customer encrypts data before uploading to S3.
- Decryption also happens on the client-side before using the data.

• Example: An app encrypts files locally with AES-256 before uploading to S3.

Access Logging in Amazon S3

Logging helps monitor and audit what's happening inside buckets.

1. Server Access Logs

- S3 can generate **detailed logs** about requests to a bucket.
- Logs include requester, bucket name, request time, action (GET/PUT), error codes, etc.
- Useful for security analysis and troubleshooting.

2. AWS CloudTrail

- Provides API-level logging across AWS services, including S3.
- Tracks who made a request, from where, and what actions were taken.
- Example: Detect if someone tries to delete objects or change permissions.

Amazon EFS (Elastic File System)

Amazon EFS is a **serverless**, **fully managed**, **elastic file storage system** for AWS. It provides a shared file system that automatically grows and shrinks as you add or remove files. It is designed to be **highly available**, **scalable**, **and simple to use**.

Core Characteristics

- **Protocol**: Supports **NFSv4.1**, so applications using standard Linux file system commands can access it.
- Scalability: Automatically scales from gigabytes to petabytes without manual intervention.
- Concurrent Access: Thousands of EC2 instances, containers, or onpremises servers can mount it at the same time.
- Availability: By default, it is multi-AZ (region-wide) → highly durable and fault-tolerant.

Performance Modes

1. General Purpose Mode

- Best for low-latency applications.
- Up to 7,000 file operations per second.
- Use cases: Content management systems, home directories.

2. Max I/O Mode

- Supports higher throughput and IOPS, but with slightly more latency.
- Use cases: Big Data, high-performance computing, large-scale workloads.

Throughput Modes

1. Bursting Mode

- · Default mode.
- Throughput scales automatically with file system size.
- Good for most general-purpose workloads.

2. Provisioned Mode

- You can set a fixed throughput independent of file system size.
- Useful when workloads require consistent performance.

3. Elastic Throughput (Newer)

- Automatically scales throughput based on workload demand.
- No need to configure size or throughput in advance.

Storage Classes

- Standard: For frequently accessed data.
- Infrequent Access (IA): For data that is not accessed often but must remain available → saves up to 50% cost.
- Archive: For rarely accessed data (long-term storage at very low cost).

ightharpoonup Lifecycle management policies can automatically move files between Standard ightharpoonup IA ightharpoonup Archive.

Use Cases

1. Content Management

• Web serving, shared content repositories.

2. Web Applications

• Store shared application data across multiple servers.

3. Big Data Analytics

 Shared datasets accessible by multiple compute nodes for parallel processing.

4. Container Storage

Persistent shared storage for EKS, ECS, or Kubernetes pods.

Amazon FSx:

- FSx for Windows File Server → Native Windows/SMB workloads that need Active Directory, NTFS permissions, DFS.
- FSx for Lustre → High-performance/HPC, analytics, ML training; POSIX Linux clients; S3 as a data lake.
- FSx for NetApp ONTAP → Mixed protocols (NFS/SMB/iSCSI), snapshots/clones/dedupe/compression; lift-and-shift NetApp.
- FSx for OpenZFS → Low-latency NFS for Linux/Unix apps; copy-on-write snapshots/clones; dev/test and databases.

FSx for Windows File Server

What it is: Fully managed Windows file shares over SMB, with NTFS permissions, Active Directory join, and DFS namespaces. Multi-AZ HA option, automatic backups, quotas, data deduplication support.

Key capabilities

- SMB/NTFS ACLs, Windows semantics, shadow copies (VSS) support.
- Integrates with AD (managed or self-managed), DFS, antivirus scanners.
- Throughput/IOPS scale with storage & throughput capacity settings.

Short example

You're moving a company's on-prem Windows file server (user home drives and departmental shares) to AWS. Create an FSx for Windows file system, join it to your existing AD domain, define shares and NTFS ACLs, map \files.company.local\DeptShare to Windows EC2 and on-prem clients via VPN/Direct Connect.

Tricky questions AWS may ask

- When do you use **DFS Namespaces** vs multiple FSx shares? What does DFS actually solve?
- How does Multi-AZ failover work, and what happens to SMB sessions during failover?
- Difference between bucket policies/posix vs NTFS ACLs—where do you enforce access?
- How do you migrate at scale (AWS DataSync vs robocopy) and preserve
 ACLs and timestamps?
- When to enable data deduplication and how it impacts performance/space?

FSx for Lustre

What it is: Managed Lustre (POSIX, parallel file system) for HPC/analytics with sub-millisecond latencies and very high aggregate throughput. Deep S3 integration (import/export and data repository tasks).

Key capabilities

- Deployment types: Scratch (temporary, highest performance, no replication) and Persistent (long-lived, replicated).
- Link to S3: lazy import (load on first read), and export results back to S3.
- Massive parallelism for compute clusters (EC2, EKS).

Short example

A rendering farm reads frames for a movie from S3 into FSx for Lustre, computes in parallel on hundreds of Spot instances, writes results back to the file system, and exports final renders back to S3 for archiving and distribution.

Tricky questions AWS may ask

- Scratch vs Persistent: durability, use cases, and cost trade-offs.
- How S3 data repository linking works (lazy load vs explicit import/export).
- Client compatibility (Lustre client on Linux; what about Windows?).
- How do you size **metadata vs data** and plan for stripe count/size?
- Why FSx for Lustre over EFS for analytics workloads?

FSx for NetApp ONTAP

What it is: Managed NetApp ONTAP in AWS with NFS, SMB, and iSCSI. Supports NetApp features: Snapshots, SnapMirror, FlexClone, deduplication/compression, thin provisioning, and tiering to capacity pools.

Key capabilities

- Multi-protocol access (Linux NFS, Windows SMB, block via iSCSI).
- Space efficiency (inline dedupe/compression/compaction).
- Snapshot/replication policies; **SnapMirror** to on-prem or other regions.

Short example

A company running NetApp on-prem migrates to FSx for ONTAP using SnapMirror. Linux app servers mount NFS exports, Windows users access SMB shares, and a database test environment consumes an iSCSI LUN—all from the same FSx for ONTAP deployment.

Tricky questions AWS may ask

- Difference between volume, qtree, and LUN; which protocol uses which?
- How FabricPool tiering works (hot data on performance tier, cold data on capacity tier).
- SnapMirror/SnapVault behavior and RPO/RTO expectations.
- How ONTAP's FlexClone enables near-instant copy for CI/CD and how it saves space.
- Choosing export policy/SMB share permissions vs NTFS/UNIX mode bits where to enforce what?

FSx for **OpenZFS**

What it is: Managed OpenZFS file systems over NFS with copy-on-write snapshots, clones, compression, checksumming, and high IOPS potential.

Key capabilities

- POSIX semantics, low latency, excellent for build farms and DB workloads needing fast NFS.
- Instant, space-efficient snapshots/clones for dev/test.
- Adjustable recordsize, atime behavior, compression (e.g., Iz4), quotas/reservations.

Short example

A CI build farm writes and reads millions of small files. FSx for OpenZFS provides low-latency NFS with ZFS caching; nightly snapshots protect against mistakes, and clones give developers instant test environments.

Tricky questions AWS may ask

- Which NFS versions are supported and when to tune recordsize for databases?
- Differences between **snapshots** and **clones** (space usage, performance).
- How ZFS checksums and copy-on-write affect corruption protection and write amplification.
- Comparing FSx for OpenZFS vs EFS for latency and small-file performance.
- Best practices for **exports**, root squash, and client caching.

Quick chooser (common interview angle)

- Need Windows/SMB + AD/NTFS → FSx for Windows
- HPC/Analytics/ML, S3 data lake, extreme throughput → FSx for Lustre
- Mixed protocols (NFS/SMB/iSCSI) + NetApp data services → FSx for ONTAP
- Low-latency NFS with ZFS features (snapshots/clones) → FSx for OpenZFS

Networking and Content Delivery {#networking}

AWS VPC: Fundamentals to AdvancedComplete Guide

1. What is a VPC? (Basics)

- Definition: A Virtual Private Cloud (VPC) is your own private network in AWS, where you design IP ranges, subnets, routing, and security.
- Analogy: Think of it as building your own mini data center in AWS.
- · You control:
 - Which servers are public vs private
 - ✓ How they talk to each other
 - Which security rules apply
 - How they connect to Internet or On-prem

Keyword to remember:

2. Core Building Blocks

2.1. CIDR Blocks (IP Addressing)

• CIDR defines IP range.

Example: $10.0.0.0/16 \rightarrow 65,536 \text{ IPs.}$

- Divide into subnets → smaller chunks like 10.0.1.0/24.
- You can add secondary CIDRs later, but not resize an existing one.

When to use what:

- Small lab/test: /24 (256 IPs).
- Production enterprise: /16 (large pool, many subnets).
- 💰 Cost: Free defining CIDR has no cost.

2.2. Subnets

Public Subnet → route to Internet Gateway (IGW).

[&]quot;Your private isolated network in AWS."

• **Private Subnet** → no direct Internet; can use NAT for outbound.

When to use:

- Public subnet → Web servers, Load balancers, Bastion hosts.
- Private subnet → Databases, application servers, backend services.

Keyword:

Subnet type = Defined by route table, not AWS label.

2.3. Route Tables

- Default: Local routes (all subnets talk inside VPC).
- Add entries for IGW, NAT, Peering, TGW, VPN.

★ Exam Trap:

A subnet is not public unless it has a route to IGW.

Cost: Free.

2.4. Internet Gateway (IGW)

- · Required for outbound Internet access.
- Inbound allowed only if SG/NACL permit.
- \checkmark Use case → Public apps, websites.
- **Cost:** Free.

2.5. NAT Gateway vs NAT Instance

- NAT Gateway (managed, HA, scalable).
 - AZ-specific → deploy per AZ for HA.
 - 45/hour + per GB processed (~\$32/month baseline + data charges).
- NAT Instance (EC2 manually configured).
 - Cheaper, but must manage HA, scaling.

When to use:

- Prod → NAT Gateway.
- Test/cheap → NAT Instance.

Keyword:

"NAT = Outbound only for private subnets."

2.6. Security Groups (SG)

- Instance-level firewall.
- Stateful → return traffic allowed automatically.
- **Best practice:** Least privilege (start with deny, allow only required ports).
- **6** Cost: Free.

2.7. Network ACLs (NACLs)

- Subnet-level firewall.
- Stateless → must allow both inbound & outbound.

When to use:

- Add extra layer of security at subnet boundary.
- **Cost:** Free.

2.8. Elastic IP (EIP)

- Static IPv4 address.
- Chargeable if **not attached** to running resource.

★ When to use:

- When you need a fixed IP for whitelisting.
- **Solution** Solution

 Solutio

2.9. **VPC Peering**

- Private connectivity between 2 VPCs.
- · Same or different AWS accounts/regions.
- · Not transitive.

When to use:

Few VPCs only → Peering.

- Many VPCs → Transit Gateway.
- **Cost:** Data transfer charged at inter-AZ or inter-region rates.

3. Advanced VPC Features

3.1. VPC Endpoints

- Private connection to AWS services without Internet/NAT.
- Gateway Endpoint: S3 & DynamoDB only. Free.
- Interface Endpoint: Elastic Network Interface (ENI). Paid hourly + per GB.

⊀ When to use:

- Gateway → S3/Dynamo (default choice).
- Interface → API Gateway, SNS, SSM, etc.

Keyword:

"VPC Endpoints = Stay inside AWS network."

3.2. Transit Gateway (TGW)

- Hub-and-spoke model → connects 1000s of VPCs & on-prem.
- Supports route propagation.
- Simplifies vs Peering Mesh.
- 💰 Cost: Per-hour attachment + per-GB data.

When to use:

Enterprises, multi-account/multi-VPC setups.

3.3. VPN & Direct Connect

- VPN: IPsec over Internet. (~\$36/month per connection).
- Direct Connect (DX): Dedicated line. Cheaper per GB, high bandwidth, stable.
- \checkmark Use VPN → Quick, flexible.
- \checkmark Use DX → Production hybrid workloads.

Keyword:

3.4. **VPC Flow Logs**

- Logs IP traffic metadata (not packet data).
- Send to S3/CloudWatch.
- ✓ Use for: Troubleshooting, security audits.
- 👸 Cost: Pay for storage/log ingestion.

3.5. AWS PrivateLink

- Private access to AWS or partner services via ENI.
- Keeps traffic inside AWS backbone.
- ✓ Use for: APIs or SaaS without exposing public endpoints.
- 💰 Cost: Same as interface endpoints (hourly + per GB).

3.6. Multi-account Design

- **Networking account** → owns TGW.
- Security account → logs, Flow Logs, Traffic Mirroring.
- Workload accounts → separate VPCs per team.
- ✓ Use AWS Organizations + SCPs for governance.

4. Common Architectural Patterns

Pattern	Description	When to Use
Public + Private Subnets	Public = NAT/Bastion, Private = App/DB	Standard 3-tier apps
Peering Mesh	Few VPCs peer directly	Small orgs, dev/test
Hub-and-Spoke (TGW)	Central TGW, all VPCs connect	Enterprises, scaling
Hybrid Cloud	VPN or DX + TGW	Extend on-prem to AWS
Microservice Isolation	Each team/service = VPC	Multi- account/microservices

5. Security & High Availability

- Deploy NAT Gateways in each AZ.
- SG default inbound = deny, outbound = allow.
- NACL default inbound/outbound = allow (lock down manually).
- Use **VPC Endpoints** to reduce exposure to Internet.
- For redundancy, use Multi-AZ design for all tiers.

6. Exam Tips & Traps

♦ Remember:

- Public Subnet = Route to IGW.
- Private Subnet = Route to NAT.
- NAT Gateway > NAT Instance (preferred).
- SG = Stateful, NACL = Stateless.
- VPC Peering is NOT transitive.
- Use Gateway Endpoint for S3/Dynamo → Cheaper + safer.
- Flow Logs = Metadata only.
- PrivateLink = Private API access.

7. Cost Optimization Tips

- NAT Gateway is expensive → If only S3/Dynamo needed, use Gateway Endpoint (free).
- Use one NAT Gateway per AZ (not per subnet) to balance cost & HA.
- Direct Connect cheaper than VPN for heavy workloads.
- VPC Flow Logs → Use S3 (cheaper) instead of CloudWatch if logs are massive.
- Use smaller CIDR blocks in test/dev environments to conserve IPs.

8. Keywords to Memorize (Quick Recall)

- IGW → Internet access.
- NAT GW → Outbound Internet only.
- SG → Stateful, instance-level.
- NACL → Stateless, subnet-level.
- VPC Peering → Non-transitive.
- TGW → Hub-and-spoke, scalable.
- Endpoint (Gateway) → S3/Dynamo.
- Endpoint (Interface) → Other services, via ENI.
- PrivateLink → Private service/API access.
- Flow Logs → Metadata only.

1. BEST PRACTICE Private IP Ranges (RFC1918)

- **10.0.0.0/8** → Large address space
- 172.16.0.0/12 → Medium
- **192.168.0.0/16** → Small

★ When to Use

- 10.0.0.0/16 → Large-scale applications needing many subnets (e.g., enterprise apps).
- 192.168.0.0/24 → Small isolated test/dev environments.

2. Subnets

- Public Subnet → Requires IGW (Internet Gateway) + subnet route 0.0.0.0/0 →
 IGW .
- Private Subnet → No direct internet; use NAT Gateway for outbound.

Practice:

Use NAT Gateway (managed, scalable, HA) instead of NAT Instance.

3. Security

- Security Groups → Stateful (allow inbound → response automatically allowed).
- NACLs → Stateless (need explicit inbound + outbound).

When to Use

- SG → Instance-level security.
- **NACL** → Subnet-level extra protection.

4. Connectivity

- VPC Peering → One-to-one, not transitive.
- Transit Gateway → Hub-and-spoke for many VPCs.

Example

- If VPC-A \leftrightarrow VPC-B \leftrightarrow VPC-C \rightarrow Peering won't work for A \leftrightarrow C.
- Use **Transit Gateway** for scale.

5. Endpoints

- Gateway Endpoint → S3, DynamoDB (cheaper, avoids internet/NAT).
- Interface Endpoint (PrivateLink) → Private API access to AWS or partner services.

When to Use

- Gateway Endpoint → Access S3 without internet charges.
- Interface Endpoint → Access SaaS APIs privately.

6. VPC Flow Logs

- Capture metadata only (source/dest, ports, allow/deny).
- No packet payload captured.
- \checkmark Use Case → Troubleshooting, auditing, traffic analysis.

7. Quick Exam Facts

• **Default VPC**: 1 subnet per AZ, route table, SG, NACL included.

- Max VPCs per Region: 5 (soft limit).
- Max CIDR blocks per VPC: 5.
- CIDR change: Cannot be changed after creation (only add secondary).
- Reserved IPs per subnet: 5 (first 4 + last 1).

Summary — When to Use What

- 10.0.0.0/16 → Enterprise / Large workloads.
- 192.168.0.0/24 → Small lab/test.
- Public Subnet → Internet-facing apps (web servers).
- Private Subnet + NAT → App/DB servers needing outbound only.
- NAT Gateway → Default choice (instead of NAT instance).
- Security Groups → For most use cases (stateful).
- NACL → Add extra subnet-level protection.
- VPC Peering → Simple 1:1 VPC link.
- Transit Gateway → Multi-VPC scale.
- Gateway Endpoint → Cost-efficient S3/DynamoDB access.
- Interface Endpoint → PrivateLink for APIs.
- VPC Flow Logs → Monitor traffic (not packet content).

Side-by-Side Comparison

Feature	VPC Peering	Site-to-Site VPN	Direct Connect
Scope	$VPC \leftrightarrow VPC$ (AWS only)	$\begin{array}{c} \text{On-prem} \leftrightarrow \text{AWS} \\ \text{VPC} \end{array}$	On-prem \leftrightarrow AWS (private link)
Network	AWS backbone	Public internet (IPsec tunnel)	Private dedicated line
Encryption	No (but private AWS network)	Yes (IPsec)	Optional (MACsec, or add VPN)
Bandwidth	High (instance- limited)	~1.25 Gbps per tunnel	1–100 Gbps

Feature	VPC Peering	Site-to-Site VPN	Direct Connect
Latency	Low (AWS internal)	Variable, internet- dependent	Very low, predictable
Cost	Data transfer charges only	Low hourly + data transfer	High (port-hour + data transfer)
Setup Complexity	Simple (few VPCs), hard to scale	Easy	Complex (contract, provisioning)
Best Use Case	Connect AWS VPCs	Small/medium hybrid connectivity	Large-scale hybrid connectivity

Amazon VPC (Virtual Private Cloud)

Core Components

CIDR Blocks

- CIDR (Classless Inter-Domain Routing) defines the IP address range of your VPC.
- Allowed private ranges:
 - \circ 10.0.0.0/8 \rightarrow very large (16M IPs)
 - 172.16.0.0/12 → medium (1M IPs)
 - \circ 192.168.0.0/16 → smaller (65K IPs)
- Subnet sizes: from /16 (65,536 IPs) down to /28 (16 IPs).
- AWS reserves 5 IPs in each subnet (first 4 + last 1) for internal use.

Subnets

- A subnet is a slice of the VPC CIDR inside one Availability Zone (AZ).
- Public subnet → Has a route to the Internet Gateway → servers can talk to the internet.
- **Private subnet** → No direct internet → used for internal servers/databases.

Route Tables

- Tell traffic where to go.
- Default routes:

- local → traffic stays within the VPC.
- \circ 0.0.0.0/0 → internet traffic.
- Can add **custom routes** (e.g., send to another VPC, on-premises, etc.).

Internet Connectivity

Internet Gateway (IGW)

- Gateway that connects the VPC to the public internet.
- Attach 1IGW per VPC.
- Must update route tables for subnets to use it.

NAT Gateway

- Lets **private subnet resources** (like databases or backend servers) access the internet **outbound only** (for software updates, API calls).
- Highly available (deploy in each AZ).
- Very fast (up to 100 Gbps).
- Charged per hour + data usage.

NAT Instance (legacy method)

- Old way: run an **EC2 instance with NAT software**.
- · You manage updates, scaling, failover.
- Not recommended now (NAT Gateway replaces it).

VPC Connectivity

VPC Peering

- Connect two VPCs privately.
- Traffic goes over AWS backbone (not internet).
- Limitations:
 - \circ Not transitive (VPC A \leftrightarrow VPC B, B \leftrightarrow C, but A ≠ C).
 - CIDR blocks cannot overlap.

Works across accounts & regions.

Transit Gateway

- Think of it as a hub for networking.
- Connect thousands of VPCs + on-premises data centers.
- Supports route tables, cross-region peering, multicast.
- Use case: large, complex enterprise networks.

Security

Security Groups (SGs)

- Stateful firewall (remembers allowed connections).
- Applied at instance level.
- Only allow rules (no deny).
- Default: block all inbound, allow all outbound.

Network ACLs (NACLs)

- Stateless firewall (doesn't remember connections).
- Applied at subnet level.
- Can have both allow and deny rules.
- Default: allow all.

Monitoring & Logs

VPC Flow Logs

- Capture who is talking to whom inside your VPC (source, destination, ports, etc.).
- Can send logs to S3, CloudWatch, or Kinesis.
- Useful for security analysis & debugging network issues.

Private Service Access

VPC Endpoints

Allow AWS services to be accessed **privately**, without going over the internet.

1. Gateway Endpoints

- Work only for S3 and DynamoDB.
- · Free of cost.
- Add route entries in subnet route table.

2. Interface Endpoints

- Work for most other AWS services (e.g., EC2 API, SSM, Secrets Manager).
- Create an Elastic Network Interface (ENI) with a private IP.
- Cost: hourly + data usage.
- Supports private DNS names (e.g., s3.amazonaws.com resolves to private IP).

Load Balancing:

Application Load Balancer (ALB)

- Layer: 7 (HTTP/HTTPS → Application Layer).
- Features:
 - Path-based routing (e.g., $\frac{1}{|Ap|} \rightarrow \text{one service}, \frac{1}{|App|} \rightarrow \text{another}$).
 - Host-based routing (e.g., shop.example.com → shop servers, blog.example.com
 → blog servers).
 - Query string & header-based routing (e.g., ?version=v2 goes to a different backend).
 - WebSocket & HTTP/2 support.
- Targets: EC2, ECS tasks, Lambda, or IPs.
- SSL/TLS: Offload SSL at the ALB → backend doesn't need to manage certificates.

Use Case:

Best for web applications, microservices, or APIs.

- Example:
 - www.mysite.com/api/* routes to API servers.
 - www.mysite.com/images/* routes to image servers.

Network Load Balancer (NLB)

- Layer: 4 (Transport Layer → TCP/UDP).
- Performance: Handles millions of requests per second with ultra-low latency.
- Static IP: One per Availability Zone.
- Preserves client IP (backend sees the original IP).

Use Case:

- Best for real-time, high-performance workloads.
- Example:
 - Gaming servers (need low latency).
 - IoT device communication.
 - Financial trading systems with high throughput.

Gateway Load Balancer (GWLB)

- **Purpose**: Deploy **third-party virtual appliances** (firewalls, intrusion detection, traffic analyzers).
- **Protocol**: Uses **GENEVE** on port 6081 (tunnels traffic to appliances).

Use Case:

- Best when you want **security inspection** in your traffic.
- Example:
 - Send all traffic through a **third-party firewall** (like Palo Alto or Fortinet) before reaching applications.

Classic Load Balancer (CLB)

Layers: Both Layer 4 and Layer 7.

- Status: Legacy service (not recommended for new workloads).
- Recommendation: Use ALB (if app-layer routing) or NLB (if extreme performance).

Use Case:

- · Older workloads still using it.
- Example: A legacy web app built 10 years ago that hasn't been migrated.

Advanced Features

Sticky Sessions

- Purpose: Send requests from the same client to the same backend target.
- Implementation: Uses cookies.
- Available on: CLB, ALB, NLB.

Example:

 Shopping cart app → ensures user's cart stays on the same server during browsing.

Cross-Zone Load Balancing

- **Distributes traffic evenly across all AZs** (not just within the AZ of the load balancer node).
- ALB: Enabled by default, free.
- NLB & CLB: Disabled by default, extra cost if enabled.

Example:

 If one AZ has 2 targets and another has 4 → without cross-zone, load is uneven. With cross-zone, load is balanced fairly.

Connection Draining (a.k.a Deregistration Delay)

- Purpose: When removing or updating a target, allow existing requests to finish gracefully.
- Timeout: 1–3600s (default 300s).

Example:

 During a rolling deployment → old EC2 instance is removed, but ongoing user downloads finish before the instance is shut down.

🔑 Summary by Use Case

- ALB → Best for web apps, APIs, microservices. (Layer 7 smart routing).
- NLB → Best for extreme performance, gaming, IoT, financial apps. (Layer 4 speed).
- GWLB → Best for security appliances (firewalls, IDS/IPS).
- **CLB** → Use only for **legacy apps** (migrate when possible).

Amazon CloudFront:

Core Features

- Global Network → 400+ edge locations worldwide ensure content is delivered closer to users.
- Content Types:
 - Static content → images, CSS, JS.
 - Dynamic content → APIs, personalized web pages.
 - Live/On-demand streaming → video, audio.
- Origins (where CloudFront pulls content from):
 - **S3** (common for static websites, files).
 - EC2 / ELB (for dynamic content/apps).
 - Custom HTTP server (any web server).
- Caching: Control how long objects are cached with TTL (Time to Live) and cache behaviors (different rules for different paths).

W Example:

 A news website → static images cached for 24h, homepage HTML cached only for 5 mins (so it's always fresh).

Distributions

• Web Distribution → Used for websites, APIs, videos (the standard one).

 RTMP Distribution → Used for Adobe Flash media streaming (deprecated, not recommended anymore).

Price Classes

Helps you **control cost vs performance** by choosing which edge locations to use.

- Price Class All → Uses all edge locations → Best performance (most expensive).
- Price Class 200 → Uses most AWS regions (excludes very expensive ones).
- Price Class 100 → Uses only least expensive regions (lowest cost, higher latency).

Example:

A startup serving mostly US & Europe users → Price Class 100 (ignore expensive Asia/Pacific regions).

Security Features

- Origin Access Control (OAC) → Only CloudFront can access your S3 bucket → users can't bypass CloudFront.
- AWS WAF Integration → Protect against SQL injection, XSS, DDoS at the edge.
- Field-Level Encryption → Encrypt sensitive form fields (like credit card numbers) before they hit your origin.
- Signed URLs/Cookies → Restrict access to paid content or premium videos.

Example:

 Video streaming service → Only paid users can access video → issue signed URLs that expire after 1 hour.

Performance Optimization

 Compression → CloudFront auto-compresses content with gzip or Brotli → faster downloads.

- HTTP/2 support → Better performance with multiplexing & header compression.
- Origin Shield → Adds an extra caching layer between origins & edge → reduces load on origin (great for live events or high-traffic apps).

Example:

Global ecommerce site → reduces repeated origin fetches with Origin
 Shield, improving speed during sales events.

Amazon Route 53

DNS Record Types

DNS records tell the world where your resources live.

- A → IPv4 address mapping.
 - Example: example.com → 192.0.2.44
- AAAA → IPv6 address mapping.
 - Example: example.com → 2001:db8::1234
- CNAME → Canonical name (alias to another domain).
 - Example: www.example.com → app.example.com
- ALIAS (AWS-specific) → Like CNAME, but can be used at root domain.
 - Example: example.com → CloudFront distribution
- MX → Mail servers for email delivery.
 - Example: example.com → mail1.example.com, mail2.example.com
- TXT → Store text (commonly for SPF, DKIM, domain verification).
 - Example: v=spf1 include:_spf.google.com ~all
- NS → Nameservers for delegation.
 - Example: ns-123.awsdns.com, ns-456.awsdns.org
- SOA → Start of authority, contains metadata about the zone.

Routing Policies

These define how Route 53 answers queries when multiple resources exist.

Simple Routing

- One record, one resource.
- · Returns all values in random order.
 - **Example:** example.com → 192.0.2.44 (single web server).

Weighted Routing

- Split traffic between multiple resources by weight.
- Great for A/B testing or gradual deployments.

Example:

- 80% traffic → old app server.
- 20% traffic → new app server.

Latency-Based Routing

Send users to the region with lowest latency.

Example:

- User in India → ap-south-1 EC2.
- User in US → us-east-1 EC2.

Failover Routing

- Active-passive setup (primary + backup).
- Uses health checks → if primary fails, traffic goes to backup.

Example:

- Primary site → us-east-1.
- Backup site → us-west-2.
- If east fails → west takes over automatically.

5 Geolocation Routing

Route based on user's country or continent.

Example:

- Users in Europe → EU servers.
- Users in Asia → Singapore servers.

Geoproximity Routing

- · Like geolocation, but with bias.
- Can "pull" more traffic toward a region.

Example:

- Normally India users → Mumbai.
- Add bias to send 20% of Indian traffic → Singapore (maybe for load testing).

Multivalue Answer Routing

- Simple load balancing.
- Returns multiple healthy IPs (like round-robin DNS).

Example:

- Query → returns 192.0.2.44, 192.0.2.55, 192.0.2.66.
- · Client picks one at random.

Health Checks

- Route 53 can check resource health:
 - Ping HTTP endpoints (e.g., http://app.example.com/health).
 - Combine multiple checks → one "calculated" health.
 - Trigger CloudWatch alarms.
 - For private VPC resources → must use CloudWatch alarms (since R53 health checks can't reach private IPs).

Example:

If /health API returns 500 → Route 53 marks resource unhealthy → routes users elsewhere.

DNS Features

- Hosted Zones: Containers for your domain records.
 - Example: Hosted Zone for example.com contains A, MX, TXT, etc.
- **Private DNS**: Internal DNS for VPC (resolves names like db.internal.local).
- DNS Failover: Works with health checks to switch to backup.
- Traffic Flow: Visual editor → lets you combine multiple routing policies (like weighted + geolocation).

🔑 Extra Insights (Important in Real Use)

1. Alias vs CNAME:

- CNAME = cannot use on root domain.
- Alias = AWS trick → can point root domain directly to S3, ELB, or CloudFront.

2. TTL (Time to Live):

- Controls DNS caching.
- Low TTL = faster changes but more DNS lookups.
- High TTL = cheaper, but slower propagation.

3. Disaster Recovery:

 Commonly done with Failover Routing → one primary region, one backup.

4. Hybrid Cloud:

Private Hosted Zones → resolve AWS + on-premise resources together.

🔽 Quick Example Scenario

Imagine you run myshop.com:

- ALIAS record → points myshop.com to CloudFront.
- Weighted routing → 90% traffic → old version, 10% → new version.
- Latency-based routing → US users → us-east-1, India users → ap-south-1.
- Failover routing → if US server fails, switch to Europe backup.
- Health check → monitors /health endpoint.

WAF + Route 53 → block malicious traffic at DNS level.

Database Services {#databases}

Amazon RDS (Relational Database Service)

Supported Engines

- Amazon Aurora
 - MySQL and PostgreSQL compatible
 - Up to 5x faster than MySQL, 3x faster than PostgreSQL
 - Global Database → cross-region replication (<1s lag)
 - Aurora Serverless v2 → auto-scaling capacity
 - Storage auto-scales up to 128TB
- MySQL → Open-source relational database
- PostgreSQL → Advanced open-source database
- MariaDB → MySQL fork with extra features
- Oracle → Enterprise DB, supports BYOL and License Included
- Microsoft SQL Server → Microsoft enterprise DB

Deployment Options

Multi-AZ Deployments

- Purpose: High availability & disaster recovery
- Mechanism: Synchronous replication to standby instance in a different AZ
- Failover: Automatic (~60–120 sec) to standby if primary fails
- Performance: No read/write scaling (standby is not accessible)
- Exam Tip: Multi-AZ ≠ Read Replicas. Aurora uses replicas instead of standby.

Read Replicas

Purpose: Read scalability & offloading read traffic

- Mechanism: Asynchronous replication
- Cross-Region: Supported (for DR & global apps)
- **Promotion**: Can be promoted to standalone DB (manual)
- Limits: Up to 15 read replicas per source DB
- Chaining: Can create replica of a replica
- Exam Tip: Not automatic failover (unlike Multi-AZ).

Blue/Green Deployments (New Exam Topic #/)

- **Blue** = current production, **Green** = staging copy (kept in sync)
- · Use cases:
 - Major/minor engine upgrades
 - Schema changes & parameter updates
 - Testing before production cutover
- Switchover: Usually <1 min downtime
- Benefits: Safe testing, rollback possible, minimal downtime

Storage Types

- General Purpose SSD (gp2):
 - o 3 IOPS/GB, burstable
 - Max 16,000 IOPS
- General Purpose SSD (gp3):
 - Cheaper than gp2
 - Configure up to 16,000 IOPS & 1,000 MB/s throughput independently
- Provisioned IOPS (io1/io2):
 - Consistent high performance
 - Up to 256,000 IOPS
- Magnetic (Standard): Legacy, low-cost, not recommended
- Max Storage: 64TB for RDS, 128TB for Aurora (auto-scaling)

Backup and Recovery

Automated Backups:

- PITR (Point-in-Time Recovery)
- Retention: 1–35 days
- Stored in S3 (high durability)
- Deleted when DB is deleted (unless final snapshot taken)

Manual Snapshots:

- o User-initiated, kept until deleted
- Can be copied across regions/accounts
- Encrypted snapshots require same KMS key

• Restore Process:

Always creates a new DB instance

• Exam Tip:

- Backups don't affect performance (storage-level snapshots).
- To encrypt an existing unencrypted DB → snapshot → copy with encryption → restore.

Security Features

Encryption at Rest:

- Managed via KMS
- Applies to storage, backups, and snapshots
- Must be enabled at creation (or via encrypted snapshot restore)

• Encryption in Transit:

SSL/TLS required

Network Isolation:

- Always runs in **VPC**
- Controlled via security groups
- IAM Database Authentication (MySQL & PostgreSQL):

- Uses temporary IAM tokens (15 min)
- Removes need for password management
- Transparent Data Encryption (TDE):
 - Supported for Oracle & SQL Server
- Audit & Compliance:
 - Integrated with CloudTrail, AWS Config, GuardDuty

Monitoring and Maintenance

- CloudWatch Metrics: Standard DB performance metrics
- Enhanced Monitoring:
 - OS-level metrics (1s, 5s, 10s intervals)
- Performance Insights:
 - Query-level analysis
 - Retention: 7 days (free) or up to 2 years (paid)
- Maintenance Windows:
 - AWS applies patches/upgrades
 - Auto minor version upgrades (optional)
 - Major upgrades = manual
 - Some maintenance events trigger failover

Exam Tips & Gotchas

- Multi-AZ = HA/DR, not scaling
- Read Replicas = scaling, not HA (unless promoted)
- Blue/Green Deployments = safe upgrade/testing with minimal downtime
- Aurora Global Database = best for cross-region DR (<1s lag)
- Snapshots & backups inherit encryption state
- Cannot encrypt an existing unencrypted DB directly
- RDS integrates with KMS, IAM, VPC, CloudTrail, Config

• Use **RDS Proxy** to manage DB connections for serverless/Lambda apps

Amazon Aurora

Architecture Benefits

- Cloud-Native Design: Built specifically for AWS cloud.
- Performance:
 - Up to 5x faster than MySQL
 - Up to 3x faster than PostgreSQL
- Storage:
 - Auto-scales in **10GB increments** up to **128TB**.
 - Shared, distributed, fault-tolerant storage.
- High Availability:
 - Maintains 6 copies of data across 3 AZs.
 - Data automatically self-heals (continuous background checks).
 - Designed to survive 2 AZ failures.
- Fast Recovery: Point-in-time recovery (PITR) down to the second.
- **Durability**: 99.99% availability target.

Aurora Replicas

- Aurora Replicas (for read scaling):
 - Up to 15 replicas.
 - Same replication engine as Aurora, faster than RDS read replicas.
 - Replication latency typically < 10 ms.
- MySQL Read Replicas: You can also create up to 5 standard MySQL read replicas from Aurora.
- Failover: Automatic failover to Aurora Replica in ~30 sec.

Aurora Serverless (v1 & v2)

Use Cases:

- Infrequent, intermittent, or unpredictable workloads.
- Development, test, and new applications.

Scaling:

- Aurora Serverless v1 → scales in "capacity units" (ACUs), can take seconds to minutes.
- Aurora Serverless v2 → scales in fine-grained increments instantly, production ready.
- Cost Model: Pay per second for ACUs consumed.
- Connectivity:
 - Uses Data API → connect over HTTPS without managing persistent connections (great for Lambda & serverless apps).

Aurora Global Database

- Purpose: Global-scale apps with low-latency reads & DR.
- Replication:
 - Dedicated replication infrastructure → <1 sec typical lag.
 - Cross-region replication with **no performance impact on primary**.
- Read Replicas: Up to 16 read replicas per secondary region.
- Failover (DR):
 - RPO (Recovery Point Objective) < 1 sec
 - RTO (Recovery Time Objective) < 1 min
- Exam Tip: Best choice for disaster recovery across regions.

Aurora Machine Learning

- Integration:
 - Amazon SageMaker (custom ML models).
 - Amazon Comprehend (NLP).

- **Use Cases**: Fraud detection, sentiment analysis, personalized recommendations.
- **Execution**: Call ML models directly from **SQL functions** without app-level integration.

Security

- Encryption:
 - At rest (KMS) & in transit (SSL/TLS).
 - Snapshots and backups inherit encryption state.
- IAM Authentication: IAM tokens (15 min validity).
- VPC Only: Always deployed inside a VPC.

Monitoring

- Performance Insights: Deep query-level monitoring.
- Enhanced Monitoring: OS-level metrics.
- CloudWatch: Basic performance metrics.
- Audit: Integrated with CloudTrail.

Unique Aurora Features vs RDS

- · Auto-healing distributed storage layer.
- Replication is faster and more reliable (built-in).
- · Serverless scaling option.
- Aurora Global Database for cross-region replication.
- Supports up to 15 replicas (vs 5 for RDS).
- Storage auto-scales to 128TB (RDS max is 64TB).

Exam Tips & Gotchas

- Aurora multi-AZ is built-in (6 copies across 3 AZs).
- Aurora Replicas provide both read scaling and HA failover.

- Aurora Serverless = pay-per-use DB (best for intermittent workloads).
- Aurora Global Database = best for low-latency global apps & DR.
- For serverless apps (e.g., Lambda), use Aurora Serverless + Data API.
- If exam mentions RPO < 1s and RTO < 1min across regions → Aurora Global DB.

Amazon DynamoDB

Core Characteristics

- Type: Fully managed NoSQL key-value & document database.
- Performance: Single-digit millisecond latency at scale.
- Scalability: Virtually unlimited throughput and storage.
- Availability: Multi-AZ by default with automatic failover.
- Durability: Data replicated across 3 AZs in a region.

Data Model

- **Tables** → Collections of items.
- Items → Records (max size 400KB).
- Attributes → Properties of an item.
- Primary Keys:
 - Partition Key only → Simple primary key.
 - Partition + Sort Key → Composite primary key (enables range queries).
- Indexes:
 - Local Secondary Index (LSI) → Same partition key, different sort key, created at table creation (max 5).
 - Global Secondary Index (GSI) → Different partition & sort keys, created anytime, supports query flexibility.

Capacity Modes

Provisioned Mode

- Allocate Read Capacity Units (RCUs) & Write Capacity Units (WCUs).
- Auto Scaling available (adjusts based on demand).
- Best for predictable workloads.
- Cost-efficient if traffic is stable.

On-Demand Mode

- Pay-per-request (read/write requests).
- Instantly scales to handle unpredictable workloads.
- More expensive per request vs provisioned.
- Best for spiky or unpredictable traffic.

Consistency Models

- Eventually Consistent Reads (default, lower cost).
- Strongly Consistent Reads (must be specified, higher latency & cost).

Global Tables

- Multi-region, Active-Active replication.
- Writes in one region are replicated to others.
- Conflict resolution = last writer wins.
- Use Cases:
 - Low-latency global apps.
 - Disaster recovery (multi-region HA).

DynamoDB Streams

- Purpose: Capture item-level data changes (insert, update, delete).
- Retention: 24 hours.
- Integration: Can trigger Lambda or be consumed by Kinesis.
- Use Cases: Replication, auditing, event-driven apps.

Advanced Features

DynamoDB Accelerator (DAX)

- **In-memory cache** → microsecond latency.
- Compatible with existing DynamoDB API calls.
- No app code changes required.
- Best for read-heavy workloads or hot keys.

Transactions

- ACID transactions across multiple items/tables.
- Ensures all-or-nothing consistency.
- Use Case: Financial apps, inventory systems.

Time-to-Live (TTL)

- Expire items automatically (background deletion).
- · Useful for session data, temporary records.

Backup and Restore

- On-Demand Backups: Full manual backups (kept indefinitely).
- Point-in-Time Recovery (PITR): Continuous backups with up to 35 days recovery window.
- Cross-Region Backup: Via AWS Backup integration.

Security

- Encryption:
 - At rest (default, KMS-managed).
 - In transit (TLS).
- Fine-Grained Access Control:
 - IAM policies down to item and attribute level.

- VPC Endpoints:
 - Private access without internet exposure.

Monitoring

- CloudWatch Metrics: RCU/WCU usage, throttled requests, latency.
- CloudTrail: Logs all API calls.
- Contributor Insights: Identify hot keys & traffic patterns.

Exam Tips & Gotchas

- Max item size = 400 KB.
- LSI must be created at table creation; GSIs can be added later.
- Global Tables = multi-region active-active (unlike RDS read replicas).
- **DAX** = caching (not persistence).
- PITR is disabled by default (must enable).
- If workload is spiky/unpredictable → On-Demand Mode.
- If workload is steady → Provisioned + Auto Scaling.
- DynamoDB does not support joins → denormalize data & use indexes.
- Hot Partitions (uneven key distribution) can throttle performance → use good partition key design.

Amazon ElastiCache

Engines

Redis

- Features:
 - Advanced data structures (lists, sets, sorted sets, hashes).
 - Pub/Sub messaging.
 - Persistence options (AOF, RDB snapshots).
 - Clustering with sharding.

High Availability:

- Multi-AZ with automatic failover (via Redis replication groups).
- Supports failover in < 30s.

• Backup & Restore:

- Manual & automatic snapshots supported.
- Encryption at-rest & in-transit supported.

Use Cases:

- Session storage (shopping carts, logins).
- Real-time leaderboards.
- Pub/Sub for messaging systems.
- Caching DB queries & API responses.

Memcached

Features:

- Simple key-value store.
- Multi-threaded architecture.
- No persistence, no replication, no failover (must rehydrate on restart).

Scaling:

- Horizontal scaling → Add/remove nodes.
- Supports sharding client-side.

• Use Cases:

- Simple caching layer.
- Database query caching.
- Object caching (e.g., HTML fragments).

Caching Strategies

· Lazy Loading:

Load into cache only when requested.

- Pro: Only requested data cached.
- Con: Cache miss penalty on first request.

Write-Through:

- Write to both cache & DB simultaneously.
- Pro: Always fresh cache.
- Con: Higher write latency.

• TTL (Time-To-Live):

- Expire stale data automatically.
- Reduces cache bloat & ensures data freshness.

ElastiCache Use Cases

- Session Store: Sticky, fast session handling (Redis preferred).
- Database Caching: Reduce RDS/DynamoDB load.
- Real-Time Analytics: Sub-millisecond retrieval for dashboards.
- Gaming Leaderboards: Sorted sets in Redis for rankings.
- Pub/Sub Messaging: Lightweight message queues with Redis.



Other AWS Database Services

Amazon DocumentDB

- Compatibility: MongoDB API-compatible (not MongoDB itself).
- Scaling: Up to 15 read replicas.
- Storage: Automatically grows in 10GB increments, up to 64TB.
- Durability: Replicates 6 copies across 3 AZs (similar to Aurora).
- Use Cases:
 - Content management systems.
 - User profiles.

Product catalogs.

Amazon Neptune

- Type: Graph database (fully managed).
- Query Languages:
 - Gremlin (property graph).
 - SPARQL (RDF).
- High Availability:
 - 6-way replication across 3 AZs.
 - Up to 15 read replicas.
- Use Cases:
 - Social networks (friendship graphs).
 - Fraud detection (relationship analysis).
 - Recommendation engines.
 - Knowledge graphs.

Amazon Keyspaces (for Apache Cassandra)

- Type: Serverless Cassandra-compatible DB.
- Scalability: Virtually unlimited, auto-scaling.
- High Availability: Multi-AZ by default.
- Use Cases:
 - o IoT time-series data.
 - Sensor logging.
 - User activity tracking.
- Notes:
 - Query model limited to partition key design.
 - Supports CQL (Cassandra Query Language).

Amazon QLDB (Quantum Ledger Database)

- Type: Immutable ledger database.
- Features:
 - Transparent, cryptographically verifiable transaction log.
 - Append-only journal → No deletion/update of history.
 - Provides full history of data changes.
- Scaling: Serverless, fully managed.
- Use Cases:
 - Financial transactions.
 - Supply chain tracking.
 - Vehicle history, compliance, auditing.
- Exam Tip: If question mentions immutable, verifiable history → QLDB (not DynamoDB).

🔽 Exam Tips & Gotchas

- · Redis vs Memcached:
 - Redis = advanced features, persistence, HA, clustering.
 - Memcached = simple, fast, ephemeral cache only.
- Session Store → Redis (because of persistence + HA).
- Leaderboard / Ranking → Redis (sorted sets).
- MongoDB workloads → DocumentDB (but not open-source Mongo).
- Graph relationships → Neptune.
- Cassandra workloads → Keyspaces.
- Immutable ledger with verification → QLDB.
- If question says "cryptographically verifiable history of changes" → always QLDB.
- If exam says "multi-region active-active database" → DynamoDB Global Tables (not ElastiCache).

Security and Compliance {#security}

AWS Key Management Service (KMS)

Key Types

- AWS Managed Keys (AWS/service-managed CMKs)
 - Automatically created and managed by AWS services (e.g., S3, EBS).
 - Minimal visibility for customers.
 - Lifecycle handled by AWS.
- Customer Managed Keys (CMKs)
 - Full control by customer.
 - You define key policies, IAM permissions, rotation, lifecycle.
 - Can enable automatic key rotation (1 year).
 - Can use for cross-account encryption.

AWS Owned Keys

- Completely abstracted.
- Used internally by AWS to encrypt things like CloudWatch Logs, S3
 objects in services you don't explicitly configure.
- Not visible to you.

Key Policies

- Default Policy:
 - Grants full access to the root user.
 - Recommended to delegate access using IAM roles/policies instead of root.

Custom Policies:

- JSON-based fine-grained control.
- Example: Allow only a specific IAM role to use a key for Encrypt but not Decrypt.

Cross-Account Access:

- Grant another AWS account access via key policies + IAM role.
- Common exam scenario: "App in Account A needs to decrypt data from Account B" → Use KMS cross-account key policy.

Encryption Operations

• Envelope Encryption:

- Data is encrypted with a **Data Key** (DEK).
- DEK is encrypted with a Customer Master Key (CMK).
- Benefits: Large data encrypted locally → only DEK sent to KMS.

• Data Key Generation:

- GenerateDataKey → Returns plaintext + encrypted key.
- Application encrypts data with plaintext key, then discards it, and stores only ciphertext key.

• Multi-Region Keys:

- New feature → Replicate keys across regions.
- Same key material & key ID across multiple regions.
- Useful for multi-region apps, DR, cross-region replication.

Integration

- AWS Services: Over 100+ integrations (S3, EBS, RDS, Redshift, Lambda, Secrets Manager).
- Client-Side: AWS Encryption SDK (Java, Python, JS).
- Performance: KMS uses HSMs (FIPS 140-2 validated) under the hood.

Other Features

- · Key Rotation:
 - AWS Managed CMKs: Auto rotation every 3 years.
 - Customer Managed CMKs: Optional auto rotation every 1 year.

- **Deletion**: Keys scheduled for deletion (7–30 days, default 30). **Immediate** deletion not allowed.
- Grant Tokens: Temporary access delegation to a CMK.

Exam Tips for KMS

- For cross-account access → Modify key policy.
- For large files → Use envelope encryption, not direct Encrypt.
- If exam says "must ensure compliance with regulatory control of keys" →
 Consider CloudHSM over KMS.
- If keys must be used across regions → Multi-Region KMS keys.



AWS CloudHSM

Overview

- Fully managed Hardware Security Module (HSM).
- Dedicated, single-tenant appliance per customer (unlike shared KMS).
- You manage key material entirely → AWS cannot access.
- FIPS 140-2 Level 3 validated.

Use Cases

- Regulatory compliance (PCI DSS, HIPAA, banking, government).
- Custom key management (import your own keys, control lifecycle).
- Applications needing low-latency, high-volume cryptographic operations.

Deployment

- Cluster Architecture: Multi-AZ, resilient HSMs.
- Client Integration: Apps integrate using CloudHSM SDK (PKCS#11, JCE, CNG).
- **Backup**: Automatic daily backups to Amazon S3 (encrypted).

Exam Tips for CloudHSM

- · CloudHSM vs KMS:
 - Use CloudHSM when you need full control of keys or meet compliance.
 - Use KMS for managed, integrated encryption.
- Exam often says: "Org needs to manage encryption keys, AWS staff must not have access" → Answer = CloudHSM.



AWS Certificate Manager (ACM)

Certificate Types

- Public Certificates:
 - Issued by AWS's public CA (Amazon Trust Services).
 - Used for public websites.
 - Free of charge.
- Private Certificates:
 - Issued by ACM Private CA (paid).
 - Used for internal apps & corporate intranet.
- Imported Certificates:
 - Upload 3rd-party certs into ACM.
 - You manage renewal manually.

Validation Methods

- DNS Validation (recommended):
 - Add a CNAME record to Route 53 (or external DNS).
 - Fully automated renewal.
- **Email Validation:**

- Sends approval emails to domain's WHOIS or admin contacts.
- Manual → less recommended.

Integration

- Works natively with:
 - Load Balancers (ALB, NLB, CLB).
 - CloudFront (must be in us-east-1).
 - API Gateway (REST, HTTP, WebSocket APIs).
 - App Runner & Elastic Beanstalk.

Automation

- Auto-Renewal:
 - ACM auto-renews public & private certs (if DNS validation is configured).
 - Imported certs must be renewed manually.
- Notifications:
 - CloudWatch Events or EventBridge → Notify before expiry.

Exam Tips for ACM

- For public websites → Always use ACM public certs.
- CloudFront with ACM → Cert must be issued in us-east-1.
- If compliance requires internal PKI → Use ACM Private CA.
- If external CA certs already exist → Import into ACM.

Quick Comparison Table

Feature	KMS	CloudHSM	ACM
Purpose	Managed key mgmt + encryption	Dedicated HSM for full control	SSL/TLS certificates

Feature	KMS	CloudHSM	ACM
Who Manages Keys	AWS (with CMK control for customers)	Customer fully	ACM / You (if imported)
Compliance	FIPS 140-2	FIPS 140-2 Level 3	Industry SSL/TLS standards
Cost	Pay per API request	Dedicated, costly	Public = free, Private CA = paid
Use Case	Encrypt data in AWS services	Regulated industries, custom PKI	SSL/TLS for apps/websites

AWS WAF (Web Application Firewall)

Rule Types

- IP Match → Allow/block based on IP (IPv4/IPv6).
- Geo Match → Block/allow requests from specific countries.
- String/Regex Match → Match patterns in headers, URIs, query strings.
- **SQL Injection (SQLi)** → Detect & block SQL injection attempts.
- Cross-Site Scripting (XSS) → Detect & block malicious scripts.
- Rate-Based Rules → Limit requests per IP (e.g., 1000 reqs/5 min).

Managed Rule Groups

- AWS Managed Rules → Preconfigured rules (e.g., SQLi, XSS).
- Marketplace Rules → Security vendors (Imperva, F5, etc).
- Custom Rules → User-defined.

Integration

- CloudFront → Global edge protection (best for CDN/websites).
- Application Load Balancer (ALB) → Regional app-level protection.
- API Gateway → REST & WebSocket API protection.
- AppSync → Protect GraphQL APIs.

★ Exam Tip:

If the attack is Layer 7 (app-level) \rightarrow use WAF.

AWS Shield

Shield Standard

- Cost → Free, enabled by default.
- Protection → Automatic mitigation of common L3/L4 DDoS attacks (SYN/UDP floods).
- Coverage → CloudFront, Route 53, ELB.

Shield Advanced

- **Cost** → \$3,000/month per org.
- Enhanced Protection:
 - Against sophisticated DDoS attacks.
 - Covers EC2, ELB, CloudFront, Global Accelerator, Route 53.
 - 24/7 DRT Support (AWS DDoS Response Team).
 - DDoS Cost Protection → Credits for scaling costs during attacks.
 - Detailed Attack Diagnostics.

Exam Tip:

If the scenario says "mission-critical app, must prevent expensive scaling during $DDoS" \rightarrow$ Shield Advanced.

Amazon Guard Duty

Detection Capabilities

- Machine Learning + Threat Intel → Detect anomalies, malicious IPs/domains.
- Data Sources:
 - CloudTrail → Suspicious API calls (IAM misuses, privilege escalations).
 - VPC Flow Logs → Unusual network traffic.
 - DNS Logs → Suspicious domain queries.

Finding Types

- Reconnaissance → Port scanning, brute force login.
- Compromised EC2 → Malware, crypto-mining, unusual outbound traffic.
- **S3 Compromises** → Unauthorized access, data exfiltration.

Integration

- CloudWatch Events/EventBridge → Trigger alerts & automation.
- Lambda → Automatic remediation (e.g., isolate EC2).
- **Security Hub** → Aggregate findings.
- **∳** Exam Tip:

If the question says "detect compromised EC2 / unusual IAM activity / S3 $exfiltration" \rightarrow GuardDuty$.

AWS Config

Configuration Management

- **Resource Inventory** → Tracks all AWS resources and configurations.
- Change History → Timeline of config changes.
- Compliance Auditing → Evaluate against compliance standards (CIS, HIPAA, PCI).

Rules

- AWS Managed Rules → Prebuilt (e.g., "S3 buckets must not be public").
- Custom Rules → Built using Lambda.
- Remediation → Auto-fix non-compliance (e.g., re-enable encryption).

Use Cases

- Auditing → Compliance reporting.
- **Security Analysis** → Detect misconfigured resources.
- Troubleshooting → Root cause of config drift.
- **∳** Exam Tip:

If the question says "must track changes to resources & enforce compliance continuously" \rightarrow **AWS Config**.

AWS Secrets Manager

Features

- Automatic Rotation → Credentials auto-rotated (e.g., RDS every 30 days).
- Cross-Region Replication → For DR and multi-region apps.
- **Fine-Grained IAM Access** → Attribute-level permissions.
- **Encryption** → Always encrypted with KMS.

Integrations

- Direct with → RDS, DocumentDB, Redshift.
- Can also store → API keys, app credentials, OAuth tokens.

Use Cases

- Database Credentials → Rotate RDS passwords automatically.
- API Keys → Manage third-party service access.
- Application Secrets → Store config/connection strings.

∲ Exam Tip:

If the question says "must securely store & rotate database credentials automatically" \rightarrow Secrets Manager.

(\triangle Don't confuse with **SSM Parameter Store** \rightarrow cheaper, basic storage, but no auto-rotation).

© Quick Exam Decision Guide

- WAF → Block SQLi/XSS/Bad bots (L7 attacks).
- Shield → Protect against DDoS (L3/L4).
- GuardDuty → Detect compromised resources (EC2, IAM, S3).
- Config → Continuous compliance & resource change tracking.

Secrets Manager → Secure storage + rotation of secrets.

AWS Systems Manager Parameter Store

Parameter Types

- String → Plain text values (e.g., DB endpoint, API URL).
- StringList → Comma-separated values (e.g., list of subnets or security group IDs).
- SecureString → Encrypted values using AWS KMS (ideal for passwords, tokens, API keys).

Parameter Tiers

- Standard Tier
 - Up to 10,000 parameters.
 - Max 4 KB size per parameter.
 - Free.
- Advanced Tier
 - Up to 100,000 parameters.
 - Max 8 KB size per parameter.
 - Supports parameter policies (e.g., TTL expiration, no-change alert).
 - Extra cost.

Use Cases

- Configuration Management → Store app config values like DB names, feature flags.
- Secrets Storage → Simple secrets (use Secrets Manager for rotation/complex secrets).
- **License Keys** → Store product keys or API licenses.
- **Environment Variables** → Store per-environment configs (dev, test, prod).

Access & Security

- Integrated with IAM Policies (fine-grained access control).
- For <u>SecureString</u>, you can specify a custom KMS CMK (instead of AWS-managed default key).
- Parameters can be versioned (you can roll back to older versions).

Integration

- Works natively with AWS SDKs, CLI, and CloudFormation.
- Can be pulled into EC2 User Data, Lambda, ECS tasks, and CodePipeline.
- Supports cross-account and cross-region access (with resource policies).

Exam Tips

- Secrets Manager vs Parameter Store:
 - Use Secrets Manager for automatic rotation & cross-service secret integrations.
 - Use Parameter Store for simple config + lightweight secrets.
- Standard = 10,000 params (4 KB each, free).
- Advanced = 100,000 params (8 KB each, policies, paid).
- SecureString requires **KMS** for encryption.

Monitoring and Management {#monitoring}

Amazon CloudWatch

Metrics

- Default Metrics
 - Provided automatically for most AWS services (e.g., EC2 → CPU Utilization, Disk Reads/Writes, NetworkIn/Out).
 - Granularity = 5 minutes (basic monitoring).
- Detailed Monitoring

- Granularity = 1 minute.
- Extra cost (per metric per instance).
- Required for Auto Scaling policies to react faster.

Custom Metrics

- Push your own app-level metrics via CloudWatch API or Agent.
- Resolution: 1 second (high-resolution custom metrics).
- Example: request latency, memory usage (since EC2 memory isn't default).

Metric Math

- Perform math on metrics (e.g., CPUUtilization across multiple instances).
- Useful for creating aggregate metrics.

∳ Exam Tip:

Memory, Disk space, and Swap usage are NOT default metrics for EC2.
 You need CloudWatch Agent for them.

Logs

- Log Groups: Container for log streams.
- Log Streams: Sequence of log events (per instance, Lambda, or container).
- Log Retention: Configurable (1 day → 10 years, or never expire).
- Log Insights: SQL-like query engine for analyzing logs.
- Subscription Filters: Stream logs to Lambda, Kinesis, or Elasticsearch for real-time processing.

★ Exam Tip:

- If asked about analyzing logs in real-time → use CloudWatch Logs + Subscription Filters.
- If asked about searching historical logs → use CloudWatch Logs Insights.

Alarms

Metric Alarms → Trigger on a single metric crossing a threshold.

- Composite Alarms → Combine multiple alarms using Boolean logic (AND/OR).
- Actions:
 - Notify via SNS.
 - Scale via Auto Scaling.
 - Recover EC2 instance (for some metrics).

♦ Exam Tip:

• EC2 Auto-Recovery works with alarms monitoring system health checks.

Events (Amazon EventBridge)

- Formerly CloudWatch Events.
- Event Sources: AWS service events, custom events, SaaS integrations.
- Rule Types:
 - Event pattern (e.g., EC2 state change).
 - Schedule (cron/interval).
- Targets: Lambda, Step Functions, SNS, SQS, Kinesis, etc.
- Cross-Account: Can forward events across AWS accounts.

Exam Tip:

If requirement mentions decoupled event-driven architecture →
 EventBridge is the right choice.

Dashboards

- Widgets: Graphs, numbers, logs, text.
- Sharing: Can be shared (read-only) publicly.
- Auto-Refresh: Real-time updates.

AWS CloudTrail

Event Types

- Management Events (Control Plane): Create/Delete/Update API calls (default enabled).
- Data Events (Data Plane): S3 object-level (GetObject, PutObject), Lambda invoke events (must be enabled separately, cost \$\$\$).
- Insight Events: Detect unusual activity (e.g., sudden increase in API calls).

Trail Configuration

- **Single Region**: Records events in one region.
- Multi-Region: Recommended for auditing (default for org trails).
- Organization Trail: Centralized logging for all AWS accounts in an Org.

Log File Integrity

- Signed with SHA-256 hash digest.
- Used to detect tampering.

Integration

- Send to CloudWatch Logs for near real-time monitoring.
- Store in **S3** for long-term retention.
- Query via Athena (serverless SQL queries).

∳ Exam Tip:

- CloudTrail = Who did what, when, and from where.
- If question is about API activity auditing → CloudTrail.
- If question is about **performance metrics** → CloudWatch.

AWS X-Ray

Distributed Tracing

- Helps debug microservices & serverless applications.
- Generates a **service map** to visualize dependencies.
- Tracks end-to-end latency and errors across multiple services.

Integration

- · Works with:
 - Lambda (tracing enabled by flag).
 - o API Gateway.
 - ECS/EKS.
 - Elastic Beanstalk.

Sampling

- Fixed Rate → e.g., sample 10% of requests.
- Rate and Burst → e.g., capture first N requests, then sample.
- Custom Rules → service-specific.

∳ Exam Tip:

- If question mentions tracing user requests through multiple microservices
 → AWS X-Ray.
- If question mentions monitoring application performance metrics → CloudWatch.

Quick Exam Cheat-Sheet

- CloudWatch Metrics/Logs → Performance + logs.
- CloudWatch Alarms → Automated actions on threshold breach.
- EventBridge (CloudWatch Events) → Event-driven automation.
- CloudTrail → API auditing (security + compliance).
- X-Ray → Tracing microservices (debugging).

AWS Systems Manager (SSM)

AWS Systems Manager is a unified service for **operations**, **configuration**, **application**, **and resource management** across AWS and on-premises environments.

Operations Management

Session Manager

- Secure, browser-based shell/CLI access to EC2 and on-prem servers (no SSH keys needed).
- Can log session activity to CloudWatch Logs / S3.

Run Command

- Execute remote shell commands or scripts across multiple instances.
- Useful for ad-hoc tasks (e.g., install software, restart service).

Patch Manager

- Automates patching of **OS & applications** across EC2 and on-prem servers.
- Can define patch baselines (approved/denied patches).

Maintenance Windows

- Define recurring schedules for maintenance tasks.
- Ensures automation tasks only run during approved windows.

Configuration Management

• Parameter Store

- Centralized storage for config data and secrets.
- Types: String , StringList , SecureString (KMS encrypted).
- Tiers:
 - Standard: 10,000 params, 4KB each.
 - Advanced: 100,000 params, 8KB, parameter policies.
- Use cases: Config management, secrets, license keys.

State Manager

- Enforces consistent system configurations.
- Example: Ensure antivirus always installed/enabled, IAM role attached,
 OS settings.

Compliance

- Continuously checks if resources comply with desired state.
- Integrates with AWS Config for compliance reporting.

Application Management

- Application Manager
 - Provides a single view of applications (resources, logs, monitoring).
 - Integrates with CloudFormation, AWS Config, and tagging.

AppConfig

- Deploys application configurations safely and quickly.
- Example: Feature flags, config toggles without redeploying app.
- Supports rollback if issues detected.

OpsCenter

- Centralized console for operational issues (OpsItems).
- Aggregates issues from CloudWatch, AWS Config, Trusted Advisor.

Extra Notes (Exam Tips)

- SSM Agent must be installed on instances (preinstalled on Amazon Linux 2 & Windows AMIs).
- Systems Manager works across multiple AWS accounts & regions.
- Integrates with CloudWatch Events/Alarms for automation triggers.
- Can manage on-prem servers via Hybrid Activations.
- Prefer Secrets Manager for rotation-heavy or complex secrets, but
 Parameter Store is cheaper for simple secrets.

Application Integration {#integration}



- Think of SQS like a post office for software systems:
 - Producers drop messages into the mailbox (queue).
 - Consumers pick them up later.
 - The sender and receiver don't talk directly this makes systems decoupled (they don't wait on each other).

This is **super important in AWS** because you want scalable, fault-tolerant apps where if one part is slow/down, the rest still works.

Queue Types

1. Standard Queue (default type)

- Unlimited throughput → You can send/receive as many messages per second as you want.
- **Best-effort ordering** → Usually order is preserved, but not guaranteed.
- At-least-once delivery → Messages may arrive more than once (duplicates possible).
- **Use cases**: When order doesn't matter (e.g., logging events, processing tasks, sending notifications).

Example:

A website where users upload photos. You don't care if you process photo
 #2 before photo #1, as long as all are processed eventually.

Extra exam points

- Max message size = 256 KB. If bigger, store in S3 and put a link in SQS.
- Max in-flight messages = 120,000. ("in-flight" = picked up but not deleted yet).
- Data durability = messages are stored across **multiple AZs**.

2. FIFO Queue (First-In-First-Out)

- **Throughput**: Up to 3,000 messages/sec with batching (300 without batching).
- Strict ordering → Message #1 always processed before #2.

- Exactly-once processing → No duplicates.
- Deduplication → Avoids duplicate messages by using a deduplication ID or content-based deduplication.
- P Example:
- Bank transactions. You must process deposit \$100 before withdraw \$50.
- Extra exam points
 - Message Group ID → allows multiple groups to be processed in parallel (scaling). Each group preserves order internally.
 - Supports per-message delay (up to 15 mins).

Key Features (applies to both queues)

1. Visibility Timeout

- When a consumer picks up a message, it becomes invisible for a time (default 30s).
- If consumer deletes it → it's gone
- If consumer crashes or doesn't delete in time → message reappears → another consumer can retry.
- This prevents losing messages but can cause **duplicates** if not managed.
- You can **extend** visibility with ChangeMessageVisibility.

2. Dead Letter Queue (DLQ)

- If a message keeps failing (e.g., corrupted or poison message), after X retries (e.g., 5), it's moved to a **DLQ** for debugging.
- Prevents infinite loops.

3. Long Polling

- Instead of constantly asking SQS "Do you have messages?" (which costs money), long polling lets the consumer wait up to 20 seconds until a message arrives.
- Saves cost + reduces empty responses.

4. Message Retention

- Messages can stay from 1 minute to 14 days (default 4 days).
- · Useful if consumers are down for a while.

Extra exam points

- Delay queues → delay delivery of all messages (0–15 mins).
- Per-message delay → only that message is delayed.
- Encryption → messages can be encrypted with KMS.
- Access control → IAM + SQS access policies (for cross-account).
- **VPC Endpoints** → access SQS privately without internet.

Access Patterns

Polling

Consumers actively request messages. (SQS never "pushes").

Batching

 Can receive/send up to 10 messages in one API call → saves cost + increases throughput.

Auto Scaling

 You can use CloudWatch metrics (queue length, oldest message age) to auto scale consumers (EC2, Lambda, ECS).

Extra exam points

- Lambda integration → Lambda polls SQS automatically, scales with backlog.
 - Standard + FIFO both supported.
 - For FIFO, Lambda respects MessageGroupId → processes one group sequentially.
- Use CloudWatch metrics like:
 - ApproximateNumberOfMessagesVisible (backlog)
 - ApproximateNumberOfMessagesNotVisible (in-flight)

Practical Nuggets & Exam Gotchas

- If question mentions ordering + exactly-once → FIFO Queue.
- If question mentions scale / millions of TPS / order doesn't matter →
 Standard Queue.
- Always design consumers to be idempotent (safe to process same message twice).
- Visibility timeout tuning is critical → too short = duplicates, too long = delays.
- Always use DLQ for poison messages.
- For large payloads >256 KB → S3 + SQS Extended Client.
- Security → use KMS + TLS + VPC endpoints.
- Cost → optimize with long polling + batching.

Quick Analogy Recap (for memory)

- Standard Queue = Huge post office, unlimited letters, but mail sorting may mix up order, and sometimes you get duplicate copies.
- **FIFO Queue** = Small, slower post office that guarantees letters are delivered in order, once and only once.
- Visibility Timeout = Like borrowing a book → if you don't return it in time, the library makes it available for others again.
- **DLQ** = A "problem shelf" in the library for damaged books.
- Long Polling = Instead of asking the librarian every second, you wait at the desk until a book arrives.

Amazon SNS (Simple Notification Service)

Message Delivery

Push Model: Publish once, deliver to multiple subscribers.

- Subscribers: Email, SMS, HTTP/S, SQS, Lambda, mobile push.
- Message Filtering: JSON-based message filtering.

Extra Points for Exam:

- Messages are **stored redundantly** across AZs until delivered.
- Supports delivery retries with exponential backoff (for HTTP/S).
- **DLQ integration** available via SQS subscription.

Explanation:

SNS is like a **megaphone** . When you publish a message, SNS **pushes** it out instantly to all subscribers.

- Subscribers can be humans (email/SMS) or AWS services (Lambda, SQS, HTTP endpoints).
- Filtering means subscribers can choose what they want (e.g., "only receive messages where region=us-east-1).
- Reliability: Messages are stored across multiple AZs and retried if delivery fails.
- If messages still fail (e.g., subscriber down), you can send them to a DLQ (SQS) for later analysis.

Topic Types

- **Standard Topics**: High throughput, best-effort ordering.
- **FIFO Topics**: Strict ordering, exactly-once delivery.

Extra Points for Exam:

- FIFO topics require FIFO queues as subscribers.
- FIFO topics use Message Group ID for ordering.

Explanation:

SNS has two "flavors":

- Standard Topic → Super fast, millions of messages per second, but ordering is not guaranteed, and duplicates are possible.
- FIFO Topic → Slower (up to 3000 msgs/sec with batching), but guarantees strict ordering and exactly-once delivery. Useful for financial transactions

Features

- Mobile Push: iOS, Android, Windows, Baidu notifications.
- Message Attributes: Metadata for message filtering.
- Large Payloads: Store messages in S3, send reference via SNS.

Extra Points for Exam:

- Message durability: Messages not persisted, but stored until delivery or retry exhausted.
- Protocol fan-out: One message can be sent to multiple protocols simultaneously.

Explanation:

- SNS integrates with mobile push notification services (APNS, FCM, Baidu Cloud Push).
- Message Attributes = key-value pairs → let you filter messages at the subscriber level.
- For big messages (over 256KB), SNS can store the full data in S3 and just send the **link**.
- Messages are temporarily stored (durability), but SNS is not a message queue → it doesn't keep messages forever.

Fan-Out Pattern

- Design: SNS topic with multiple SQS queue subscribers.
- Benefits: Decouple publishers from subscribers, reliable delivery.
- **Use Cases**: Parallel processing, multiple service integration.

Extra Points for Exam:

- Exam often tests SNS + SQS fan-out for:
 - Guaranteed delivery to multiple systems.
 - Decoupling + retries.
 - Handling slow/unreliable subscribers with SQS buffering.

Explanation:

The **Fan-out Pattern** = one publisher, many consumers.

- Example: An e-commerce app publishes an OrderPlaced event to SNS.
- SNS pushes that message to:
 - An SQS queue for shipping system.
 - An SQS queue for billing system.
 - A Lambda function for sending email.
- Each subscriber gets its own **copy** of the message.

This is powerful because:

- If one consumer is slow/unreliable → no problem, its SQS queue buffers messages until it's ready.
- Other consumers are not affected.

Exam Tip (SNS vs SQS vs EventBridge)

- **SQS** = Queue → buffering + decoupling (pull model).
- SNS = Pub/Sub → push notifications, instant delivery to many.
- EventBridge = Smart routing → advanced filtering + SaaS integrations.

Common Exam Trick: If you see "fan-out with retries + slow subscribers", the correct answer is SNS + SQS together.

What is Amazon EventBridge?

Imagine AWS services, apps, or even SaaS tools (like Zendesk, Shopify, Datadog) are **shouting events** like:

- "A new file was uploaded to S3"
- "An EC2 instance started"
- "A DynamoDB table was updated"
- "A SaaS app detected a new ticket"

EventBridge is like a **smart event router** that listens to those events and delivers them to the **right target**, based on rules you define.

Key Building Blocks

1. Event Sources

Where events come from:

- AWS Services (S3, EC2, DynamoDB, etc.)
- Your Applications (you can put custom events into EventBridge using APIs)
- SaaS Applications (e.g., Zendesk, PagerDuty, Shopify)

2. Event Bus

- A **pipeline** where events travel.
- AWS gives you a default event bus (for AWS services).
- You can create custom event buses for your apps.
- SaaS apps may get their own event bus.

3. Rules

Rules decide what happens to an event:

- You define patterns → "If event type = S3:ObjectCreated"
- You can use schedule rules → "Run every day at 6PM" (cron/interval).

4. Targets

Where events are sent after matching a rule:

- Lambda (most common, to run code)
- SQS (queue events for later)
- Step Functions (kick off workflows)
- Kinesis (send to streaming pipeline)
- ECS tasks, API Gateway, Event buses, HTTP endpoints

Why Use EventBridge?

- **Decoupling**: Producers (services/apps) don't care who consumes events.
- **Filtering**: Only forward events that match your rules (SNS filtering is more limited).
- Fan-out: One event can go to multiple targets.
- Retry logic + DLQ: If delivery fails, EventBridge retries; if still fails, you can send to a Dead Letter Queue (SQS/SNS).
- Cross-account delivery: Send events between AWS accounts.
- Archive & Replay: Keep past events and replay them (great for testing).

Example

- Say you run an e-commerce site:
 - 1. A customer places an order → DynamoDB emits an event.
- 2. EventBridge rule matches event type = "NewOrder".
- 3. Targets:
 - Lambda → Send confirmation email.
 - Step Functions → Run payment workflow.
 - SQS → Send order details to shipping system.

All this happens **automatically**, in real time.

Comparison (Exam Tip)

- **SQS** = Queue → Buffering + decoupling (pull model).
- **SNS** = Pub/Sub → Broadcast notifications (push model).
- EventBridge = Smart event router → Filtering + routing + SaaS/AWS integration.
- **Exam Trick**: If you see "advanced filtering, SaaS integration, event routing"
- → Answer = **EventBridge**.
- \leftarrow If you see "simple notifications" \rightarrow **SNS**.
- \leftarrow If you see "queueing for consumers" \rightarrow **SQS**.

Advanced Features

- Custom Event Buses: Isolate events by application or environment.
- Schema Registry: Discover and manage event schemas.
- Archive and Replay: Store events and replay for testing.

Extra Points for Exam:

- Cross-account event routing supported.
- EventBridge is serverless, highly available, and scales automatically.
- Exam often asks difference between SNS vs EventBridge:
 - SNS = pub/sub, simple notification.
 - EventBridge = event bus + complex routing + filtering.

∳ Exam Tip:

- **SQS** = decoupling, buffering.
- **SNS** = push to multiple subscribers.
- **EventBridge** = routing/filtering events with complex rules.
- Often tested in fan-out, ordering, retries, and DLQ scenarios.

6 AWS Step Functions

Think of **Step Functions** as a **workflow engine**.

It lets you design workflows that connect multiple AWS services together.

- You write the workflow as a state machine (like a flowchart).
- Each step (called a state) does some work, decides what to do next, or waits.
- AWS Step Functions make sure each step runs in order, handles errors, and retries if needed.

It's like building a **pipeline of tasks** where AWS makes sure things don't break halfway.

State Machine Types

1. Standard Workflows

- Think of them as long-running, reliable workflows.
- They keep **detailed history** of every execution.
- Exactly-once execution → AWS guarantees each step runs only once in order.
- Duration: Up to 1 year.
- Use Case: Order processing, loan approval, multi-step transactions.

2. Express Workflows

- Think of them as fast and lightweight workflows.
- Very cheap and scales to millions of executions per second.
- Less detailed history (you get logs, not full execution history).
- Duration: Up to 5 minutes.
- Use Case: IoT data ingestion, streaming processing, real-time ETL.

State Types (Building Blocks)

1. Task

- · Actually does the work.
- Can call AWS Lambda, ECS, Batch, DynamoDB, SNS, SQS, etc.
- Example: "Run a Lambda function to process an order."

2. Choice

- Like an if/else statement.
- Example: "If amount > 10,000 → require approval, else → auto approve."

3. Parallel

- Runs multiple branches at the same time.
- Example: "Do fraud check and credit check in parallel."

4. Wait

Adds a delay before the next step.

Example: "Wait 1 hour before retrying."

5. **Pass**

- Doesn't do work, just **transforms data** or passes it forward.
- Example: "Add a flag to data and move to next state."

6. Fail / Succeed

• Ends the workflow with success or failure.

Error Handling

Workflows in real life fail sometimes. Step Functions handle errors for you:

- **Retry**: Automatically retries on failure (with backoff like 1s, 2s, 4s...).
- Catch: If retries fail, send the workflow to an alternate path (like a fallback).
- Dead Letter Queues (DLQ): Persist errors to SQS or SNS if things fail completely.

Use Cases

1. Data Processing (ETL)

- Example: Upload file → Validate → Transform → Save to Database.
- Step Functions handle each step, retries if one fails.

2. Application Orchestration (Microservices)

- Example: An order system with payment, inventory, shipping.
- Each service is independent, Step Functions connect them in order.

3. Human Approval

- Example: Loan application → Wait for human approval → Continue workflow.
- Workflow pauses until approval is received.

Exam Tips

Standard vs Express → Long-running vs high-volume short tasks.

- State Types → Know what each does.
- Error handling → Retry, catch, DLQ are frequently tested.
- Step Functions often appear in "serverless orchestration" questions.
- Common exam pattern: "How to coordinate multiple Lambda functions into a sequence with retries and monitoring?" → Answer: AWS Step Functions.

Amazon API Gateway

What is API Gateway?

- Think of API Gateway as a front door for your applications.
- It sits between your clients (mobile apps, browsers, IoT devices) and your backend (Lambda, EC2, DynamoDB, etc.).
- Its job is to:
 - Accept requests
 - Process them (security, rate limits, caching, transformation)
 - Send them to the right backend service
 - Return responses to the client

✓ In short: API Gateway = Managed API management service by AWS.

API Types

API Gateway supports different types of APIs, depending on use case:

- 1. REST API (Legacy, Full-featured)
 - Supports all features: request/response transformations, usage plans, caching, advanced security.
 - Good for complex, feature-rich APIs.
 - More expensive than HTTP APIs.
 - Example: An e-commerce platform with multiple endpoints like /products , /cart , /orders .

2. HTTP API (New, Lightweight, Cheaper)

- Faster and cheaper (up to 71% less cost) compared to REST API.
- Limited features but enough for most workloads.
- Great for microservices, serverless apps, and simple APIs.
- Example: A serverless blog that fetches data from DynamoDB via Lambda.

3. WebSocket API (Real-time Communication)

- Supports two-way communication (client ↔ server).
- Great for chat apps, real-time notifications, gaming.
- Example: A live chat app where messages update instantly without refreshing.

Integration Types

How does API Gateway talk to the backend?

1. Lambda Proxy Integration

- Request is passed directly to AWS Lambda.
- The Lambda function handles all the logic and returns the response.
- Example: /user endpoint calls a Lambda that fetches user details from DynamoDB.

2. AWS Service Integration

- API Gateway connects directly to AWS services without Lambda.
- Example: /items endpoint connects directly to DynamoDB PutItem action.

3. HTTP Integration

- API Gateway forwards requests to an external HTTP endpoint (outside AWS).
- **W** Example: Connect to a **third-party weather API**.

4. Mock Integration

- API Gateway returns a **static response** (no backend).
- Useful for testing or when backend isn't ready.

• V Example: /status always returns { "status": "ok" }.

Features

API Gateway provides many built-in features to make APIs **secure**, **fast**, **and reliable**:

1. Authentication / Authorization

- IAM → Restrict API usage to AWS users/roles.
- Cognito → Allow users to log in with username/password or social login (Google, Facebook).
- Lambda Authorizer → Custom auth logic using a Lambda function.
- Z Example: Only logged-in users can call /checkout.

2. Rate Limiting

- Prevents abuse (too many requests).
- Set throttling and usage plans.
- Example: Free users can make 100 requests/day, premium users 1000 requests/day.

3. Caching

- Store responses at API Gateway for a set time.
- Reduces load on backend and improves performance.
- Example: /products list is cached for 60 seconds.

4. Request/Response Transformation

- Modify the request before sending to backend.
- Modify the response before sending to client.
- Z Example: Convert XML from backend into JSON for client.

Deployment

API Gateway supports multiple environments and rollout strategies:

1. Stages

• Like different versions of your API.

- Example: dev , test , prod .
- Each stage can have its own settings (logging, caching, throttling).

2. Canary Deployment

- Roll out a new version to a small % of users first.
- Example: Release new /checkout logic to 10% of traffic, keep 90% on old version.
- Helps catch bugs before full rollout.

3. Custom Domain Names

- Use your own domain instead of the default AWS one.
- X Example: Instead of

 $\label{eq:https://abc123.execute-api.us-east-1.amazonaws.com/prod} \\ \mbox{you can use} \rightarrow \mbox{https://api.myshop.com} \; . \\$

Extra Advanced Points

- Private APIs: Can only be accessed from inside your VPC.
- API Gateway + CloudWatch: Built-in monitoring, logs, and metrics.
- Security: Supports WAF (Web Application Firewall) to block malicious traffic.
- Cost Model: Pay per request (cheaper for HTTP API).
- **Integration with Step Functions**: APIs can trigger workflows for long-running processes.

Quick Analogy:

Imagine you own a restaurant (your backend services).

- Customers (clients) can't go to the kitchen directly → they order via a waiter (API Gateway).
- The waiter checks if they're allowed (Authentication), ensures no one orders too much (Rate Limiting), remembers popular orders (Caching), and delivers the food (Response).
- The waiter can also adjust the order format (Transformation).

Analytics and Machine Learning {#analytics}

Amazon Kinesis

Think of Kinesis as AWS's real-time data streaming platform.

Instead of waiting for data to be stored in a database or file, Kinesis lets you collect, process, and analyze data the moment it's generated.

Example: Imagine Netflix \rightarrow every time a user presses play/pause, that event can be sent in real time to Kinesis for analytics (recommendations, dashboards, alerts).

1. Kinesis Data Streams (KDS)

- Purpose: **Real-time data ingestion** (you want to process data as it arrives).
- Data is stored across shards (like partitions).
- Retention: 1–365 days (default 24 hours).
- Consumers: AWS Lambda, KCL apps (Kinesis Client Library), Kinesis Data Analytics.

Key Concepts:

- Shard = unit of capacity.
 - 1 shard = 1 MB/sec write, 2 MB/sec read, 1000 records/sec.
 - More shards = more throughput (you can reshard).
- Parallelism: Each shard can be consumed by one consumer in parallel.
- Data Retention: Lets you reprocess old data (e.g., replay logs).

Example Use Case:

 Stock trading platform → process trades in real time, detect fraud, update dashboards instantly.

2. Kinesis Data Firehose (KDF)

- Purpose: Load data from streams into storage/analytics systems automatically.
- No need to write custom code for processing.

- **Destinations**: S3, Redshift, OpenSearch, Splunk.
- Transformations: Can use Lambda to clean/format data before delivery.
- **Buffering**: Groups data by size (1 MB–128 MB) or time (60–900 sec) before sending.

Example Use Case:

- Send website clickstream logs into S3 → use Athena to query logs later.
- Or → send logs into Redshift for BI dashboards.

🔷 3. Kinesis Data Analytics (KDA)

- Purpose: Analyze streaming data in real time.
- **SQL Mode**: You can write SQL queries directly on data streams.
- Apache Flink Mode: For developers → build custom, complex stream processing apps in Java/Scala.

Example Use Cases:

- Real-time dashboards → monitor how many users are watching Netflix right now.
- Real-time fraud detection → alert when suspicious transactions occur.
- IoT anomaly detection → flag devices sending unusual data.

4. Kinesis Video Streams

- Purpose: Ingest and process live video streams securely.
- Data is available for ML, analytics, or replay.
- Integrates with Amazon Rekognition Video for object/face detection.

Example Use Cases:

- Smart home camera → send video to AWS for live monitoring.
- Retail store → analyze foot traffic in real time.

Amazon EMR (Elastic MapReduce) (Beginner → Advanced)

Amazon EMR = AWS's big data processing platform.

It's basically **Hadoop/Spark clusters in the cloud** → but fully managed, so you don't worry about installing, scaling, or maintaining them.

Example: Imagine you have **100 TB of log data**. Analyzing this on one computer is impossible. With EMR, AWS spins up a cluster of hundreds of nodes to process it in parallel.

◆ Supported Frameworks (Big Data Tools you can run on EMR)

- 1. **Apache Spark** → General-purpose, in-memory computing (fast, versatile).
 - ▼ Example: ETL pipelines, ML jobs.
- 2. **Apache Hadoop** → Batch processing framework (MapReduce model).
 - V Example: Processing petabytes of logs overnight.
- 3. **Apache Hive** → SQL-like queries on top of big data.
 - X Example: Analysts query logs with SQL instead of writing code.
- 4. **Apache HBase** → NoSQL database for real-time read/write access.
 - **V** Example: Store clickstream data for millisecond lookups.
- 5. **Presto** → Fast SQL query engine across multiple data sources.
 - Z Example: Query S3, RDS, and DynamoDB at the same time with SQL.

Cluster Types

- Long-Running Clusters
 - Always on → used for continuous jobs.
 - Z Example: A financial company runs Spark ML models every hour.
- Transient Clusters
 - Temporary → spin up, run job, shut down (saves cost).

∘ ✓ Example: Process 1 TB of logs once per day, then terminate.

Notebooks

- Interactive Jupyter-like environment for development.

Storage in EMR

- HDFS (Hadoop Distributed File System)
 - Stores data across EMR cluster instances.
 - Fast, but tied to cluster lifespan.

EMRFS

- Lets EMR directly read/write to Amazon S3 as storage.
- This is the most common option (S3 = durable, scalable).
- ▼ Example: Store all logs in S3, process them with Spark via EMRFS.

Local Storage

Temporary instance storage (used for caching/processing).

Quick Exam Nuggets

Amazon Kinesis

- Data Streams = **real-time ingestion** (requires you to build consumers).
- Firehose = load into destinations automatically (no consumer code).
- Analytics = real-time SQL/stream processing.
- Video Streams = video ingestion & ML analytics.
- Shards control capacity (scaling via resharding).
- · Lambda often used as consumer.

Amazon EMR

- Managed Hadoop/Spark → process big data.
- EMRFS lets you process S3 data directly.

- Cluster types: Long-running, Transient, Notebooks.
- Frameworks: Spark, Hadoop, Hive, HBase, Presto.
- Use transient clusters for cost savings.

Analogy for remembering:

- **Kinesis = streaming river** (data flows in real time, consumers pick it up).
- EMR = factory (takes huge piles of raw material → processes into useful output, but usually in batches).

Amazon Redshift (Data Warehouse)

Imagine you have **huge amounts of business data** (sales, customers, transactions, logs). You don't just want to store it—you want to **analyze it fast** (e.g., "What were my top 10 selling products this month?").

That's where **Amazon Redshift** comes in \rightarrow it's AWS's **data warehouse** (designed for analytics, not transactions).

Architecture

- Leader Node
 - Acts as the brain.
 - Accepts queries (usually in SQL), creates a plan, and distributes work to compute nodes.
- Compute Nodes
 - Store the data and run the queries.
 - Multiple compute nodes = queries run in parallel (faster).

Node Types

- DC2 (Dense Compute) → SSD storage, high performance, but limited storage.
- DS2 (Dense Storage) → HDD, large storage, but slower.
- RA3 (Managed Storage) → AWS manages storage separately, you scale compute independently (most popular now).

- ✓ Analogy: Think of Redshift as a restaurant →
 - Leader node = head chef (decides who does what).
 - Compute nodes = **kitchen staff** (do the actual cooking).

Performance Features

- Columnar Storage
 - Instead of storing data row by row (like RDS), Redshift stores by column.
 - This makes analytical queries (SUM, AVG, GROUP BY) much faster.
- Compression
 - Saves space by compressing similar data (e.g., all "yes/no" values).
- Massively Parallel Processing (MPP)
 - Splits big queries across multiple nodes, runs them simultaneously.
- Result Caching
 - Frequently run queries are cached → faster response.

Redshift Spectrum

- Lets you query data directly in S3 without loading it into Redshift.
- Supports data formats: Parquet, ORC (best performance), JSON, CSV, Avro.
- Elastic → scales compute separately from storage.
- ✓ Use Case: You have **some data in Redshift**, but **archived logs in S3**. Instead of loading all into Redshift, you can run queries across both seamlessly.

Backup & Recovery

- Automated Snapshots → Daily, incremental backups.
- Manual Snapshots → Trigger anytime you want.
- Cross-Region Snapshots → For disaster recovery (DR).

Amazon Athena (Serverless Query Engine)

Athena lets you query S3 data directly using SQL, without setting up servers.

Serverless SQL

- No infrastructure → You don't manage servers or clusters.
- Pay-per-Query → Charged per TB scanned (\$5 per TB).
- Uses ANSI SQL (standard SQL).
- $extbf{ extbf{ iny Cloud Trail logs in S3}}$ with Athena, you can write:

SELECT user, COUNT(*) FROM cloudtrail_logs WHERE eventName='StartIn stances';

And get results immediately.

Performance Optimization

Because Athena charges per data scanned, you want to reduce scanned data:

- Columnar Formats (Parquet, ORC) → Much smaller, faster queries.
- Partitioning → Store data by date, region, etc., so queries scan only relevant files.
- Compression → Reduces size and costs.
- Result Caching → Repeated queries are cached (faster, cheaper).

Integration

- QuickSight → Build dashboards on top of Athena queries.
- Glue Data Catalog → Store schema (Athena uses it).
- CloudTrail → Directly analyze AWS logs.

AWS Glue (ETL Service)

Glue is AWS's data integration service.

It helps move and transform data between sources (Extract \rightarrow Transform \rightarrow Load).

 \checkmark Example: You want to take **raw logs from S3**, clean them, and load them into **Redshift** for reporting \rightarrow Glue does this automatically.

ETL Service

- **Serverless** → No servers to manage.
- Runs on Apache Spark (distributed engine → handles big data).
- Supports **Python and Scala** code for transformations.

Data Catalog

- Think of it as the metadata library for all your data.
- Crawlers → Automatically scan data in S3/RDS/Redshift to detect schema (columns, types).
- Integration → Works seamlessly with Athena, EMR, Redshift.

Job Types

- ETL Jobs → Take raw data, clean/transform, load to destination (e.g., S3 → Redshift).
- Streaming Jobs → Real-time ETL (e.g., Kinesis → S3).
- Python Shell → For lightweight tasks (not big Spark jobs).

Quick Exam Nuggets

Redshift

- Data warehouse for **analytics (OLAP)**, not transactions.
- Columnar storage = best for queries like SUM, AVG.

- **Redshift Spectrum** = query S3 directly.
- Snapshots (auto/manual, cross-region for DR).
- Node types: DC2 (fast SSD), DS2 (cheap HDD), RA3 (scale compute/storage separately).

Athena

- Serverless SQL over S3.
- Pay per TB scanned → optimize with Parquet/ORC, partitions, compression.
- Works with Glue Catalog for schema.
- Integrates with QuickSight for dashboards.

Glue

- ETL service (serverless, Spark-based).
- **Data Catalog** = schema repository.
- Crawlers = auto schema detection.
- Job Types = ETL, Streaming, Python Shell.

Analogy for remembering:

- Redshift = Library (organized, structured, high-performance queries).
- Athena = Librarian (asks questions directly to S3, no database needed).
- Glue = Organizer (takes messy books, labels them, puts them in right shelves).

Amazon QuickSight (Explained for Beginners)

Amazon QuickSight is a Business Intelligence (BI) tool from AWS.

Think of it like **Google Data Studio / Power BI / Tableau**, but **serverless** and built into AWS.

Its main purpose: turn raw data into dashboards, charts, and insights.

Business Intelligence (BI) with QuickSight

• Serverless (Fully managed BI service)

- You don't need to install software, manage servers, or worry about scaling.
- AWS takes care of the backend (storage, compute, scaling).
- Example: If your company's sales data grows from thousands to millions of records, QuickSight will scale automatically.

SPICE Engine (In-memory calculation engine)

- SPICE = Super-fast, Parallel, In-memory Calculation Engine.
- Instead of querying data directly from a database every time,
 QuickSight can cache it in memory for super-fast results.
- Imagine you have a sales dashboard with 1 million records instead of waiting minutes for SQL queries, QuickSight SPICE shows results in seconds.

Machine Learning Features

- AutoNarrative: Automatically writes plain English summaries of data (e.g., "Sales increased by 15% in Q1 compared to Q4").
- Anomaly Detection: Flags unusual patterns (e.g., "Why did sales in California suddenly drop by 30%?").
- **Forecasting:** Uses ML models to **predict future trends** (e.g., "Next month's revenue is expected to be \$50,000").
- Benefit: You don't need to be a data scientist; QuickSight gives Alpowered insights automatically.

Data Sources

QuickSight can connect to data stored in many places:

- AWS Services (native integration)
 - RDS (Relational Database Service MySQL, PostgreSQL, etc.)
 - Redshift (Data warehouse for big data analytics)
 - S3 (Store CSV, JSON, Parquet files, etc.)

- Athena (Query S3 data directly using SQL without servers)
- SaaS Applications (cloud apps your company might already use)
 - Salesforce (customer data)
 - ServiceNow (IT service management data)
 - Twitter (social media analytics)

On-Premises Databases

- SQL Server, Oracle, MySQL (via Direct Connect or VPN to AWS).
- Example: If your company still keeps financial data in a local Oracle DB,
 QuickSight can still pull reports.
- Key takeaway: QuickSight can pull data from almost anywhere cloud, SaaS, or on-premises.

Sharing and Collaboration

This is how teams actually use QuickSight:

Dashboards

- Visual, interactive charts and graphs.
- Users can apply filters (e.g., show sales only in India, or only for 2023).
- Example: A CEO sees global sales trends; a regional manager sees sales in their territory.

Stories (Guided analytics narratives)

- Like a slideshow built from dashboards.
- Walks stakeholders through data step-by-step ("First, revenue trends...
 then customer churn... then forecast").
- Useful for board meetings or presentations.

Embedded Analytics

- QuickSight dashboards can be embedded inside your company's apps or portals.
- Example: If you're building an e-commerce site, you can show vendors their sales performance inside your platform using QuickSight.

Benefit: You don't have to build a separate analytics tool.

Why QuickSight Matters (for SAA exam)

- Fully managed Bl service (serverless, auto-scales).
- Uses **SPICE engine** for speed.
- Integrates deeply with AWS data services + SaaS + on-prem.
- Provides ML-powered insights (without data science knowledge).
- Enables collaboration via dashboards, stories, and embedding.
- In the AWS SAA exam, QuickSight questions usually test:
- When to use QuickSight vs Athena/Redshift vs EMR.
- Knowing that QuickSight is for visualization/BI, not for storing data.
- Knowing SPICE is in-memory and improves performance.
- Recognizing ML features like anomaly detection and forecasting.
- So in simple words:

QuickSight = AWS's BI tool → Connect data → Analyze → Visualize → Share insights.

AWS Machine Learning & AI Services (SAA Exam Notes with Explanations)

Amazon SageMaker

(AWS's main ML platform)

Full ML Lifecycle →

SageMaker helps you with the **entire process** of machine learning:

- 1. **Build** → Prepare and clean your data.
- 2. **Train** → Use algorithms to train models.

- Deploy → Put the trained model into production so applications can use it.
- Example: Suppose you want to predict house prices. With SageMaker, you can upload your dataset, train a model, and then deploy it so your app can show predictions.

Keyword: *ML lifecycle, build-train-deploy*

Jupyter Notebooks →

SageMaker comes with **Jupyter notebooks** (interactive coding environment).

You can write Python code, test models, and visualize results in one place.

Think of it as Google Colab but fully managed by AWS.

Keyword: Managed Jupyter IDE

• Built-in Algorithms →

AWS provides **pre-built ML algorithms** (like XGBoost, Linear Learner, Random Cut Forest, etc.).

You don't always need to write ML code from scratch.

Example: Use a pre-built sentiment analysis model to detect if a tweet is positive/negative.

Keyword: *Pre-built ML algorithms*

Model Endpoints →

After training, you can deploy models as **endpoints (APIs)**.

Applications can send requests to these endpoints and get predictions in **real-time** or **batch mode**.

Example: A web app can call SageMaker's endpoint to recommend a product instantly.

Keyword: Deploy ML models as API

Al Services

(AWS Pre-trained ML models for specific use cases → **no ML expertise needed**)

1. Rekognition

- Image & video analysis.
- Can detect **faces**, **objects**, **scenes**, **activities**, and even unsafe content.
- Used in security cameras, photo tagging, and identity verification.
- Example: Facebook photo auto-tagging is similar to Rekognition.

Keywords: Face detection, Object recognition, Video analysis, Content moderation

2. Textract

- Extracts text, handwriting, tables, forms from scanned documents/PDFs.
- More advanced than OCR → because it understands structure (forms, invoices, receipts).
- Example: Scan a medical form and automatically digitize patient info.

Keywords: OCR+, Document processing, Structured data extraction

3. Comprehend

- NLP (Natural Language Processing) service.
- Can find sentiment (positive/negative/neutral), detect entities (names, places, dates), and classify text.
- **Example:** Analyze customer reviews to see overall sentiment.

Keywords: NLP, Sentiment analysis, Entity recognition, Text classification

4. Polly

- Converts text → speech (human-like voices).
- Supports multiple languages and accents.
- Example: News apps reading articles aloud.

Keywords: TTS (Text-to-Speech), Human-like speech

5. Transcribe

Converts speech → text.

- Can add timestamps, detect multiple speakers, and recognize domainspecific terms.
- Example: Convert customer service call recordings into searchable text.

Keywords: Speech recognition, Speech-to-text, Transcription

6. Translate

- Language translation at scale.
- Supports many languages, used for real-time chat translation or global websites.

Keywords: Language translation, Multilingual, Real-time translation

7. Lex

- Build chatbots and voice assistants.
- The same tech behind Alexa.
- Integrates with Polly (voice) and Lambda (logic).
- Example: A bank chatbot that answers balance queries.

Keywords: Conversational AI, Chatbots, Alexa technology

8. Personalize

- Recommendation engine.
- Trains on user behavior + product data to provide personalized suggestions.
- No ML expertise required.
- Example: Netflix-style movie recommendations or Amazon product suggestions.

Keywords: Recommendation engine, Personalization, User behavior analysis

© Quick Exam Tip (Summary Table)

Service	Use Case	Keyword
SageMaker	Build/train/deploy ML models	ML lifecycle, Jupyter, endpoints
Rekognition	Image & video analysis	Face detection, moderation
Textract	Extract text from documents	OCR+, forms, invoices
Comprehend	NLP (text analysis)	Sentiment, entities
Polly	Text-to-speech	Human-like voice
Transcribe	Speech-to-text	Call center transcripts
Translate	Language translation	Multilingual, real-time
Lex	Chatbots, voice assistants	Alexa, conversational Al
Personalize	Personalized recommendations	Recommendation engine

Disaster Recovery and Migration {#dr-migration}

Oisaster Recovery (DR) Strategies — The Ultimate Guide

Think of Disaster Recovery like backup plans for your life:

If your house (primary environment) is gone, what do you have ready at another place?



Analogy:

Like keeping your important documents in a locker far away.

When your house burns down, you'll travel there, open the locker, and rebuild your house from scratch. Takes time, but it's the cheapest option.

Property Details:

- RTO (Recovery Time Objective): Hours to Days ∑
 (Because you must restore from backups, install apps, reconfigure.)
- Cost: Lowest

- Implementation: Automated backups → stored in Amazon S3. Restore when disaster strikes.
- Use Cases: Dev/test, non-critical apps (blog, hobby projects).

When to Choose:

- Choose this if you want minimum cost
- App can tolerate downtime of hours-days.
- Example: College project DB, blog, reporting DB.

🔼 Pilot Light 🥠 🢡



Analogy:

Think of a gas stove's pilot flame \rightarrow always a tiny flame burning.

When you want to cook (recover), you just turn the knob, and it's instantly ready.

Property Details:

- RTO: Tens of minutes to hours
- **RPO**: Minutes to hours
- Cost: Low → Medium
- Implementation:

Keep core infrastructure ON (DB, minimal servers) in the DR region.

Rest of the system is off, but can be quickly launched from templates/AMIs.

• Use Cases: Business-critical apps where downtime is not OK, but still tolerable (HR system, ERP).

When to Choose:

- · Choose this if:
 - Need **faster recovery** than Backup/Restore.
 - Willing to pay a bit more.
 - Don't need everything instantly running 24/7.

Warm Standby 🤊

Analogy:

Like having a second house with minimal furniture.

If your main house burns, you can **move in quickly**, but maybe need to buy extra furniture to live comfortably.

Property Details:

• RTO: Minutes to Hours

• **RPO**: Minutes

• Cost: Medium to High 1991

• Implementation:

Keep a **scaled-down version** of production (fewer servers, smaller DB instance).

When needed \rightarrow scale up quickly.

Use Cases: Mission-critical apps (online banking, hospital systems).

★ When to Choose:

- Choose this if:
 - You want fast failover (minutes).
 - ✓ You can afford extra cost.
 - ✓ Downtime must be minimal but not zero.

🔼 Multi-Site Active/Active 🌍 🦩

Analogy:

Like owning two identical fully furnished houses.

Even if one burns, you're already living in the other. **No disruption.**

Property Details:

• RTO: Seconds to Minutes

• RPO: Near Zero

Cost: Highest ** ** ** ** **

• Implementation:

Full production environments running in multiple regions.

Uses DNS routing (Route53), load balancers, replication.

• Use Cases: Apps that cannot tolerate downtime (Amazon, Netflix, stock trading platforms).

★ When to Choose:

- Choose this if:
 - **▼** Business is **mission-critical with zero downtime tolerance**.
 - You're willing to pay premium cost.
 - Customer-facing apps with global reach.

Memory Tricks

Like increasing **firepower**:

- Pilot Light → Small flame (low-medium cost, faster)
- Warm Standby → Warm house \(\strice{\cost} \) (medium-high cost, faster)

Decision Matrix

Strategy	Cost 🐇	RTO 📆	RPO 🚡	Best For
Backup & Restore	Lowest	Hours → Days	Hours	Non-critical, cheapest
Pilot Light	Low-Med	Minutes → Hrs	Minutes → Hrs	Business apps, moderate RTO
Warm Standby	Med-Hi	Minutes → Hr	Minutes	Mission-critical apps
Multi-Site Active	Highest	Seconds → Min	Near Zero	Zero downtime apps

Final Shortcut (When to choose what)

- Budget tight → Backup/Restore
- Some critical → Pilot Light
- Mission critical → Warm Standby
- Zero downtime → Multi-Site

AWS Migration Services (SAA Exam Notes – Explained)

Think of migration as moving your stuff (databases, apps, servers) from your home (on-premises or other cloud) into your new house (AWS cloud).

AWS gives different **moving trucks/tools** depending on what exactly you're moving.

AWS Database Migration Service (DMS)

Helps you move databases into AWS.

1. Homogeneous Migration

- Meaning: Same database engine at both ends.
- Example: Oracle → Oracle OR MySQL → MySQL.
- Only the data is transferred, no need to change database features.
- Easy and fast because it's like copying files between two folders.

2. Heterogeneous Migration

- Meaning: Different database engines.
- Example: Oracle → PostgreSQL OR SQL Server → Aurora.
- Harder, because the structure and SQL commands might be different.
- AWS provides a Schema Conversion Tool (SCT) that automatically converts database structure (tables, views, stored procedures).

So:

- Homogeneous = Simple copy.
- Heterogeneous = Needs translation (SCT does that).

3. Continuous Replication

- DMS doesn't just do a one-time move.
- It can keep syncing new changes (inserts, updates) from source → target until you're ready to switch over.
- This reduces downtime.

DMS – Supported Sources (Where data can come from):

- On-premises databases (Oracle, SQL Server, MySQL, PostgreSQL).
- AWS RDS databases (all engines).
- Amazon S3 (if you stored database dumps there).
- NoSQL like MongoDB or Amazon DocumentDB.

🔽 DMS – Supported Targets (Where data can go to):

- Amazon RDS or Aurora (for relational workloads).
- Amazon Redshift (for analytics/data warehouse).
- Amazon DynamoDB (for NoSQL apps).
- Amazon S3 (for backup/data lake).
- Amazon OpenSearch Service (for search and analytics).
- Basically, you can move from anywhere → to almost anywhere in AWS.

AWS Application Migration Service (MGN)

Helps you move entire applications & servers (not just databases).

It's called the lift-and-shift tool.

Think: You want to move your **whole company server** (Windows/Linux apps, configs, storage) from your data center → AWS.

1. Lift-and-Shift Migration

- Move the app as it is without rewriting.
- No code changes, just move it to AWS.

2. Continuous Replication

- The service keeps copying your VM/server changes from on-premises → AWS.
- Ensures the AWS copy is always up to date.

3. Non-Disruptive Testing

- You can launch test servers in AWS to see if everything works fine.
- Testing doesn't affect your real, live on-prem servers.

4. Automated Conversion

- It automatically converts your physical or virtual server into an EC2 instance in AWS.
- Handles networking, storage, OS, etc. for you.

AWS Server Migration Service (SMS) –Deprecated

- This was the old tool for migrating VMs (on-prem → AWS).
- Now replaced by MGN (Application Migration Service) because MGN is more advanced (real-time replication, testing, automation).

for the exam:

- If you see SMS, remember it's deprecated.
- The recommended service is MGN.

Quick Analogy (for easy memory)

- MGN (Application Migration Service) = Moving your whole office (servers + apps) into AWS with no redesign.

Exam Tips (SAA-C03)

- If the question says "migrate databases with minimal downtime" → Answer: AWS DMS (with continuous replication).
- If the question says "lift-and-shift whole apps/servers" → Answer: AWS
 MGN.
- If you see SMS, pick MGN (because SMS is deprecated).
- Remember: SCT is only for heterogeneous migrations.

AWS Data Transfer Services (SAA Exam Notes)

When working with AWS, you often need to **move data** between **on-premises systems** and the **AWS Cloud** (or between AWS services).

AWS offers multiple tools for this — each suited for **different sizes**, **speeds**, and **scenarios**.

AWS DataSync

- What it is: A fully managed service to move data online (over the internet or Direct Connect) between on-premises storage and AWS.
- Sources:
 - On-premises NFS (Linux file servers)
 - SMB (Windows file servers)
 - AWS S3, EFS, FSx
- Destinations:
 - AWS **S3**, **EFS**, **FSx**
- Features:
 - One-time migration or scheduled transfers (good for ongoing syncs)
 - Automatic data verification (ensures no corruption)
 - Bandwidth throttling (control network usage)

Secure transfer (TLS encryption)

When to use DataSync:

- You want fast + automated online transfers.
- You have a steady internet/Direct Connect link.
- You need incremental syncs (not just one-time copies).
- Example: Migrate 50 TB of file server data into S3 with automatic scheduling.

AWS Snow Family

Used when internet bandwidth is too small or too expensive for moving large data volumes. These are physical devices AWS ships to you \rightarrow you load data \rightarrow ship back \rightarrow AWS uploads into S3.

1. AWS Snowcone

- Smallest device in the Snow Family.
- Capacity:
 - 8 TB HDD or
 - 14 TB **SSD**
- Use Cases:
 - **Edge computing** (run apps close to where data is generated).
 - Collect data in remote/harsh environments (factories, ships, oil rigs, military bases).
- Connectivity: Wi-Fi, wired, cellular (portable device).

When to use Snowcone:

- Small-scale transfers (a few TB).
- Field or edge environments with low power and limited connectivity.
- Example: Scientists in the Arctic collect 5 TB of sensor data and ship it later.

2. AWS Snowball Edge

Medium-scale device for petabyte migrations.

- Comes in two types:
 - Storage Optimized: 80 TB HDD
 - Compute Optimized: 42 TB storage + built-in compute (EC2 + Lambda at the edge)

• Use Cases:

- Large data migrations (10s 100s of TB).
- Edge computing with local processing (e.g., ML inference on video data).
- Temporary storage + batch transfer to AWS.

When to use Snowball Edge:

- You need to transfer tens to hundreds of TBs but don't have enough bandwidth.
- You also want compute power at the edge (pre-processing before sending data).
- Example: Oil exploration company processes seismic data locally, then sends 200 TB to AWS.

3. AWS Snowmobile

- · Huge truck-sized device.
- Capacity: 100 PB (per truck).
- Use Cases:
 - Exabyte-scale data migrations (entire data centers).
- Security: Encrypted, GPS tracking, 24/7 monitoring, escort vehicles.

When to use Snowmobile:

- You need to move massive amounts (PB-EB) of data.
- Moving over the internet would take years.
- Example: A bank migrates its entire data center (70 PB) to AWS in weeks using Snowmobile.

When to Use What (Quick Guide)

Service	Best For	Data Size	Mode	Example
DataSync	Fast online transfer	GB – 100s of TB	Online (Internet/DC)	Migrate 50 TB from on-prem NFS to S3
Snowcone	Small edge data collection	< 14 TB	Offline (ship device)	Collect 5 TB research data in a desert
Snowball Edge	Large migration + edge compute	10s – 100s of TB	Offline (ship device)	Move 200 TB + preprocess video files locally
Snowmobile	Data center migration	100 PB+	Offline (truck)	Move entire data center (70 PB)

W Exam Tip:

- If you see petabytes/exabytes + no bandwidth → Snow Family (Edge/Snowmobile).
- If you see ongoing sync / automation → DataSync.
- If you see tiny portable device / edge computing → Snowcone.
- If you see "truck" or "exabyte" → Snowmobile.

Hybrid Connectivity & Business Continuity (AWS SAA Notes)

Hybrid connectivity = connecting **on-premises environments** with **AWS Cloud**.

Business continuity = ensuring your applications **survive disasters** with minimal downtime/data loss.

AWS Storage Gateway

A **hybrid storage service** that connects on-premises apps with AWS cloud storage.

1. File Gateway

Protocol: Supports NFS (Linux) and SMB (Windows).

• What it does: Provides file shares where files are stored in Amazon S3.

Features:

- Local cache for frequently accessed data (low latency).
- Files appear as normal NFS/SMB shares, but they're actually stored in S3.
- Use Case: Replace on-prem file servers with S3-backed storage.

When to use: You want apps to keep using NFS/SMB without rewriting them, but actually save data in S3.

2. Volume Gateway

Provides block storage volumes to on-premises apps, with cloud integration.

Stored Volumes:

- Primary data stored locally.
- Asynchronous backup to S3.
- Best if you want local low-latency performance + cloud backup.

Cached Volumes:

- Primary data stored in S3.
- Frequently accessed data cached locally.
- Best if you want to reduce on-premises storage costs.

When to use: Need block storage for apps but want to extend into AWS for backup/DR.

3. Tape Gateway

- What it does: Emulates a Virtual Tape Library (VTL).
- Integration: Works with existing backup applications (NetBackup, Veeam, etc.).

Storage:

- Virtual tapes stored in S3.
- Archived to Glacier / Glacier Deep Archive for long-term retention.

- Use Case: Replace expensive physical tape infrastructure with cloudbacked tapes.
- When to use: You want to retire tape libraries but keep using the same backup software.

AWS Direct Connect (DX)

A dedicated, private network connection between on-premises and AWS.

- Bandwidth: 1 Gbps → 100 Gbps.
- Benefits:
 - Lower latency than internet.
 - More **predictable** performance.
 - Lower data transfer cost than VPN/internet.
- When to use: Large, constant data transfers or low-latency, reliable network needs (e.g., hybrid apps, financial trading).

Virtual Interfaces (VIFs)

- Private VIF: Connects to your VPC resources (EC2, RDS, etc.).
- Public VIF: Access AWS public endpoints (S3, DynamoDB, etc.) over DX instead of internet.
- Transit VIF: Connect Direct Connect → Transit Gateway, enabling multiple VPC/region connections.

Exam Tip:

- Private VIF = VPC
- Public VIF = Public services (S3, DynamoDB)
- Transit VIF = Many VPCs / Regions

Business Continuity Planning (BCP)

Ensures business keeps running during failures or disasters. Two key metrics:

1. RTO (Recovery Time Objective)

- **Definition**: Maximum acceptable **downtime**.
- **Example**: If RTO = 1 hour, system must be restored within 1 hour of failure.
- Factors: Criticality of the app, business impact.

2. RPO (Recovery Point Objective)

- Definition: Maximum acceptable data loss (in time).
- **Example**: If RPO = 5 minutes, backups/replication must ensure no more than 5 minutes of data is lost.
- Factors: Data change rate, backup frequency.

Memory Trick:

- RTO = "Time" until system is up again.
- RPO = "Point" in time you can roll back to.

Multi-Region Strategies

- Active-Passive:
 - One region is active. Secondary is on standby.
 - Failover happens during disaster.
 - Cost-effective, but slower recovery.

Active-Active:

- Multiple regions handle traffic at the same time.
- V Higher availability, but more complex and expensive.

Read Replicas:

- Only **read traffic** distributed across regions.
- Useful for apps with global users (e.g., read from closest region).

When to choose:

- Active-Passive → Balanced cost + availability.
- Active-Active → Mission-critical, zero downtime.
- Read Replicas → Global apps needing low-latency reads (not full DR).

Quick Decision Guide

Need	Service/Strategy
Replace file servers with S3 but still use NFS/SMB	File Gateway
Block storage with hybrid backup/DR	Volume Gateway
Replace tape backup infra	Tape Gateway
Private, consistent network to AWS	Direct Connect
Max downtime allowed	RTO
Max data loss allowed	RPO
Balanced DR with low cost	Active-Passive
Zero downtime global availability	Active-Active
Low-latency reads for global users	Read Replicas

Cost Optimization {#cost-optimization}



AWS Cost Optimization (SAA Notes)

AWS Pricing Models

1. On-Demand Pricing

- What it is: Pay only for what you use, no upfront or long-term commitment.
- Billing: Per hour or per second (depending on service).
- · Best for:
 - Short-term workloads
 - Unpredictable or spiky usage
 - Development & testing environments

Think "flexibility, but most expensive."

2. Reserved Capacity

Reserved Instances (EC2)

Standard RI:

- Highest discount (up to 72% off)
- Fixed instance type, region, and OS → less flexibility

Convertible RI:

- Lower discount (up to 54%)
- Can change instance attributes (family, OS, tenancy)

Scheduled RI:

- Reserve capacity for recurring time windows (e.g., every day 9 AM-5 PM)
- **When to use**: Steady-state workloads (databases, production apps).

Savings Plans (newer, more flexible than RIs)

- Compute Savings Plans:
 - Up to 66% discount
 - Applies across EC2, Fargate, and Lambda (most flexible).
- EC2 Instance Savings Plans:
 - Up to 72% discount
 - Locked to instance family + region (less flexible).
- **Commitment**: 1-year or 3-year terms.
- **When to use:** Long-running apps but want more **flexibility than RIs**.

3. Spot Pricing

- **Discount**: Up to **90% cheaper** than On-Demand.
- Drawback: AWS can terminate with 2-minute warning if capacity is needed.
- Best for:
 - Fault-tolerant workloads
 - Batch jobs
 - CI/CD pipelines

- Big data & ML training
- Think "super cheap, but unreliable."

Cost Monitoring & Management

1. AWS Cost Explorer

- Purpose: Visualize cost and usage trends.
- Features:
 - Interactive reports (by service, region, tags, accounts).
 - Forecasting future spend.
 - Recommendations for RIs & Savings Plans.
- **Best for:** Exploring trends & optimizing spend.

2. AWS Budgets

- Budget Types: Cost, Usage, RI, Savings Plans.
- Alerts: Send via Email or SNS when thresholds are exceeded.
- Actions: Can trigger automated responses (e.g., stop EC2, restrict IAM actions).
- **Best for**: Staying under budget & proactive alerts.

3. AWS Cost and Usage Report (CUR)

- What it is: The most detailed cost & usage dataset.
- Delivery: Sent to an S3 bucket daily.
- Integration: Can query with Athena or visualize in QuickSight.
- **Best for: Deep dive analysis** & chargeback/accounting use cases.

4. AWS Cost Anomaly Detection

- Powered by ML: Detects unusual spending patterns automatically.
- Root Cause Analysis: Shows which service/account caused the spike.
- Alerts: Notify via Email or SNS.

Quick Decision Guide

Need	Best Option
Short-term, unpredictable workloads	On-Demand
Steady-state workloads with known usage	Reserved Instances (Standard/Convertible)
Flexible long-term commitment across EC2, Lambda, Fargate	Savings Plans
Fault-tolerant, cheap computing	Spot Instances
Explore trends & cost forecasting	Cost Explorer
Stay within budget & trigger alerts/actions	AWS Budgets
Most detailed billing data for analysis	Cost & Usage Report (CUR)
Detect and alert on unusual spend	Cost Anomaly Detection

AWS Cost Optimization Strategies (Improved Notes)

Right-Sizing (Avoid Overpaying for Resources)

- CPU Utilization → Monitor with CloudWatch metrics. If consistently low → downsize instance.
- Memory Usage → Choose instance family that matches workload (e.g., R-series for RAM-heavy apps).
- Storage Performance → Pick right storage tier (e.g., gp3 instead of provisioned IOPS if not needed).
- Tip: Always start smaller, scale up only if needed.

Auto Scaling (Scale with Demand)

 Horizontal Scaling → Add/remove multiple smaller instances instead of one large one.

- Predictive Scaling → Uses ML forecasting to scale ahead of traffic (good for seasonal spikes).
- Target Tracking → Example: Keep EC2 CPU at ~50% automatically.
- Tip: Prevents paying for idle resources.

Storage Optimization

S3 Storage Classes

- Lifecycle Policies → Automatically move objects to IA (Infrequent Access),
 Glacier, or Deep Archive.
- Intelligent-Tiering → Best if access patterns are unknown/unpredictable.
- Delete Incomplete Multipart Uploads → Avoid paying for unfinished uploads.

EBS Optimization

- gp3 vs gp2 → gp3 gives higher baseline performance at lower cost.
- Volume Types → Use st1/sc1 for throughput workloads, io2 for high IOPS.
- Snapshots → Clean up unused snapshots to save cost.
- **Tip:** Match storage class to usage pattern.

Data Transfer Optimization

- CloudFront (CDN) → Cache data at edge → reduces origin fetch & transfer costs.
- VPC Endpoints → Private connection to AWS services → avoids NAT Gateway charges.
- Cross-AZ Traffic → Costs \$\$ → Place workloads in same AZ if possible.
- **Tip:** Data transfer can secretly drive up bills optimize it early.

Reserved Capacity Planning

Usage Analysis → Use Cost Explorer RI recommendations before buying.

- Commitment Strategy → Mix of On-Demand + RIs/Savings Plans for flexibility.
- Convertible Options → If future requirements may change, use Convertible RIs.
- Tip: Lock predictable workloads into 1yr/3yr discounts.

♦ Tagging Strategy

Cost Allocation Tags

- Purpose: Track spend by team, project, or environment.
- AWS-Generated Tags → Automatically applied (e.g., createdBy).
- User-Defined Tags → Custom business tags (Project=FinanceApp).

Tag Policies

- **Enforcement**: Require specific tags (e.g., CostCenter).
- Standardization: Consistent tagging across all accounts.
- Automation: Use AWS Config + Service Control Policies (SCPs).
- **Tip:** No tagging = no visibility into cost drivers.

Multi-Account Cost Management

AWS Organizations

- Consolidated Billing → One bill for all accounts.
- Volume Discounts → Usage aggregated → higher discounts.
- RI & Savings Plans Sharing → Unused discounts applied across accounts.

Cost Categories

- Custom Grouping → Group costs by department or project.
- Rules-Based → Automatically categorize based on tags, accounts, or services.
- Reporting → Custom cost reports per category.

Quick Exam-Ready Decision Guide

Optimization Area	Key Action	AWS Tools
Right-Sizing	Adjust instance/storage sizes	CloudWatch, Trusted Advisor
Auto Scaling	Scale based on demand	Auto Scaling Groups
Storage Optimization	Use cheapest storage tier	S3 Lifecycle, gp3, delete snapshots
Data Transfer	Reduce transfer costs	CloudFront, VPC Endpoints
Reserved Capacity	Commit for steady workloads	RIs, Savings Plans
Tagging	Track costs properly	Cost Allocation Tags, Tag Policies
Multi-Account	Optimize org-wide	AWS Organizations, Cost Categories

AWS Organizations – Complete Notes

1. What is AWS Organizations?

- A service for managing multiple AWS accounts in a centralized way.
- Helps with:
 - Consolidated billing (all accounts under one payer).
 - Policy control (govern what accounts can/can't do).
 - Centralized security & compliance.
 - Easier scaling (separate accounts for teams, projects, or environments).
- Think of it like a company hierarchy:
 - **Root** → Company
 - Organizational Units (OUs) → Departments
 - Accounts → Teams/Projects

Service Control Policies (SCPs) → Rules

2. Key Terminology

1. Root

- The top-level container in your AWS Organization.
- Every account belongs under the Root.
- Default root is created when you enable AWS Organizations.

2. Account

- An individual AWS account.
- Two types:
 - Management (Master) Account → Full admin, manages billing, creates OUs, applies policies.
 - Member Accounts → Created/added under Organization, controlled by management account.

3. Organizational Units (OUs)

- A group of accounts inside the Organization.
- Can be **nested** (OUs inside OUs).
- Helps organize by:
 - Function (Dev, Test, Prod)
 - Department (HR, Finance, IT)
 - Compliance levels

4. Service Control Policies (SCPs)

- Policies that define what actions accounts can/cannot perform.
- They don't grant permissions, they **restrict** them.
- Example:
 - An IAM policy says: "You can create EC2."
 - An SCP says: "EC2 creation is denied."
 - Result: EC2 creation is denied.

SCPs are applied at OU or Account level.

5. Consolidated Billing

- All accounts under the Org share one bill.
- Benefits:
 - Bulk discounts (Reserved Instances, Savings Plans).
 - Easier management (only one payment method).
- Each account still sees its own usage, but payer sees all.

3. AWS Organization Features

🔑 Account Management

- Create accounts via the Org console or API.
- Invite existing AWS accounts to join.
- Remove accounts (they become standalone accounts again).

Security Management

- Apply SCPs for compliance (e.g., block regions, deny root access).
- Centralize IAM permission boundaries.

== Billing

- One bill for the whole Org.
- Shared discounts across accounts.
- Detailed per-account cost breakdown.

Policies

- Types of Policies (only SCPs are common in exam):
 - SCPs (Service Control Policies) → Control permissions.
 - Al services opt-out policies → Opt accounts out of AWS AI/ML data usage.

Multi-Account Strategy

- Best practice: one account per workload/environment.
 - Dev, Test, Prod → Separate accounts.
 - Security, Networking, Logging → Separate accounts.
- Improves security isolation and billing clarity.

4. Service Control Policies (SCPs) - Details

- Default behavior:
 - Root has a FullAWSAccess SCP (allows everything).
- Explicit Deny wins:
 - If SCP denies something, IAM can't override it.
- No effect on management account:
 - SCPs do not apply to the management account.
- Use Cases:
 - Deny usage of certain regions.
 - Deny deletion of CloudTrail logs.
 - Restrict EC2 instance types.

Example SCP (deny use of N. Virginia region us-east-1):

AWS Solutions Architect Associate 137

}

5. When to Use AWS Organizations

- Use it when:
 - You have multiple teams/projects needing separate accounts.
 - You want centralized billing.
 - You need compliance/security enforcement.
 - You want discount sharing across accounts.
- Don't need it when:
 - You're a single AWS account user (small use cases).
 - You don't need strict separation of workloads.

6. Exam Tips (AWS SAA Level)

- SCPs do NOT grant permissions → they only restrict.
- SCPs don't apply to management account.
- · Consolidated billing is default benefit.
- Best practice → separate accounts for Dev/Test/Prod.
- RI/Savings Plans discounts are shared across all accounts.
- Al services opt-out policies → important in AI/ML use cases.

7. Real-World Example

Imagine you're running a company:

- Root: "Pillai Tech"
- OUs:
 - Production OU
 - Account: Prod App
 - Account: Prod Database

Dev/Test OU

Account: Dev Team

Account: QA Team

Security OU

Account: Logging

Account: Security Monitoring

SCPs applied:

- Dev/Test OU → Deny EC2 m5.24xlarge (too costly).
- Production OU → Deny access to us-west-1 (only allowed ap-south-1).
- Security OU → Deny deletion of CloudTrail.

Well-Architected Framework {#well-architected}

What is the Well-Architected Framework (big picture)

Think of the Framework as AWS's checklist for building reliable, secure, efficient, and cost-effective cloud systems. It's a set of *principles + best practices* organized into six pillars (Operational Excellence, Security, Reliability, Performance Efficiency, Cost Optimization, Sustainability). AWS provides a Well-Architected Tool to run reviews and get prioritized improvement suggestions.

Use it like a recipe: before you launch or when you refactor, run a Well-Architected review to find risks and fixes.

Quick mnemonic to remember the pillars

OSRPCS → "Oscar Says Real People Can Scale"

(O)perational, (S)ecurity, (R)eliability, (P)erformance, (C)ost, (S)ustainability

Pillar-by-pillar — plain-language explanation + every point you listed

Operational Excellence

Goal: Operate systems and processes that allow you to deliver business value and improve over time.

Design principles (explained)

- **Perform operations as code:** Automate runbooks, deployments, and infra changes using code (CloudFormation/CDK/Terraform + CI/CD).
- Make frequent, small, reversible changes: Deploy small updates often easier to rollback and reason about.
- Refine operations procedures frequently: Treat runbooks as living code update after incidents.
- Anticipate failure: Design and test failure modes (chaos engineering, game days).
- Learn from failures: Post-incident blameless reviews and incorporate lessons.

Best practices (Prepare / Operate / Evolve)

- Prepare: Build observability (metrics, logs, traces), document runbooks, setup alerts.
- Operate: Use automation for incident response (Lambda playbooks), automated scaling and healing.
- **Evolve:** Collect postmortems, update runbooks/automation, run periodic reviews and drills.

Key services

- CloudFormation / CDK infra as code (IaC)
- AWS Config view drift and compliance
- CloudTrail API audit logs
- CloudWatch metrics, alarms, dashboards

X-Ray — distributed tracing

Action checklist (operational)

- Store runbooks as code (Git).
- Setup dashboards + alerts for SLI/SLOs.
- · Automate common remediation actions.
- Run periodic restore tests and game days.

Exam angle

Questions will often focus on automation, observability, and iterative improvement. If problem mentions "reduce human error" think Operational Excellence → automation.

Security

Goal: Protect data, systems, and assets while enabling business needs.

Design principles

- **Strong identity foundation:** Single source of truth (IAM/SSO), least privilege, MFA.
- Apply security at all layers: Defense in depth network, host, app, data.
- Automate security best practices: Automate patching, scanning, rotation.
- Protect data in transit & at rest: TLS + KMS encryption.
- **Keep people away from data:** Use encryption, tokenization, and role-based access.
- Prepare for security events: Incident response playbooks and automation.

Best practice areas

- Identity & Access Management: IAM roles, policies, SMP, permission boundaries, SSO, MFA
- Detective controls: GuardDuty, CloudTrail, CloudWatch logs & anomaly detection
- Infrastructure protection: VPC configs, Security Groups, NACLs, WAF, Shield

- Data protection: KMS, encryption, key rotation, backups
- Incident response: Automated containment, forensic logs

Key services

- IAM / IAM Identity Center
- KMS / CloudHSM
- GuardDuty
- WAF / Shield
- Secrets Manager

Action checklist (security)

- Enforce least privilege and MFA for principals.
- Protect sensitive data with KMS; ensure backups are encrypted.
- Enable CloudTrail multi-region and send logs to a central account.
- Use GuardDuty and automated responses for suspicious activity.
- Maintain an incident response playbook and test it.

Exam angle

If the question mentions "must ensure no one can read customer data, even AWS staff," think **CloudHSM** or customer-managed keys. If it's about blocking SQL injection or web attacks \rightarrow **WAF**. If it's about detecting threats \rightarrow **GuardDuty**.

Reliability

Goal: Ensure a system recovers quickly from infrastructure or service disruptions and meets availability expectations.

Design principles

- Automatically recover from failure: Auto-heal, backups, retries.
- **Test recovery procedures:** Practice restores & failovers.
- Scale horizontally: Add parallel units rather than scaling a single node.
- Stop guessing capacity: Use Auto Scaling, use real metrics.

 Manage change in automation: Deploy using pipelines and canary/bluegreen deployments.

Best practice areas

- Foundations: Quotas, multi-AZ networking, service limits.
- Workload architecture: Decouple components (SQS), design idempotent operations.
- Change management: CI/CD, canary deploys, feature flags.
- Failure management: Backups, snapshots, cross-region DR, monitoring.

Key services

- Route 53 (health checks & failover)
- Elastic Load Balancer
- **S3** (11 nines durability)
- Auto Scaling
- RDS Multi-AZ, DynamoDB global tables, SQS (decoupling)

Action checklist (reliability)

- Use multi-AZ for stateful services; multi-region for DR needs.
- Decouple workloads and build retries with backoff.
- Automate backups & test restores.
- Implement health checks and automated failover.

Exam angle

If they ask about durable storage for backups \rightarrow S3. For decoupling and buffering between systems \rightarrow SQS. For low RTO requirements across regions \rightarrow multi-region active-active or global DB (Aurora Global DB).

Performance Efficiency

Goal: Use resources efficiently to meet system requirements and maintain that efficiency as business needs evolve.

Design principles

- Democratize advanced tech: Use managed services (RDS, Redshift, SageMaker) rather than DIY where sensible.
- Go global in minutes: Use global services/edge (CloudFront) to reduce latency.
- Use serverless architectures: Reduce operations & scale automatically.
- Experiment often: Provision test resources cheaply and measure.
- Mechanical sympathy: Choose services/instances that match workload characteristics.

Best practice areas

- **Selection:** Pick the right instance family, DB engines, or storage types.
- Review: Continuously test and right-size.
- Monitoring: Use CloudWatch metrics & traces.
- Tradeoffs: Balance throughput, latency, cost, consistency.

Key services

- CloudWatch, X-Ray monitoring & tracing
- Lambda serverless compute
- ElastiCache in-memory caching
- CloudFront CDN
- Auto Scaling performance scaling

Action checklist (performance)

- Profile workloads; choose optimized compute (Graviton) or accelerator.
- Use caching (ElastiCache / CloudFront) to reduce backend load.
- Prefer managed services for complex workloads.
- Use A/B testing and load test before production.

Exam angle

If they ask about low latency for global users \rightarrow CloudFront. If they ask about speeding up database reads \rightarrow caching (ElastiCache) or read replicas.

Cost Optimization

Goal: Run systems that deliver business value and avoid unnecessary cost.

Design principles

- Cloud financial management: designated people + processes.
- Adopt consumption model: pay for what you use.
- Measure efficiency: track business metric per dollar.
- Stop undifferentiated heavy lifting: use managed services.
- Analyze & attribute cost: tagging and chargeback.

Best practice areas

- Cloud financial management budgets, cost ownership
- Expenditure awareness Cost Explorer, Budgets
- Right-sizing & purchasing Rls, Savings Plans, Spot
- Optimize over time continuous cost review

Key services

- Cost Explorer
- AWS Budgets
- Trusted Advisor
- S3 Intelligent-Tiering
- Compute Optimizer

Action checklist (cost)

- Tag resources and enable consolidated billing (Organizations).
- Use Savings Plans/RIs for steady loads; Spot for fault-tolerant jobs.
- Right-size instances and storage; lifecycle policies for S3.
- Monitor anomalies and set budgets/alerts.

Exam angle

If the task is to reduce spend on dev/test noncritical instances \rightarrow recommend spot or scheduled RIs, or shut down unused resources. For long-term steady workloads \rightarrow RIs/Savings Plans.

Sustainability

Goal: Minimize environmental impacts of running workloads in the cloud.

Design principles

- Understand your impact: Measure carbon and energy use.
- Set goals: Define sustainability KPIs.
- Maximize utilization: Right-size & eliminate idle resources.
- Adopt new tech: Use energy-efficient instances (Graviton) and managed services.
- **Use managed services:** Shared infrastructure more efficient than singletenant on-prem.
- Reduce downstream impact: Reduce data movement; compress & archive.

Best practice areas

- Region selection (prefer greener regions).
- Software & architecture efficiency (minimize compute).
- Data minimization (store only what you need).
- Use efficient instance types (ARM/Graviton), serverless where possible.

Key services

- AWS Graviton (energy efficient)
- S3 Storage Classes (reduce storage footprint)
- Lambda / Fargate (serverless)
- Containers (ECS/EKS) for density

Action checklist (sustainability)

- Turn off non-prod resources when idle.
- Use batch/spot & serverless to increase utilization.

- Choose efficient regions/instances.
- Reduce data transfers and compress datasets.

Exam angle

Sustainability is emerging — expect questions about efficiency and choices (serverless vs always-on VMs), and picking Graviton or appropriate storage classes.

Well-Architected Tool — how to use it and what to expect

Purpose: Guided review of a workload against best practices. Generates prioritized recommendations and an improvement plan.

Assessment process (practical steps)

- 1. **Define Workload** describe architecture, goals, SLAs, data classification.
- Answer Questions the tool asks pillar-specific questions about your workload (operational practices, security controls, fault tolerance, etc.).
- 3. **Review Recommendations** tool scores and gives remediation advice per pillar, prioritized by risk.
- 4. Create Improvement Plan choose actions, owners, deadlines.
- 5. **Implement & Track Progress** re-run review when changes are made to show progress.

Benefits

- Standardized review process
- Prioritized, concrete recommendations
- Progress tracking and governance (especially useful for audits)

Best practice for using it

- Run before production and after major changes.
- Include stakeholders (security, ops, dev, business owners).
- Treat recommendations as prioritized backlog items.

Cross-pillar practices (things that help everywhere)

- Infrastructure as Code (CloudFormation/CDK/Terraform) repeatable, testable.
- Centralized logging & monitoring (CloudWatch, X-Ray) needed by Ops, Security, Reliability, Performance.
- Automation & CI/CD consistent deployments and rollback.
- Tagging & org structures cost, security, compliance management.
- **Testing & Game Days** validate recovery, security, performance.

One-page checklists you can memorize (quick study)

Operational

- IaC + CI/CD? < ✓
- Runbooks in Git?
- Dashboards + alerts? V
- Automated remediation?

Security

- Least privilege + MFA?
- KMS & encrypted backups?
- CloudTrail multi-region enabled?
- GuardDuty & alerting?

Reliability

- Multi-AZ or multi-region for critical?
- Backups & tested restores? <a>

- SQS / SNS for decoupling?
- Health checks & failover?

Performance

- Profiling & right instance types?
- Caching & CDN used?
- Autoscaling configured?
- Tracing to find bottlenecks?

Cost

- Tags + Cost Explorer configured?
- RIs/Savings plans where appropriate?
- Spot for batch?
- Lifecycle policies and snapshots pruning?

Sustainability

- Idle resources off?
- Graviton/Serverless considered?
- Data transfer minimized?

Common exam traps and how to answer

- Trap: "Use SCP to grant permissions" wrong. SCPs restrict actions; they
 don't grant privileges. Look for IAM + SCP interplay.
- Trap: "S3 vs EBS for durable backups?" Use S3 for long-term durable backups; EBS snapshots are stored in S3 but typically used differently.
- Trap: "Need global low-latency reads for DB" consider read replicas or a global DB like Aurora Global Database.
- **Trap:** "Reduce cost for dev/test" answer: automated scheduled shutdown or use Spot and smaller instance types.

Sample scenario + model answer (practice)

Q: A web app must remain available during a regional AZ outage with minimal manual intervention. What do you do?

A (short model): Deploy across multiple AZs with auto-scaling groups behind an ELB, use RDS Multi-AZ (or Aurora replicas), store state in S3/DynamoDB, implement health checks and Route 53 failover where needed; automate deployments (CloudFormation/CDK) and test failover.

Practice Questions and Exam Tips {#exam-tips}

Question Analysis Strategy

Read Carefully

- Identify Key Requirements: Cost-effective, highly available, secure, scalable
- Note Constraints: Region limitations, compliance requirements, budget
- Understand Context: Company size, technical expertise, current architecture

Elimination Process

- Rule Out Obviously Wrong: Services that don't fit requirements
- Consider Feasibility: Overly complex solutions are usually wrong
- AWS Best Practices: Choose solutions following AWS recommendations
- Cost-Effectiveness: When multiple options work, choose most costeffective

Common Question Patterns

High Availability Scenarios

- Multi-AZ Deployments: RDS Multi-AZ, ALB across AZs
- Auto Scaling: ASG for compute, Aurora for databases
- Load Balancing: ELB for traffic distribution

• Backup Strategies: Automated backups, cross-region replication

Security Scenarios

- Least Privilege: IAM roles with minimal necessary permissions
- **Defense in Depth**: Multiple security layers (VPC, security groups, NACLs)
- **Encryption**: At-rest and in-transit encryption
- **Compliance**: Service-specific compliance features

Performance Scenarios

- Caching: CloudFront, ElastiCache, DAX
- Database Performance: Read replicas, appropriate instance types
- **Network Optimization:** Placement groups, enhanced networking
- Content Delivery: CloudFront global distribution

Cost Optimization Scenarios

- Reserved Capacity: RIs and Savings Plans for predictable workloads
- Spot Instances: For fault-tolerant, flexible workloads
- Storage Tiering: S3 lifecycle policies, EBS volume types
- Right-Sizing: Match resources to actual requirements

Service Comparison Quick Reference

Compute Services

Service	Use Case	Key Feature
EC2	General compute	Full control, multiple instance types
Lambda	Event-driven, serverless	Pay per request, automatic scaling
ECS	Container orchestration	Docker support, AWS integrated
EKS	Kubernetes	Managed Kubernetes control plane
Batch	Batch processing	Automatically managed compute resources
Elastic Beanstalk	Web app deployment	Platform as a service, multiple languages

Storage Services

Service	Туре	Use Case
S3	Object storage	Web content, backup, data archiving
EBS	Block storage	EC2 instance storage, databases
EFS	File storage	Shared storage across EC2 instances
FSx	Managed file systems	High-performance workloads, legacy apps

Database Services

Service	Туре	Use Case
RDS	Relational	Traditional relational databases
Aurora	Cloud-native relational	High-performance, cloud-optimized
DynamoDB	NoSQL	High-performance, serverless applications
ElastiCache	In-memory cache	Database caching, session storage
DocumentDB	Document database	MongoDB-compatible applications
Neptune	Graph database	Social networks, recommendation engines

Exam Day Tips

Time Management

- 2 minutes per question average (130 minutes for 65 questions)
- Flag difficult questions and return later
- Don't spend too long on any single question
- Review flagged questions if time permits

Answer Strategy

- · Read all options before selecting
- Look for AWS-specific solutions over generic ones
- Consider scenario context (startup vs enterprise, budget vs performance)
- Choose most specific correct answer

Common Mistakes to Avoid

• Don't overthink simple scenarios

- AWS-managed is usually better than self-managed
- Security is paramount choose most secure valid option
- Serverless when possible for cost and scalability

Final Study Recommendations

Hands-On Practice

- Create AWS Free Tier account for practical experience
- Follow AWS Workshops for guided learning
- Build sample architectures from course examples
- Break things safely to understand failure modes

Additional Resources

- AWS Documentation: Service-specific deep dives
- AWS Well-Architected Framework: Architecture best practices
- AWS Whitepapers: Disaster recovery, security, cost optimization
- AWS FAQs: Common guestions about services
- Practice Exams: Identify knowledge gaps

Review Strategy

- Focus on weak areas identified in practice tests
- Review architecture patterns from real-world scenarios
- Understand service limits and quotas
- Know integration patterns between services

Conclusion

This comprehensive study guide covers all essential topics for the AWS Solutions Architect Associate certification. The exam tests your ability to design resilient, high-performing, secure, and cost-optimized architectures on AWS.

Key Success Factors:

- Understand AWS services deeply not just what they do, but when and why to use them
- 2. **Think like an architect** consider trade-offs, best practices, and real-world constraints
- 3. **Practice hands-on** theoretical knowledge must be supplemented with practical experience
- 4. Focus on scenarios the exam is scenario-based, not just factual recall
- 5. Stay current AWS evolves rapidly; ensure your knowledge is up-to-date

Remember: The goal is not just to pass the exam, but to become a competent AWS Solutions Architect who can design and implement robust cloud solutions. Use this guide as your foundation, but continue learning and practicing beyond the exam.

Good luck with your certification journey! 🚀

This study guide is designed to be comprehensive yet focused on exam requirements. Regularly review AWS service updates and new features, as the cloud landscape evolves continuously.