

AWS Certified CloudOps Engineer – Associate (SOA-C03) Exam Study Notes

General Tips for AWS CloudOps Engineer – Associate Exam

General Best Practices & AWS Default Choices

- **Prefer managed and automated AWS services** → e.g., use **Amazon RDS instead of self-managed databases on EC2** for easier maintenance, or **AWS Systems Manager** to automate admin tasks instead of manual scripts.
- **Use Infrastructure as Code (IaC)** → Deploy with **AWS CloudFormation or CDK** rather than clicking in the console. This ensures **consistent, repeatable environments** and avoids configuration drift.
- **Enable comprehensive monitoring by default** → **Amazon CloudWatch** for metrics/alarms and **AWS CloudTrail** for API logs **should be turned on in all accounts**. This provides visibility and audit trails for troubleshooting.
- **Design for failure and recovery** → Assume things *will* fail. Use **Multi-AZ deployments, health checks, and auto-recovery** features (e.g., EC2 Auto Recovery, Auto Scaling health checks) to minimize downtime.
- **Eliminate single points of failure** → Opt for **distributed architectures with load balancing and redundant components** (e.g., multiple AZs, multiple servers behind an ELB) for high availability by default.

Cost & Performance Optimization Strategies

- **Right-size and use the appropriate pricing model** → Regularly run **AWS Compute Optimizer** to get instance sizing recommendations. Use **Spot Instances for non-critical or batch workloads** (can save up to ~90% cost) and **Savings Plans/Reserved Instances** for steady-state workloads.
- **Optimize storage costs** → Use **S3 lifecycle policies** to move data to cheaper tiers (e.g., Glacier) for infrequently accessed data. **Prefer Amazon S3 Intelligent-Tiering** when access patterns are uncertain to automatically save costs. **Default to gp3 volumes for EBS** – they offer high performance at lower cost; only use **io2 Block Express** for mission-critical IOPS needs (databases) where ultra-low latency is required.
- **Leverage caching and CDNs** → Implement **Amazon CloudFront** to cache content at edge locations (reduces latency *and* offloads origin servers). Use **Amazon ElastiCache (Redis/Memcached)** or **DynamoDB DAX** to cache frequent read queries and reduce database load. Caching improves performance and lowers direct database costs by reducing throughput needs.
- **Use managed services that scale automatically** → Wherever possible, use services that handle scaling for you (e.g., **DynamoDB with on-demand or auto scaling capacity**, **Lambda** for serverless auto-scaling, **Aurora Serverless v2** for auto-scaling databases). This ensures optimal performance under load and avoids over-provisioning resources during low usage.
- **Monitor and optimize idle resources** → Configure **AWS Trusted Advisor** and **Cost Explorer** reports to find underutilized resources (e.g., idle EC2s, underused EBS volumes, unassociated Elastic IPs).

Delete or downsize idle resources to cut costs. Always consider the **most cost-effective storage or compute option** that meets requirements (e.g., use EFS or S3 instead of keeping large data on expensive EBS if shared access or object storage suffices).

Monitoring & Incident Response

- **Set proactive alarms and notifications** → Configure **CloudWatch Alarms** on critical metrics (CPU, memory, errors, latency) for all major resources. Tie alarms to **Amazon SNS** to get immediate notifications or trigger **AWS Systems Manager Automation** documents for self-healing. Early detection via alarms is key to preventing small issues from becoming big outages.
- **Centralize logs and enable log analysis** → Use **CloudWatch Logs** to collect logs from EC2, Lambda, VPC Flow Logs, etc., and **AWS CloudTrail** for API logs. This centralized logging makes it easier to search across logs during incidents. For container environments (ECS/EKS), ensure the **CloudWatch Agent** or Fluent Bit is forwarding container logs to CloudWatch. Leverage **CloudWatch Log Insights** or **Amazon Athena** to query logs for troubleshooting.
- **Implement automated remediation** → Where possible, automate common response actions. For example, use **Amazon EventBridge (CloudWatch Events)** rules to detect specific events or alarm states and trigger responses (like a Lambda function or Systems Manager runbook). E.g., if a CloudWatch alarm fires for “EC2 status check failed”, an EventBridge rule could automatically trigger a Systems Manager Automation to reboot the instance or a Lambda to send an alert to Ops.
- **Have a runbook for triage** → Develop a standard **troubleshooting checklist**. For instance, if an EC2 instance is unreachable: **1)** check CloudWatch **Status Checks** (to see if AWS has hardware/network issues), **2)** verify **Security Groups** (ensure inbound rules allow the traffic), **3)** check **Network ACLs** (no deny rules blocking), **4)** confirm **Route Table** and **Internet Gateway/NAT** setup for connectivity, **5)** review any recent **CloudTrail events** or **Config changes** that could have modified network settings. Following an ordered approach ensures no obvious causes are missed.
- **Practice incident response** → Regularly simulate failures (e.g., shutting down an instance, revoking IAM permissions) and practice the recovery steps. This helps ensure that in a real incident, you’re familiar with the CloudWatch dashboards, log locations, and know the correct remediation steps quickly.

Automation & Infrastructure as Code

- **Automate OS and software patching** → Use **AWS Systems Manager Patch Manager** to automatically apply patches to groups of instances on schedule, instead of manual SSH updates. **SSM State Manager** can ensure certain configurations or agents are present and correct on your instances (automated compliance).
- **Use AWS Systems Manager for operations** → Systems Manager offers a suite of automation tools: **Run Command** to execute scripts/commands on multiple instances at once, **Automation runbooks** for multi-step workflows (e.g., detect an unhealthy instance and replace it), **Parameter Store** for centralizing config settings or secrets, and **Session Manager** for secure shell access to instances *without* opening ports or managing SSH keys. **Favor Session Manager over bastion hosts** for accessing instances – it’s more secure (no exposed SSH) and auditable.
- **Infrastructure deployment pipelines** → Implement CI/CD style deployments for your infrastructure and code. For example, use **AWS CodePipeline** with **CloudFormation/CDK** to push infrastructure updates, or third-party tools (Terraform, Ansible) if your organization uses them. The key is to treat Ops changes as version-controlled code for repeatability and rollback.

- **Event-driven operations** → Take advantage of event integrations to eliminate cron jobs on servers. For example, schedule routine tasks with **EventBridge Scheduler** (instead of cron on EC2). Use **S3 Event Notifications** to trigger processing (via Lambda or Step Functions) when objects are created, rather than polling. This reactive automation reduces the need for always-running caretaker scripts.
- **Automate multi-account setups** → In AWS Organizations with multiple accounts, automate resource sharing and deployments. Use **CloudFormation StackSets** to deploy stacks across multiple accounts/regions in one go (e.g., standard CloudTrail logging, Config rules in all accounts). Use **AWS Resource Access Manager (RAM)** to share resources like Transit Gateways or License Manager across accounts rather than duplicating them.

Security & Access Management

- **Enforce least privilege everywhere** → Always start with minimal IAM permissions and add only what's needed. Use **IAM roles for EC2/Lambda access to AWS resources instead of embedding AWS keys**, and prefer **IAM Identity Center (AWS SSO)** or Federation for user login over creating long-term IAM users. Regularly review and remove unused credentials.
- **Understand IAM vs Resource policies vs SCPs** → **IAM policies** grant permissions to users/roles within an account. **Resource-based policies** (on S3, KMS, SQS, etc.) grant cross-account access or public access directly on the resource. **Service Control Policies (SCPs)** are applied at the AWS Organizations level to *restrict* what accounts (within an OU or the whole org) can do (SCPs cannot grant, only filter/deny actions account-wide). For example, to ensure no one in dev accounts can launch high-cost instances, an SCP could deny those actions regardless of any IAM permissions.
- **Lock down accounts and use guardrails** → Implement an **AWS Organizations** multi-account structure for better isolation (prod vs dev accounts, etc.) and attach SCPs for critical restrictions (e.g., deny deletion of CloudTrail, restrict resource creations to specific regions for compliance). Enable **CloudTrail and AWS Config in all accounts**, with logs aggregated to a central S3 in the security/logging account. This ensures a tamper-proof audit log across the org.
- **Data encryption is a must** → **Encrypt data at rest with AWS KMS** managed keys or customer-managed keys for sensitive data (S3 buckets, EBS volumes, RDS instances, etc. all support encryption). **Encrypt in transit** by using **SSL/TLS** everywhere (use **AWS Certificate Manager (ACM)** to provision TLS certs for your Load Balancers or CloudFront). Know how to troubleshoot if something isn't encrypted (e.g., check ACM certificate expiration, trust store issues, or KMS key policy if a service can't access a key).
- **Use AWS security services to monitor and protect** → Set up **Amazon GuardDuty** (threat detection from CloudTrail/DNS/Flow Logs) to alert on suspicious activity (like crypto-mining or unusual logins). Use **Amazon Inspector** to automatically scan EC2 instances and container images for vulnerabilities. **AWS Security Hub** can aggregate findings from GuardDuty, Inspector, IAM Access Analyzer, etc., into a single pane and even run automated remediation playbooks. Regularly review Security Hub's dashboard for any **HIGH/CRITICAL findings** and address them (e.g., unencrypted S3, open SG ports).
- **Protect secrets and keys** → Avoid storing secrets in plain text. Use **AWS Secrets Manager** for database passwords, API keys, etc., which provides encryption, rotation, and fine-grained access control. For less sensitive configuration data, **SSM Parameter Store** (with KMS encryption for secure strings) is a lightweight option. Never hardcode secrets in Lambda code or AMIs; retrieve them from Secrets Manager/Parameter Store at runtime via secure API calls.

📖 High Availability & Resiliency

- **Always design multi-AZ for critical resources** → Assume an AZ outage will happen. Deploy workloads across at least two AZs (e.g., span Auto Scaling Groups across AZs; use **RDS Multi-AZ** option for databases). Multi-AZ ensures that even if one data center goes down, your application remains available from the other.
- **Implement health checks and failover** → Use **Elastic Load Balancers** (ALB/NLB) to distribute traffic and automatically remove unhealthy instances. For global failover or latency-based routing, use **Amazon Route 53** with health checks – e.g., active-passive failover DNS records to shift traffic to a standby region if the primary fails. Regularly test these failovers to validate your RTO (Recovery Time Objective).
- **Use auto scaling for resilience and scaling** → **EC2 Auto Scaling Groups** not only scale out for performance, they can automatically replace an instance if it fails a health check (useful for self-healing). Similarly, enable **DynamoDB auto scaling** or Aurora auto scaling so the service adjusts capacity as needed. This helps absorb traffic spikes seamlessly while also replacing unhealthy nodes.
- **Backup regularly and design for recovery** → Set up **automated backups/snapshots** for data: EBS snapshots, RDS automated backups, DynamoDB point-in-time recovery, etc. Consider using **AWS Backup** to centralize and manage backup policies across services (and even cross-region for DR). More importantly, practice **restores** – know how to restore an RDS from a snapshot or recover a deleted S3 object (with versioning). Your **Recovery Point Objective (RPO)** will dictate backup frequency, and your **Recovery Time Objective (RTO)** will dictate the restore strategy (e.g., pilot light vs warm standby vs multi-site – see DR strategies below).
- **Plan Disaster Recovery (DR) according to RTO/RPO** → For critical systems, define how much downtime/data loss is acceptable. **If RTO/RPO is very low**, consider active-active multi-region or a **warm standby** (a scaled-down version of your stack running in another region, ready to scale up). **If some downtime is acceptable**, a **pilot light** (critical databases replicated and a minimal infrastructure running) or simple backup-and-restore strategy might suffice. *Example:* If RTO is minutes, use warm standby or multi-site; if RTO is hours, backups with manual restore might be okay. Always document and rehearse the DR plan.

• Disaster Recovery Strategy Comparison:

DR Strategy	RTO (Downtime)	RPO (Data Loss)	Cost/Complexity
Backup & Restore	High (hours+)	Data since last backup lost (could be hours)	Low cost, simple (cold standby). Manual intervention to restore backups.
Pilot Light	Medium (hours)	Minimal data loss (db is replicated continuously)	Moderate cost. Core systems (DB) always on in secondary region, but app servers off. Start remainder of infrastructure from AMIs/scripts on failover.
Warm Standby	Low (minutes)	Little to none (near-sync replication)	Higher cost. Smaller scale version of full environment running in secondary region. Scaled up on failure.

DR Strategy	RTO (Downtime)	RPO (Data Loss)	Cost/Complexity
Multi-Site Active	Very Low (near-zero)	None (fully synchronized)	Highest cost & complexity. Full production in two+ regions serving traffic (active-active). Automatic failover essentially instantaneous.

Rule of Thumb: Use **cost vs uptime requirements** to choose DR strategy – e.g., for a mission-critical app that can't be down, spend more on warm standby or multi-site. For a dev/test environment or non-critical workload, simple backups might be enough.

New Focus Areas in SOA-C03

(The SOA-C03 exam is an evolution of the old SysOps Admin exam, with expanded coverage of modern AWS operational tools.)

- **Containers & Orchestration now in scope** → Be prepared for questions involving **Amazon ECS, EKS, and ECR** (though not as deep as a developer exam). Know the basics of running containers on ECS/EKS, and how CloudWatch collects container metrics/logs (e.g., **CloudWatch Container Insights** or Prometheus for EKS). Understand that operations tasks like scaling clusters, checking pod health, or troubleshooting an EKS daemon are fair game.

- **Infrastructure as Code & CI/CD** → The exam expects you to know about deploying infrastructure programmatically. **AWS CDK (Cloud Development Kit)** is explicitly mentioned – know that it's an IaC tool where you can define cloud resources in code. You might get scenario questions on CloudFormation vs CDK or how to manage deployments across multiple accounts (StackSets). Also, a general understanding of **code pipelines/deployment strategies (blue-green, rolling)** is useful in an ops context for deployments and updates.

- **Advanced Observability** → In addition to CloudWatch, newer services like **Amazon Managed Service for Prometheus** and **Amazon Managed Grafana** may appear. Prometheus is used for metrics (especially in container environments) with a query language (PromQL), and Grafana is for dashboards (which can aggregate data from multiple sources/accounts). You don't need to know PromQL syntax, but be aware these exist as **AWS-managed monitoring tools for metrics and visualization** beyond CloudWatch.

- **Cloud Financial Management** → While cost optimization isn't a standalone domain now, **cost awareness is integrated throughout**. Expect that many scenario options will test if you know the cost-efficient choice. For example, if a question scenario highlights that a solution is getting too expensive or underutilized, the correct answer might be to use **Spot instances, or move data to Glacier, or enable Compute Optimizer**. Be familiar with cost tools like **AWS Budgets, Cost Explorer, and AWS Cost and Usage Reports (CUR)** for visibility. CloudOps engineers are expected to keep operations **cost-efficient** as well as reliable.

- **Multi-Account Operations** → There's emphasis on operating in multi-account setups (organizations). Know how to **monitor across accounts** (e.g., CloudWatch cross-account dashboards & metrics sharing, CloudTrail organization trails logging to central S3), how to **share resources** (RAM), and how to use organizations features like **SCPs** and **AWS Control Tower** guardrails. A CloudOps engineer should be comfortable with the idea of a landing zone/multi-account environment and tooling for governance.

Exam Strategy & Mindset

- **Read the full scenario carefully for keywords** → Exam questions are scenario-based and often long. **Identify the key requirement:** Are they asking for a *cost-optimized* solution? A *secure* solution?

A *high-performance* solution? The correct answer will directly address that key need. For example, if the question emphasizes security of cross-account access, look for an answer involving IAM roles or KMS encryption. If it mentions sudden spikes in traffic, focus on auto scaling or event-driven scaling.

- **Eliminate obviously wrong options** → Usually 1–2 answer choices can be ruled out quickly. Anything that suggests doing something *manually* that could be automated by AWS, or using an archaic/less efficient approach (e.g., a cron job on an EC2 when EventBridge could be used, or custom scripts where AWS has a service) is likely not the best choice. **AWS exams favor using managed, serverless, and automated solutions** in line with best practices.
- **Think “cloud-native” and AWS-default in solutions** → If unsure, choose the solution that **most aligns with AWS Well-Architected best practices**: e.g., monitoring = CloudWatch, automation = Systems Manager or Lambda, deployment = CloudFormation/CDK, security = IAM/KMS, logging = CloudTrail/CloudWatch Logs. The exam expects you to apply these defaults unless there’s a special requirement otherwise stated.
- **Manage your time – 65 questions in 130 minutes** → That’s ~2 minutes per question on average. Some complex scenario questions can take longer, but many will be quicker if you know the material. **Don’t get stuck too long on one hard question.** Use the flag and review strategy: answer as best as you can, flag it, and move on. It’s often effective to do a first pass answering everything you find straightforward, then use remaining time to tackle the flagged ones in depth.
- **Use context clues and double-check the specifics** → A scenario might include extraneous info. Focus on the actual problem to solve. E.g., a question might describe an architecture and then ask how to troubleshoot a specific issue (like “users cannot access an S3 website”). In that case, zero in on the S3 website setup (is it public? DNS correct? bucket policy?) rather than the other parts of the story. The exam may test if you know specific service limits or configurations (e.g., “Which CloudWatch metric would detect this issue?” or “What’s the first thing to check if X is failing?”). Make sure to review key metrics (CPUCreditBalance for burstable instances, SurgeQueueLength for ELB, etc.) and common troubleshooting steps for each service.
- **Keep security and cost in mind as hidden requirements** → Even if a question doesn’t explicitly ask about security or cost, **the best answer will not violate best practices**. For example, an option that suggests storing credentials in user data or hardcoding secrets should raise a red flag (security violation). An option that solves the problem but deploys an unnecessarily large or expensive resource when a smaller/cheaper one would do is likely not AWS-preferred. So, choose solutions that **are secure by design and cost-conscious** when all other factors are equal.
- **Be familiar with the AWS console behavior and CLI** → While the exam no longer has manual lab scenarios, it may still ask how to perform tasks. For example, “How can an Ops engineer securely connect to an EC2 in a private subnet without a bastion?” Answer: Use SSM Session Manager. Or, “What is the easiest way to collect metrics from 100+ EC2 instances?” Answer: Install and configure the CloudWatch Agent (not writing custom scripts on each). The exam favors knowing the **simplest operational method** for a given task. If you’ve practiced hands-on, these will be easier.
- **Stay calm and use the process of elimination** → If you’re unsure, eliminate what you know is wrong and then choose the *most AWS-aligned* option from the rest. There is no penalty for guessing, so never leave a question blank. Often a seemingly tough question becomes easier when you remove two wrong answers and then consider which of the remaining fits best. Trust your preparation and cloud intuition – these notes and your experience are your guide.

Domain 1: Monitoring, Logging, Analysis, Remediation, and Performance Optimization (22%)

Task Statement 1.1: Implement Metrics, Alarms, and Filters using AWS Monitoring & Logging Services

1. AWS Monitoring & Logging Services

- **Amazon CloudWatch** – The primary service for monitoring AWS resources. It collects **metrics** (CPU, network, memory via agent, etc.), **logs** (via CloudWatch Logs), and **events**. Use CloudWatch to set **Alarms** on metrics (e.g., high CPU, low free memory) and send notifications or trigger automated actions. CloudWatch also provides **CloudWatch Dashboards** to display graphs and metrics (and these dashboards can be made cross-account and cross-region to monitor multiple accounts in one view).
- **AWS CloudTrail** – Captures **API call logs** for your AWS account. CloudTrail is essential for **auditing and security analysis**. When troubleshooting “who did what” or investigating unexpected changes, CloudTrail logs will show the user, time, IP, and request details for API calls. Best practice: **Enable CloudTrail in all regions and all accounts**, with logs centralized to an S3 bucket (ideally in a separate audit account).
- **CloudWatch Logs & Logs Insights** – CloudWatch Logs can ingest logs from various sources: EC2 instances (via CloudWatch Agent), Lambda logs, VPC Flow Logs, Route 53 Resolver query logs, etc. **Create log groups** per application or service. Use **Logs Insights** (a query engine) to search and analyze logs with SQL-like queries – very useful for pinpointing errors or trends in huge log datasets.
- **CloudWatch Agent** – An agent installed on EC2 (or on-prem servers) to collect **system-level metrics** (like memory, disk, swap) and additional logs. For example, you’d install the CloudWatch Agent on an EC2 to push application logs or to get memory % since the basic EC2 metrics (CPU, network, disk IO) are provided by default but not memory. Also, in container environments (ECS tasks or EKS pods), you can run the CloudWatch Agent or Fluent Bit to send container logs and custom metrics (or use Container Insights which automates this).
- **Amazon Managed Service for Prometheus (AMP)** – A fully-managed Prometheus-compatible monitoring service. Prometheus is popular for **container and microservice metrics** (scrapes metrics in a time-series database and allows complex queries). In AWS, AMP lets you ingest Prometheus metrics without managing the servers. Likely usage: if you have EKS clusters using Prometheus metrics (e.g., from Kubernetes pods), you can use AMP to store and query those metrics.
- **Amazon Managed Grafana** – Managed Grafana for dashboards and visualization. Grafana can pull data from CloudWatch, Prometheus, and other sources to build rich dashboards. In a multi-account environment, Grafana can be used to create a single dashboard viewing metrics across accounts. For the exam, remember that Grafana is about **visualizing and analyzing metrics/logs** and is often paired with Prometheus for advanced use cases.

CloudWatch Alarms & Events:

- **Alarm Actions** – CloudWatch Alarms can be set to perform actions **automatically**. They can **notify via**

SNS, trigger an **Auto Scaling policy** (e.g., add instances when CPU > 80%), or **invoke EC2 actions** (recover, reboot, stop an instance if it fails health checks). They can also trigger **Amazon EventBridge rules** (CloudWatch Events) on alarm state changes. *E.g.* You might configure a CloudWatch Alarm on an RDS metric like FreeStorageSpace; if it's below threshold, it triggers an EventBridge rule that invokes a Lambda to clean up old data or send a detailed report.

- **Composite Alarms** – A composite alarm combines multiple alarms (AND/OR conditions). For example, you can have a composite alarm that goes to ALARM state only if **both** CPU is high **and** memory is high across an instance group, to reduce noise. Composites don't have their own metrics; they just evaluate other alarms. They are useful to make alerts more meaningful (requiring multiple signals).

- **EventBridge (CloudWatch Events)** – Think of this as the **central event bus** for AWS. It can respond not just to alarms, but to many AWS events (e.g., an EC2 instance state change, an ECR image push, an Auto Scaling event). EventBridge rules allow filtering events and routing them to targets: e.g., **invoke a Lambda, run an SSM Automation, or send to an SNS topic**. For CloudOps, EventBridge is key for building **automated remediation**: if X happens, do Y. *Example*: If an instance goes into "impaired" status, EventBridge can automatically trigger a recovery action or notification.

Monitoring Best Practices:

- **Use Tags and Filter Logs** – Ensure resources are tagged (Environment=Prod, Application=XYZ) because CloudWatch and other tools can filter/alarm based on tags. Also, create metric filters for logs if needed. For instance, you can create a CloudWatch Logs Metric Filter to count occurrences of the word "ERROR" in your app logs and create a metric from it, then alarm on that metric. This way you get notified on application errors, not just infrastructure metrics.

- **Metric Granularity** – Understand default metric intervals: Standard AWS metrics (like EC2 CPU) are 5-minute by default, but can be 1-minute with detailed monitoring (enabled by default in many services now). For high-resolution needs, custom metrics can be as low as 1-second resolution. The exam may ask about ensuring you have the right granularity (e.g., if 5 min is too coarse, enable detailed 1-min or use custom metrics).

2. Automated Remediation & Event-Driven Response

- **Amazon EventBridge for auto-remediation** – Identify critical events that should trigger an automated fix. Example scenarios: an EC2 instance enters a stopped state unexpectedly (trigger EventBridge to send an SNS alert or trigger Lambda to start it if needed), or an ACM certificate is about to expire (trigger a notification), or a specific CloudTrail security event occurs (trigger an SSM Automation to lock something down). EventBridge can ingest events from almost all AWS services and even on a schedule (Cron).
- **AWS Systems Manager Automation Runbooks** – These are predefined or custom workflows of actions in AWS. For example, the AWSSupport-TroubleshootEC2Network runbook might check common network misconfigurations on an EC2. You can have runbooks to restart services, take memory dumps, patch an instance, etc. CloudOps engineers should know how to execute automation documents (via Console, CLI, or even triggered by EventBridge). If a question scenario describes a repetitive multi-step fix, consider **SSM Automation** as the answer.
- **Example Use Case – EC2 Recovery**: If an EC2 fails a status check, you could use a CloudWatch Alarm that triggers an **Auto Recovery** action (for certain instance types) or triggers an SSM Automation that stops and starts the instance or executes predefined fixes. Know that **EC2 Auto Recovery** can be set for certain instance types which essentially is an automated instance reboot/refresh if the underlying hardware fails (no user intervention).

- **AWS Config Rules & Remediation** – Config isn't explicitly named in the domain outline, but it plays a role in continuous compliance. You might get scenarios where certain configurations must always be true (e.g., S3 buckets must not be public). **AWS Config Rules** can detect non-compliance, and with **Config Remediation** (often via SSM Automation or Lambda), automatically revert changes or fix configurations. Example: a Config rule detects an unencrypted EBS volume and triggers an automation to snapshot it and re-encrypt or alert you. This is an advanced form of automated remediation beyond CloudWatch events.
- **User Notifications** – AWS recently added a feature called **AWS User Notifications** (in the console) which can aggregate alerts and service health events. While not heavily covered, just be aware that it's another way to get notifications in the Console for certain events (like Personal Health Dashboard events). For exam purposes, SNS and email alerts are the main notification method, but keep in mind Personal Health Dashboard for AWS outages or service issues that might impact your resources.

3. Performance Optimization Techniques

- **Compute Performance** – Use metrics and AWS tools to identify performance bottlenecks: e.g., **CPUCreditBalance** for burstable T instances (to see if CPU throttling), **SurgeQueueLength** for ELB (if >0, ELB can't keep up with connection attempts), or **ConsumedRead/WriteCapacity** for DynamoDB (to see if you're throttling). For EC2, if CPU is high, scale out via ASG or scale up to a larger family. If memory is the issue (found via CloudWatch Agent metrics), consider memory-optimized families. **AWS Compute Optimizer** will recommend instance types if your workload is consistently under or over-utilized (including suggesting newer-gen instance types for better price/performance).
- **EBS Volume Optimization** – Monitor **VolumeReadOps/WriteOps** and **VolumeQueueLength** in CloudWatch. If an EBS volume is a bottleneck (high queue, maxed IOPS), consider switching to a higher-performance type: e.g., General Purpose volumes (gp3) allow provisioning additional IOPS and throughput beyond the baseline, whereas if you need extremely high IOPS with sub-millisecond latency for a database, **io2** volumes may be used. Also distribute I/O across multiple volumes if possible (RAID 0 for throughput, though rarely needed with gp3's high limits). Keep an eye on **EBS Burst Balance** for gp2 (gp3 doesn't have burst credits).
- **S3 Performance** – For very high throughput to S3, consider **S3 Transfer Acceleration** (uses CloudFront's edge network to accelerate uploads to S3), or **Multipart Upload** for large files (parallelize upload parts). S3 can handle at least 3,500 PUT/COPY and 5,500 GET per second per prefix; if you need more, use **multiple prefixes** (keys with different prefixes distribute load). In general, S3 performance these days scales automatically, but older guidance about prefix spread might appear. Also, use **S3 Intelligent-Tiering** to optimize cost/performance for unpredictable access—frequently accessed objects stay in frequent tier, others move to IA or Archive tiers automatically (ensuring you're not paying for performance you don't need).
- **Shared Storage (EFS/FSx)** – If using **Amazon EFS** for concurrent access (e.g., a fleet of web servers needing a common file share), remember performance modes: **General Purpose vs Max I/O** (GP is usually fine, Max I/O for extreme scale with some latency tradeoff). EFS has bursting throughput based on size, and allows **Provisioned Throughput** if needed. Also implement **EFS Lifecycle Management** to move infrequently used files to EFS Infrequent Access to cut costs by ~92% for those files. For high-performance computing or ML workloads needing a fast file system, use **Amazon FSx for Lustre** which provides very high throughput and integrates with S3 (good for big data). In a question, if they mention needing shared POSIX storage for multiple instances or high-performance scratch space, think EFS or FSx depending on requirements (EFS for gen-purpose, FSx for high I/O specialized needs).

- **Database Performance** – Utilize **Amazon RDS Performance Insights** (supported on many RDS engines) to identify slow queries or resource bottlenecks at the DB level. If CPU or connections are maxing out on RDS, scale up instance class or add **Read Replicas** (for read-heavy loads) or use **Aurora Auto Scaling** for replicas. Consider **Aurora Serverless** for spiky workloads (it auto-scales capacity in fine-grained increments). For DynamoDB, watch **ThrottledRequests** and **ConsumedCapacity**; use **DynamoDB auto scaling** or on-demand mode to adjust throughput. Add **DynamoDB DAX** if you have hot keys or require microsecond read latency (DAX is an in-memory cache for Dynamo). For caching in front of relational DBs, use **ElastiCache Redis** to cache frequent queries or session data (reducing burden on the DB).
- **Optimize Networking** – Performance optimization isn't just about compute and storage. If you see issues with throughput, consider **Enhanced Networking (ENA)** on EC2 (most new instance types have it by default for higher PPS). Utilize **Placement Groups** (Cluster strategy) if low network latency or high throughput is needed between certain instances (like a tightly coupled HPC workload). Use **AWS Global Accelerator** for reducing latency for global users (GA will route user traffic to the nearest AWS edge, then over AWS's backbone to your application). If an app is chatty between regions, maybe re-architect to keep traffic in one region or use caching to minimize cross-region latency.

Rule of Thumb: Always identify the **bottleneck** (CPU, memory, disk, network, or application-level) and then apply the appropriate AWS solution: scale up/out for CPU or memory, optimize IOPS or switch storage for disk, improve network path or throughput for network issues, and use caching or code optimization for application issues. AWS provides specific tools for each (Compute Optimizer, PI, CloudWatch metrics) – use them rather than guessing.

Exam Rules of Thumb – Domain 1

- **For monitoring any AWS service** → Use **CloudWatch** (collect metrics, set alarms, view dashboards). If it's about API activity or detecting *who did X*, the answer is **CloudTrail**.
- **For aggregating or analyzing logs** → Use **CloudWatch Logs Insights** or **Athena** (CloudTrail logs can be sent to S3 and queried by Athena). For container logs/metrics, think **CloudWatch Agent/ Container Insights** or **Prometheus** as the specialized tool.
- **For alerting and notifications** → Use **CloudWatch Alarms with SNS** (don't custom-build your own cron or mail system for alerts).
- **When an incident happens** → **First place to check is CloudWatch metrics and relevant logs** (e.g., if an app is slow, check EC2 CloudWatch metrics, then app logs in CloudWatch Logs; if a resource changed, check CloudTrail).
- **Automated Healing > Manual Intervention** → The correct answers often use EventBridge or Auto Scaling to automate a healing action, **instead of telling an admin to go click a restart**. E.g., use an Auto Scaling health check to replace a bad instance, or EventBridge to trigger a Lambda to fix a known issue.
- **For performance issues:**
 - High CPU/memory → **Scale out or up**, or offload work (e.g., use SQS to decouple, cache results).
 - EBS throughput issues → **Consider gp3 with provisioned IOPS, or io2 for extreme cases**. Striping multiple volumes (RAID 0) can boost throughput for big data if needed.
 - Database slow → **Identify queries with Performance Insights**, add **read replicas or caching**, scale vertically if needed, or consider **Aurora (if MySQL/Postgres)** for better scaling.
 - Large-scale analytics → Use the right tool (don't try to make RDS do Redshift's job, etc.).

- **Cost optimization is part of ops** → If a monitoring question hints at costs (e.g., “the logs are expensive to store”), a valid answer could be **use log filtering or retention settings** or **move logs to S3 Glacier via lifecycle**. Don’t always assume more data is better – operational excellence includes turning off unnecessary logging and rightsizing resources.
-

Domain 2: Reliability and Business Continuity (22%)

Task Statement 2.1: Implementing Scalability and Elasticity

1. Scaling & Elasticity

- **EC2 Auto Scaling** – Set up Auto Scaling Groups (ASG) for EC2 instances behind load balancers. **Dynamic Scaling** can be based on CloudWatch metrics (e.g., CPU > 70% add 1 instance, or custom metrics like request count). **Target Tracking scaling** is a common, easy method (e.g., keep CPU at 50% and ASG will add or remove instances to maintain that target). **Scheduled Scaling** can be used for known daily cycles (scale out at 8am, in at 8pm), and **Step Scaling** for more complex rules. An exam scenario might describe unpredictable workloads – the answer would be **target tracking auto scaling** for simplicity, or if they have specific times, then scheduled scaling.
- **Managed Service Scaling** – Many AWS managed services scale automatically or have their own scaling settings: **AWS Lambda** scales by invoking more instances of your function automatically (just ensure concurrency limits are set appropriately). **DynamoDB** can be set to **on-demand capacity** (auto scales internally, cost per request) or **provisioned with auto scaling** for predictable usage. **Aurora** (and RDS to some extent) can **Auto Scaling read replicas**, adding replicas based on load. **ECS** can do Service Auto Scaling (scale task count based on CPU, etc.). The key is: use the service’s built-in scaling. For example, don’t manually launch more RDS instances—use read replicas or Aurora Serverless (which can scale the capacity units).
- **Elastic Load Balancing** – While not “scaling” itself, ELBs (ALB, NLB, CLB) enable elasticity by distributing traffic and handling varying loads. Ensure **sticky sessions** only if needed (ALB can stick by session cookie, but that can reduce effective load distribution). For reliability, use **multiple AZs** in your load balancer target group so traffic is balanced across AZs. If an AZ goes down, ELB will stop sending traffic there. In an exam question about ensuring an application can handle a spike, an answer likely involves an **Auto Scaling Group + ALB** combination.
- **Stateless architectures** – To scale horizontally (increase instances seamlessly), apps should be stateless where possible (store session state in DynamoDB/ElastiCache, store files in S3/EFS rather than local disk). A question might imply that the current architecture is not scaling due to stateful issues; the solution could be to **externalize state to a shared store**, enabling adding more instances easily.
- **Managed Database Scaling** – For RDS, you can scale **vertically** (bigger instance) or **horizontally** (add Read Replicas for read traffic scaling). DynamoDB is horizontal by nature; ensure the partition key is well-designed to avoid hot partitions under scale. With Aurora, you get up to 15 replicas that share storage, and Aurora can also auto-scale the **capacity of the writer** with Aurora Serverless v2 (scales seamlessly in fine-grained increments). If an exam asks how to handle bursty database

workloads at low cost, **Aurora Serverless** might be the answer since it can scale down to zero (for infrequent use) and up quickly when needed.

- **Caching for scalability** – Using caches isn't just performance, it's scalability. **Amazon CloudFront** caches static content at edge locations globally, offloading origin servers (meaning your web tier or S3 doesn't get hammered on peak). **ElastiCache** (Redis/Memcached) can absorb frequently repeated read requests (e.g., instead of hitting a database every time for the same data, check Redis first). This allows you to handle more users with the same backend capacity. If a question scenario involves repeated lookups or computational results, adding a cache layer is often the right move to improve scalability and user latency.
- **Event-driven scaling** – Use queueing and messaging to handle burst loads gracefully. For example, if an application suddenly gets 1000 requests, rather than processing all in real-time, use **Amazon SQS** to queue work and have a fleet of consumers (which can auto scale based on queue depth). Or use **AWS Lambda**, which can scale out concurrency extremely quickly (Lambda can handle sudden bursts by spawning many instances in parallel). The key is decoupling: so bursts are smoothed out and not overwhelming any single component.

2. High Availability (HA) & Resiliency

- **Multi-AZ and Failover** – Many AWS services have built-in Multi-AZ. **RDS Multi-AZ** (for supported engines) keeps a synchronous secondary in another AZ; failover happens automatically on primary loss (typically ~60-120 seconds downtime). **ElastiCache (Redis)** can run in primary-replica with Multi-AZ and auto-failover. **EFS** is regional (data stored redundantly across AZs by default). For systems that don't automatically failover (like EC2 or ECS tasks), design a fallback: e.g., run redundant EC2 instances in 2+ AZs behind an ALB, or for ECS use **multiple AZs for your cluster** and enable service auto scaling to replace tasks if a node fails.
- **Route 53 for Resilience** – Use **Route 53 health checks and routing policies** for certain HA scenarios. Example: a web application deployed in two regions (active-passive) – use a **Failover Routing Policy** in Route 53 with health checks on the primary site. If primary is unhealthy, Route 53 will send traffic to the secondary. Or use **Latency-based Routing** to send users to the closest region (which also provides an HA benefit if one region goes down, users automatically go to the other because the unhealthy region's endpoints fail health checks).
- **Circuit Breaker Patterns** – Although not an AWS service, it's worth noting operational patterns: e.g., if a downstream service is failing, sometimes it's better to **fail fast and show a fallback** than to keep retrying and hanging threads. AWS doesn't have a named "circuit breaker" service, but tools like AWS App Mesh for microservices allow setting health checks and failure handling. The exam likely won't call it "circuit breaker" but might hint at "stop sending requests to the unhealthy service for a period" – which essentially is what Route 53 or ELB health checks do at the infrastructure level.
- **Fault Isolation** – Design using **AWS Fault Isolation units**: AZs are isolated units, as are Regions. If your application can tolerate regional outages (very rare but possible), consider multi-region active-active or active-backup (with Route 53). If the cost or complexity of multi-region is too high, at least ensure **regular backups to another region** (so you can restore if region-wide catastrophe occurs, meeting business continuity). For AZ level, always deploy multi-AZ. Additionally, decouple components such that a failure in one does not cascade: e.g., use SQS between microservices so if one fails, messages queue up instead of crashing everything.
- **Testing Resiliency** – Practice **chaos engineering** on a small scale: you could be asked conceptually how to test if your system is highly available. The answer might be to perform game days or

simulate failures (shut down instances, kill a zone) to see if auto healing works. Also ensure **graceful shutdown** for instances (e.g., handle SIGTERM in apps so that when an instance is removed by auto scaling, it finishes in-flight requests instead of dropping them).

3. Backup, Restore & Recovery

- **AWS Backup** – A fully managed, policy-based backup service. You can define backup plans (frequency, retention) and apply them to AWS resources like EBS volumes, RDS, DynamoDB, EFS, FSx, etc. It can coordinate cross-service backups and store them in a backup vault. This is easier than writing Lambda scripts to snapshot multiple services. If a question involves “ensure backups across all these services with central management,” **AWS Backup** is the likely answer.
- **Point-In-Time Recovery (PITR)** – Supported by some services: **DynamoDB** PITR can restore to any second in the last 35 days. **RDS** (with automated backups on) allows PITR up to the retention (by combining daily snapshot and transaction logs). If you need to recover a database to right before an accidental deletion, PITR is the feature. Ensure it’s turned on where needed.
- **Snapshots vs Backups** – Know the difference: **EBS Snapshots** are incremental block storage backups (stored in S3 internally). They can be automated with Data Lifecycle Manager or AWS Backup. **RDS Automated backups** create daily snapshots plus transaction logs for PITR. **Manual snapshots** of RDS and EBS are user-initiated and retained until deleted. **S3 Versioning** is a form of backup for S3 objects (keeps old versions when things change or get deleted). If the scenario is “we need to be able to recover an object if deleted or modified”, answer: **Enable S3 Versioning** (and perhaps MFA Delete for extra protection).
- **Restoration Methods** – If an entire instance fails, you can **restore EBS volumes from snapshot** (and reattach or recreate an instance). For RDS, you **restore from a snapshot** which creates a new DB instance – important: the restore doesn’t overwrite the existing DB, it makes a new one from that point in time. For DynamoDB, PITR restore creates a new table with the data at that time. Thus, part of DR is planning the restore: e.g., have infrastructure-as-code ready to spin up new resources from backups. The exam might test knowledge that, say, **RDS restore requires you to point your apps to the new endpoint** because the endpoint will differ (or you promote a read replica to master, which has a different endpoint, etc.).
- **Cross-Region Backups** – Consider if you need region-level disaster recovery. Many services let you copy backups to another region: e.g., **RDS snapshots** can be copied to another region, **EBS snapshots** can too (even automatically via AWS Backup or event scripts), **AWS Backup** can directly create cross-region backups if configured. DynamoDB global tables automatically replicate data to other regions (that’s another strategy for DR with near-zero RPO for Dynamo). If business continuity requirements specify surviving region outage, ensure backups or replication to a secondary region is happening.
- **AWS Disaster Recovery Whitepaper knowledge** – AWS often references strategies like Backup & Restore, Pilot Light, Warm Standby, Multi-site (as we tabled above). Recognize them in scenarios. E.g., if a question describes a minimal environment in second region (just database replication and minimal servers off), that’s **Pilot Light**. If it says a scaled-down full copy of prod in another region running at all times, that’s **Warm Standby**.

Exam Rules of Thumb – Domain 2

- **For sudden traffic spikes** → **Auto Scaling + Load Balancer** is the go-to. (EC2 ASG for web/app tier, or use serverless like Lambda which auto-scales by default). Never answer “manually launch more instances” – it should be automated.

- **For scaling databases** → **Read-heavy?** Add read replicas or caching. **Write-heavy?** Scale up instance, or sharding (in DynamoDB use partition keys, in RDS perhaps use Aurora or sharding architectures). **Highly variable workload?** Aurora Serverless or DynamoDB on-demand.
 - **High availability** → Achieved by **Multi-AZ (within a region)** and **multi-region** for extreme cases. If the question mentions an uptime SLA or “must withstand AZ outage” – ensure solution uses at least 2 AZs. If “must withstand region outage” – a multi-region active-passive or active-active approach with Route 53 is needed.
 - **Stateful vs Stateless** → If scaling is an issue and they hint at session or state, answer likely involves **externalizing state** (use S3/EFS for shared storage, ElastiCache or Dynamo for sessions). AWS best practice is design stateless servers behind ELB.
 - **Backups** → By default, enable them! If you see a scenario where data could be lost, the answer is **enable automated backups or versioning**. For example: “We accidentally deleted some data from DynamoDB” → If PITR was enabled, you can restore. So the takeaway: ensure backups (or point-in-time restore features) are enabled for critical data stores (RDS, Dynamo, etc.).
 - **RTO vs RPO tradeoff** → Low RTO (quick recovery) usually means running something pre-prepared (pilot light or warm standby). If RTO can be high (like 24 hours), then simple backup to S3 and manual restore might be fine. If RPO must be zero (no data loss), you need synchronous replication (multi-AZ or multi-region with database replication). Slight data loss acceptable (RPO minutes) you can use regular backups or async replication.
 - **Service limits in scaling** → Recognize soft limits: e.g., API Gateway has a default throttle, Lambda concurrency default limit, ASG scaling speed (cooldowns). The exam might not dive deep, but if an option suggests increasing a service quota to improve reliability (like increasing EC2 instance limit), that could be valid. Generally, though, default limits are high enough for associate scenarios except maybe “hundreds of Lambdas were throttled -> raise concurrency limit”.
 - **Use global services for global problems** → For example, if needing to distribute traffic globally or balance load across regions, consider **Route 53**, **CloudFront**, or **Global Accelerator**. If needing to share data globally, maybe **S3 Cross-Region Replication** or **DynamoDB global tables**. A wrong answer would be something like “set up your own DNS server on EC2 for failover” – Route 53 should do that.
-

Domain 3: Deployment, Provisioning, and Automation (22%)

Task Statement 3.1: Provisioning and Managing Cloud Resources (Infrastructure Deployment)

1. Infrastructure Provisioning & Deployment

- **Amazon Machine Images (AMIs)** – Custom AMIs allow faster provisioning of EC2 instances with pre-baked software. Use **EC2 Image Builder** to automate the creation and patching of custom AMIs (it can run scheduled pipelines to create an AMI with latest patches, and run tests on it). The exam might describe needing up-to-date hardened server images – the answer: **EC2 Image Builder**

(rather than manually updating instances each time). Also, know how to share AMIs across accounts (AMI permissions or copying to another account/region).

- **AWS CloudFormation** – The core AWS IaC service. Define resources in YAML/JSON templates and deploy as a **stack**. Know concepts: **Stacks** (collection of resources managed together), **StackSets** (deploying stacks to multiple accounts/regions), **CloudFormation Drift Detection** (detect if someone changed a resource outside CF). If a deployment fails, CloudFormation can roll back automatically. Common issues: attempting to delete a stack with dependencies (you might need to retain or delete resources manually), or stack updates failing due to a changed resource – you may need to **import or adjust resources**. For the exam, key point: use CloudFormation for repeatable, auditable deployments (versus manual).
- **AWS CDK (Cloud Development Kit)** – An abstraction on top of CloudFormation that lets you write code in languages like Python, Typescript, etc. to define infrastructure. The CDK code synthesizes to a CloudFormation template. You should know that CDK exists and is used when teams prefer writing infra as code in a familiar programming language with constructs, rather than raw YAML. If a question suggests developers want to manage AWS resources using Python code, CDK is the likely answer.
- **Terraform (3rd party)** – AWS allows use of third-party IaC like Terraform. While the exam won't test Terraform specifics, they acknowledge its usage. You might see an answer mentioning Terraform in a valid context like "if a company already uses Terraform, continue to use it (AWS is interoperable)" or CDK has Terraform integration, etc. But in general, AWS's answer to IaC is CloudFormation/CDK, so prefer those in answers unless the scenario clearly indicates a third-party tool is in place.
- **Resource Access Manager (RAM)** – Allows you to share certain resources across accounts (without needing to deploy duplicates). Resources like **Transit Gateway, License Manager configurations, Route53 Resolver rules, and some new AWS services** can be shared. In an org with multiple accounts, instead of provisioning separate TGWs in each, you can create one in a networking account and share it via RAM. So for cross-account resource sharing scenarios, think RAM.
- **Multi-Account Deployments** – For pushing infrastructure to many accounts, **CloudFormation StackSets** is the native solution. It uses a management account to push stack instances to target accounts (with an IAM setup for the StackSet roles). If an exam case says "Ensure all new accounts have a set of standard resources X, Y, Z," one approach is using **StackSets with AWS Organizations integration** to auto-deploy a baseline stack to each new account. Alternatives could involve Control Tower (which uses StackSets under the hood) or service catalog, but likely StackSets is expected.
- **Common Deployment Issues:**
 - *Subnet sizing:* If CloudFormation tries to create more EC2 instances or ENIs than the subnet IP range allows, it fails. Always ensure subnets have sufficient IPs for scaling (plus extras for overhead). If an exam question mentions "template works in dev but not in prod region", maybe the prod subnets are too small. Solution: enlarge subnets or adjust IP ranges.
 - *IAM permissions:* CloudFormation might fail because the current IAM role doesn't have rights to create something (like an S3 bucket or a Lambda execution role). The error usually points to AccessDenied. Solution: update the execution role or assume role with needed perms. If a stack fails, check the events – the exam may test if you know to look at CF Events to troubleshoot resource creation failures.
 - *Service limits:* E.g., trying to create too many EIPs or NAT Gateways at once and hitting account limits can cause deployment to fail. The solution is often to request a limit increase or adjust architecture

(e.g., do we need so many NATs?). But on exam, if it's a straightforward limit, the answer might be "request service quota increase via Support".

- **Deployment Strategies (for code and infrastructure):**

- **All-at-once (Big Bang)** – Deploy everything new at once (downtime likely). Not ideal for production if it can be avoided.
- **Rolling Update** – Gradually replace instances or tasks with new version, e.g., update 2 instances at a time out of 10, until all are updated. This maintains some capacity but can be slower. In ASG, you can adjust **MaxSurge** (how many extra to launch) and **MinInService** (how many to keep running) to control rolling deployments.
- **Blue/Green Deployment** – Spin up a new fleet (green) separate from the current (blue), then cut traffic over (often via ELB swap or Route 53 DNS switch). This offers near zero-downtime and easy rollback (just revert traffic to blue if issues). In AWS, **CodeDeploy** supports blue/green for ECS, Lambda, and EC2 auto scaling groups; **Elastic Beanstalk** has blue/green environment swap; for databases, blue/green might mean a new cluster and then cutover. For CloudFormation, the *stack update* itself isn't blue/green (it's in-place unless you use techniques like creating new resources and deleting old ones manually or using ChangeSets with replacements).
- **Canary or A/B** – Release to a small percentage of users then full. Usually mentioned in context of API Gateway or Lambda alias routing or just a manual process. For Ops exam, probably less focus, but if a question is around minimizing impact of a new deployment, an answer might involve deploying to a small set first then increasing (which is essentially canary).

Bottom line: To avoid downtime, prefer **blue/green or rolling deployments** over in-place updates for services that support it.

- **AWS CodeDeploy & CodePipeline** – While this exam isn't DevOps Pro, as a CloudOps engineer you should know basics: **CodePipeline** orchestrates CI/CD (source, build, deploy stages). **CodeDeploy** can deploy application code to EC2/on-prem or Lambda or ECS with different strategies (In-Place or Blue/Green). For example, CodeDeploy can do rolling updates on EC2 with health check tracking and automatic rollback on failure. A question might not go deep but could ask how to automate deployments of a new version of an app to hundreds of EC2 instances – **Answer:** Use CodeDeploy with a proper deployment group, instead of manually SSH or user-data scripting each instance.

2. Deployment Troubleshooting & Optimization

- **CloudFormation Change Sets** – Allows previewing what changes a template update will do (which resources will be replaced, modified, deleted). If worried about accidental deletions or replacements, use change sets. If exam says "Ops team is afraid to update a stack because it might replace resources," the answer: **create a Change Set** to review changes first.
- **Stack Policies** – You can set a Stack Policy to protect certain resources in a CloudFormation stack from updates (e.g., don't accidentally replace the database). If a scenario involves ensuring a critical resource isn't impacted by a CF update, mention stack policy or at least the concept of not replacing that resource (and maybe update it manually if needed).
- **Parameter Store/Secrets in CloudFormation** – One challenge is handling secrets/config in templates. Usually use **SSM Parameter Store** or **Secrets Manager** and reference those in templates instead of hardcoding. If deployment fails due to missing secrets, ensure they have been created and the CF execution role can access them.

- **Permissions for Deployment Services** – CodeDeploy and others require IAM roles (like CodeDeploy EC2 agent needs a role to pull from S3, CodePipeline needs roles for each action). If a deployment isn't working, often it's an IAM issue (CodePipeline can't invoke CodeDeploy, etc.). So check those roles and trust relationships.
- **ECS Deployments** – If deploying new container versions, know the ECS deployment types: Rolling update (default for ECS services) vs Blue/Green (if using CodeDeploy or ECS deployment circuit breaker). ECS rolling will stop/start tasks gradually. If an ECS service is not reaching desired count during deploy, could be failing health checks (maybe the new container is bad – ECS will rollback after a few failures by default if **deployment circuit breaker** is enabled).
- **Common Errors:**
- *"UPDATE_ROLLBACK_FAILED" in CloudFormation:* sometimes manual intervention (like stack import or contacting AWS Support) needed. But if exam mentions stack in inconsistent state, could suggest **deleting the stack and redeploy** as last resort (with careful planning).
- *"Timeout" on waiting for resource:* Usually an issue with dependencies – e.g., an EC2 in CF stack stuck CREATE_IN_PROGRESS if user-data script never signals success (maybe missing cfn-signal). Use **cfn-signal** or assure the instance resource creation policy is correct. The exam might not go this deep, but be aware of how CF coordinates with EC2 creation using helper scripts.
- *Permissions issues on deployment (Access Denied):* Usually means the role or the service principal doesn't have needed IAM policy. E.g., CodeDeploy's service role must allow it to read the S3 bucket with artifacts and interact with EC2 or ECS.

Task Statement 3.2: Automating the Management of Existing Resources

1. AWS Systems Manager & Resource Automation

- **AWS Systems Manager (SSM)** – This is the Ops hub. Key components to know:
- **Session Manager** – as mentioned, for shell access without SSH. If a question says "securely manage EC2 instances in private subnets without bastion," Session Manager is the answer.
- **Run Command** – Execute commands/scripts on a fleet of instances (supports targeting by tags, specific instance IDs, resource groups). This is great for ad-hoc or scheduled tasks (e.g., run a script on all web servers to clear cache). This requires the instances to have the **SSM Agent** installed (on Amazon Linux, Windows, etc.) and an IAM role (AmazonEC2RoleforSSM managed policy or similar) allowing SSM access.
- **Automation** – Already covered that you can run automation documents (AWS-provided like restarting EC2, or custom ones) to do multi-step tasks. For example, one automation document might take a snapshot of an EBS, then create a new volume from it and attach somewhere (multi-step). The exam may expect you to choose an SSM Automation document for tasks like patching, or changing resource configurations in sequence.
- **State Manager** – Define a desired state for instances (via SSM documents applied on a schedule). For example, ensure CloudWatch Agent is always running or a certain software is installed; if not, State Manager can run a document to fix it. Think "configuration management light".
- **Parameter Store** – Hierarchical key-value store for configuration data and secrets (secure string parameters are KMS-encrypted). It's often used to pass configuration to EC2 or Lambda (you retrieve the parameter at runtime). If multiple services need a common config (like a DB connection string),

storing it in SSM Parameter Store (or Secrets Manager if it's sensitive) is better than duplicating in each app's config.

- **Event-Driven Automation** – Combining services for ops tasks:
 - **S3 Event to Lambda** – e.g., automatically run antivirus scan (Lambda) when a file is uploaded to an S3 bucket, or generate a thumbnail image. If scenario: “Whenever a file arrives, we need to process it immediately,” think **S3 event notification -> Lambda (or Step Function)**.
 - **EC2 Scheduled Events** – Suppose you want to stop non-prod instances at 6 PM daily to save cost: Use **EventBridge Scheduler (Cron)** to trigger an SSM Run Command or Lambda that stops instances on schedule.
 - **Auto remediation** – If Config rule flags non-compliance, it can trigger an SSM Automation. If CloudWatch alarm triggers (e.g., CPU high), it could trigger Lambda to clear some queue or add capacity. Essentially, use the AWS “Lego blocks” to automate responses. Common combos: CloudWatch Alarm -> EC2 Action (recover/ reboot), CloudWatch Alarm -> SNS -> Lambda (custom logic), EventBridge event (like “EC2 instance terminated”) -> Lambda (do cleanup tasks).
 - **AWS Lambda use in Ops** – As a CloudOps engineer, you might write small Lambda functions for chores, e.g., housekeeping tasks triggered by events, custom metrics publisher (maybe a Lambda that runs daily to push a custom metric), or Slack notifications. The exam could have a scenario where a simple script is needed in response to something – the answer often: **AWS Lambda** (since you can run code without servers in response to events).
 - **Example Scenario:** Patching 100 servers. Option A: SSH to each – wrong. Option B: Use SSM Patch Manager with a baseline and maintenance window – correct.
Another: Cleaning up old EBS snapshots. Could schedule a Lambda triggered by EventBridge (cron) to delete snapshots older than X days. These show understanding of tying automation to tasks.

2. Configuration Management & Drift Handling

- **AWS Config** – Track config changes to AWS resources. If a resource drifts (e.g., someone opened an SG port), Config can alert. While not explicitly mentioned in domain text, Config is a key ops service. Use **Config Rules** to enforce compliance (like “no security group open to 0.0.0.0/0 on port 22” – if violated, flag it or auto-remediate). In an automated ops context, if the question is about ensuring configurations remain as intended, AWS Config (with rules and possibly auto-remediation) is a service to mention.
- **Drift Detection in CloudFormation** – CloudFormation can detect if the real-world resource differs from the stack template (for supported resource types). If drift is detected, you can decide to import those changes into CF or revert the resource. So if they say “someone changed a setting on a resource deployed by IaC outside of IaC, how do you catch it?” – answer: **CloudFormation drift detection** or **AWS Config** (both can catch it, but CF drift works only for CF-managed resources, whereas Config covers many resource types whether or not under CF).
- **Immutable Infrastructure vs Configuration Management** – Modern approach (immutable): Rather than patch servers in place, you build new ones from a fresh image and replace (like blue/green). But both approaches can coexist (you might patch monthly via SSM, and also rebuild AMIs). The exam may not debate this, but might hint that using CloudFormation and auto healing is better

than long-lived pets that you SSH into. So lean towards answers that replace infrastructure automatically over ones that require manual fixes.

Exam Rules of Thumb – Domain 3

- **Use IaC for everything possible** → If the option is to do something manually or one-off vs using CloudFormation/CDK, choose the IaC approach. AWS likes infrastructure as code for consistency and repeatability.
 - **Cross-account deployment** → **StackSets** or **AWS Organizations** solutions (like Service Catalog, Control Tower) are preferred to keep multi-account environments in sync.
 - **Change management** → The safest deployment answers involve **blue/green or canaries** for production changes. If the question's concern is "no downtime" or "quick rollback," **Blue/Green** is the key term (like CodeDeploy Blue/Green or swapping ALB target groups). If the concern is "limited impact," **Canary/One AZ at a time** or **rolling update** might appear.
 - **Automate routine tasks** → For any routine job (backups, patching, scaling, deployments), the answer is to automate it: e.g., patch using **Systems Manager**, schedule using **EventBridge**, deploy using **CodePipeline/CodeDeploy** – not by hand or clicking.
 - **Pick the right deployment service:**
 - **CloudFormation/CDK** for infrastructure.
 - **CodeDeploy** for code on EC2/Lambda/ECS (handles the rollout logic).
 - **CodePipeline** to tie together sources, tests, deploy actions.
 - **Beanstalk** if the question suggests a dev just wants to deploy code and not manage infra (Beanstalk auto handles underlying infra using CF). But Beanstalk is less likely to be answer on CloudOps exam, unless question is about quickly deploying with minimal ops.
 - **Troubleshooting deployments** → If CloudFormation fails: check the Events, identify which resource failed. Common fixes: add missing IAM permissions, correct resource dependencies, ensure no name/DNS conflicts, increase timeout for slow operations, or handle existing resources (import or delete them). If CodeDeploy fails on EC2: check the CodeDeploy Agent is running on instances and the AppSpec file instructions – maybe the application health check failed, causing rollback (so maybe the new version was flawed or needed a longer wait time).
 - **Use Systems Manager for fleet-wide changes** → E.g., update agent on all instances = Systems Manager Run Command with a target tag for all instances. It's quick and audit-logged, versus writing a script and SSH looping. If any scenario mentions applying a change to many instances, **SSM is the answer**.
-

Domain 4: Security and Compliance (16%)

Task Statement 4.1: Implementing and Managing Security and Compliance Policies

1. Identity Management & Access Control

- **IAM Users vs Roles** – By default, avoid IAM users for application access; use **IAM Roles** attached to AWS resources (like EC2, Lambda) so they can assume privileges without static creds. Reserve IAM

users for actual people who need AWS Management Console/CLI access (and even then, preferably use SSO). **Enforce MFA** on user accounts, especially anything privileged. On the exam, if they mention a root account being used, the correct advice is: **don't use root**, enable MFA on root, and use least-privilege IAM users/roles instead.

- **IAM Policies** – Know how to read a basic policy. “Effect”: Allow/Deny, “Action”: what, “Resource”: on what. If a question gives a policy and asks what it does, parse it carefully (e.g., explicit deny or condition). **Implicit deny**: anything not explicitly allowed is denied by default. Multi-account: to allow cross-account, you often need both an IAM policy on the caller and a resource policy on the target (like S3 bucket policy allowing that account).
- **SCP (Service Control Policy)** – As mentioned, SCPs apply at Org/OU/account level to restrict actions. For example, an SCP might say “Deny ec2: on region us-east-2” to prevent using that region. *If a question scenario has multiple accounts and need to ensure a rule across all (like no usage of certain services), SCP is the tool. Remember, SCPs do not grant permissions**; they only filter out. If an action is blocked by SCP, even an Admin IAM user in the account cannot do it.
- **IAM Access Analyzer** – This tool helps find resources that are shared externally (outside the account) via resource policies. It can warn, e.g., “S3 bucket X is shared with account Y” or “an IAM role can be assumed by an external account.” It’s good for discovering unintended external access. The exam might hint at “identify any S3 buckets that are accessible by other accounts” – IAM Access Analyzer can do that analysis. It can also analyze IAM policies for broad grants.
- **IAM Policy Simulator** – Useful for testing what a policy will allow or deny by simulating calls. For troubleshooting access issues, the policy simulator can be used to pinpoint which policy or condition is blocking. If a question scenario is “Admin can’t access an S3 bucket in another account, what to do?”, steps: check bucket policy, check IAM policy, possibly use Access Analyzer or Simulator to find the missing permission.
- **Federation & SSO** – Understand that with multiple accounts it’s best to use **AWS IAM Identity Center (formerly AWS SSO)** or an external IdP to federate user access rather than creating separate users in each account. Identity Center can assign roles in target accounts to user groups centrally. Federation with IAM allows SAML or OIDC providers. The main point: **centralize identity management** for easier compliance (fewer access key usage, easier user offboarding, etc.).
- **Password Policies** – If IAM users exist, set a strong password policy (min length, require numbers, etc.) and maybe enforcement of rotation. This might appear as a simple question: how to improve IAM security for console users -> implement an IAM Password Policy with complexity and rotation.

2. Data Protection & Compliance

- **Data Classification** – AWS expects you to categorize data (sensitive vs public, etc.) and apply appropriate controls. For example: *PII data* – ensure it’s encrypted, access is logged, maybe stored in specific approved locations. *Public data* – maybe in a public S3 bucket but then ensure nothing sensitive is there. The exam might not dive deep, but could ask how to identify sensitive data in S3 -> **Amazon Macie** (uses ML to find PII like names, addresses, credit card numbers in your S3). Or how to enforce that certain data never leaves a region – maybe via SCP or backend architecture.
- **KMS (Key Management Service)** – Central service for managing encryption keys. KMS keys (CMKs) can be used by many services (S3, EBS, RDS, etc.). For KMS know: **Customer Managed Keys vs AWS Managed Keys vs AWS Owned**. Customer-managed are fully under your control (you define rotation, key policy, etc.). AWS-managed (for services like S3, EBS by default) are convenient but you have less control (they rotate automatically yearly, but you can’t manage policy). Use customer CMK

when you need control or cross-account encryption, etc. Key policies: to allow cross-account use of a CMK, you must add the other account's principal to the key policy (just giving them IAM permission isn't enough – KMS is very tightly controlled).

- **Encryption in transit** – Mainly TLS for data in transit. Use **ACM** to issue certificates for your domain and attach to ELB/CloudFront, etc. Know that ACM can provide public certs (for use on AWS endpoints) for free and auto-renew. If a scenario mentions an expired certificate causing outage, the fix might be “use AWS Certificate Manager for automated renewal” or rotate the cert. Also, **AWS CLI endpoints** can be enforced to only use HTTPS (which they do by default), and VPC endpoints (Interface Endpoints) can ensure traffic to AWS services doesn't traverse the internet at all.
- **Secrets Manager vs Parameter Store** – Both store secrets, but Secrets Manager is more robust for sensitive info (automated rotation for supported databases, immediate CloudWatch Events on secret rotation, etc.). Parameter Store is free (up to limits) and good for less sensitive or when rotation isn't needed. If question: “rotate RDS credentials automatically” -> AWS Secrets Manager with Lambda rotation function is the solution. If “store application config securely” -> could be either, but likely Parameter Store if rotation not needed, or Secrets Manager if it's passwords/keys.
- **Security Services for findings:**
 - **Amazon GuardDuty** – Enabled per account/region, uses CloudTrail, VPC Flow Logs, DNS logs to detect anomalies like crypto mining, unusual data exfiltration, IAM credential compromise patterns, etc. It generates findings (e.g., “Recon:EC2 PortProbe” or “IAM calling unusual service”). The best practice is to **enable GuardDuty in all accounts** and have it publish findings to a central Security account (GuardDuty supports a multi-account central view). If a scenario is about detecting *threats or malicious activity*, GuardDuty is key.
 - **AWS Security Hub** – Aggregates findings from GuardDuty, Inspector, Macie, etc., and checks against standards (like CIS AWS Foundations, PCI DSS). It gives a score and list of failed controls (e.g., “S3 bucket not encrypted” as a failed control). It can also coordinate automated responses (with Amazon EventBridge) to findings. In a compliance context, Security Hub is your one-stop to see if you meet certain baseline and to triage issues from multiple tools.
 - **Amazon Inspector** – Runs automated vulnerability scans. **Inspector v2** can scan EC2 instances (looking at installed packages vs CVE database), container images in ECR, and Lambda function code dependencies for vulnerabilities. It's great for identifying outdated software with known vulnerabilities or misconfigurations (like SSH open on 0.0.0.0 might not be Inspector's job – that's more GuardDuty or Config). If exam mentions “ensure AMIs are regularly scanned for vulnerabilities or container images are safe,” answer: **Amazon Inspector**.
- **AWS Config** – We mentioned in Ops context, but in security: Config rules can check encryption enabled, public access blocks on S3, unrestricted ports, etc. So it's also a security compliance tool.
- **Remediating Findings** – If Security Hub/GuardDuty/Inspector finds something, how to fix? Could be: GuardDuty says “EC2 doing crypto mining” – likely the instance is compromised; remedy: isolate it (security group lockdown), investigate, maybe terminate. Inspector says “vulnerability on Apache version” – remedy: patch the instance (SSM Patch Manager). Security Hub says “S3 bucket public” – remedy: make it private (or if it should be public, maybe use block public access accounts correctly, or use CloudFront). The exam might ask what to do when a certain alert is seen. Choose the answer that **directly addresses the issue** (e.g., for a leaked key – disable that IAM credential, audit CloudTrail, etc., rather than just “change CloudWatch alarm”).
- **Compliance Programs** – Not deeply tested, but know that AWS has compliance programs (HIPAA, PCI, etc.). If a question asks how to demonstrate compliance, likely answer: **AWS Config + AWS Audit**

Manager (Audit Manager helps collect evidence for audits). Also mention **CloudTrail** for audit logs, and **Artifact** (not likely in exam, but that's AWS's portal for compliance reports).

- **Encryption Compliance** – If an org requires “all data encrypted at rest,” you as CloudOps must ensure services have encryption toggled (S3 buckets have “encryption by default” enabled, EBS volumes encrypted, RDS with storage encryption turned on, etc.), and use Config rules to flag any unencrypted resources.

Exam Rules of Thumb – Domain 4

- **Least privilege always** → If an answer choice grants full `*.*` admin access, it's probably wrong (unless question literally asks for giving someone admin). Look for narrower scopes. Use **IAM roles for AWS resource access** (no embedding keys).
- **Cross-account access** → Use **IAM roles and resource policies**, not sharing root credentials or something silly. E.g., to let Account B's EC2 read an S3 bucket in Account A: create an IAM role in A that trusts B's account (or specific role from B) and grants S3 access, then have B's EC2 assume that role. Or add a bucket policy allowing B's IAM principal. The key: secure, auditable cross-account roles > hardcoding creds.
- **If you see “public” and it shouldn't be** → S3 bucket public, security group wide open, etc., the fix is to **lock it down** (close the port, add bucket policy to restrict or turn on block public access). AWS exams frequently test recognizing an overly permissive setup and choosing the option that tightens security while still meeting requirements (e.g., allowing only a specific CIDR instead of 0.0.0.0/0).
- **Encryption by default** → Always choose to encrypt if there's an option and no reason not to. KMS integration is ubiquitous. For KMS questions: if needing to share encrypted data with another account, use a KMS CMK with a key policy that allows the other account (and they'll need to use a grant or have the key policy permission). If an EC2 needs to read from an encrypted S3 bucket, ensure the EC2's role has decrypt permission for the KMS key used on that bucket's objects. These kind of subtle things might appear.
- **Use AWS security services** → Instead of building your own intrusion detection, use GuardDuty. Instead of manually reviewing configs, use Config rules or Security Hub. Instead of writing custom vulnerability scripts, use Inspector. The correct answers often leverage these managed services.
- **Incident response** → If a key is compromised or an instance is breached (GuardDuty alert), best practice: **rotate credentials immediately**, isolate affected resources (stop instance or at least remove from network, like changing security group), then investigate via logs. Sometimes exam answers include “contact AWS Support” – generally you don't need to contact AWS for your own security issues except in extreme cases (like reporting abuse). The better action is to **take direct control: revoke, isolate, recover from backups if needed**.
- **Compliance and logging** → Ensure **CloudTrail is on**, ideally storing logs in a locked-down bucket (and maybe validated by CloudTrail log integrity). If the question is about proving something to an auditor, CloudTrail logs, Config history, and AWS CloudTrail Lake (new queryable logs) may be in play. Also, **AWS Audit Manager** can be an answer if the question explicitly is about compiling compliance evidence automatically.

Domain 5: Networking and Content Delivery (18%)

Task Statement 5.1: Implementing and Optimizing Networking Features and Connectivity

1. VPC Networking Essentials

- **VPC and Subnets** – A Virtual Private Cloud is your network boundary. Know that subnets are tied to an AZ. **Public subnet** = has a route to an Internet Gateway (IGW). **Private subnet** = no direct IGW route (often goes to NAT for outbound or no internet at all). If instances need to be accessible from the internet, they must be in a public subnet *and* have an Elastic IP or public IP and correct SG/NACL. If not accessible, check: Does it have a public IP? Is route to IGW? SG open? NACL allowing? Common exam scenario: instance in public subnet but no public IP assigned – it won't be reachable from internet (unless through a bastion or ALB). Solution: assign Elastic IP or place behind ALB.
- **Route Tables** – Control traffic routing for each subnet (each subnet associates with one route table). Key routes: `0.0.0.0/0 -> IGW` (for internet), `0.0.0.0/0 -> NAT` (for private subnets to reach out), routes to VPNs or Direct Connect (to on-prem), etc. If resources can't reach each other, check routes: e.g., EC2 in subnet A to RDS in subnet B – need to ensure they are in same VPC or peered VPC and routes cover that.
- **Security Groups (SG)** – Stateful firewalls attached to ENIs (instances). They allow traffic *inbound* or *outbound* based on rules. Default SG outbound is allow all, inbound allow none (unless you add rules). SGs are *stateful* – if inbound is allowed, response outbound is automatically allowed and vice versa. For instance communication, remember: If Instance A's SG allows inbound from Instance B's SG (by referencing the SG), then B can talk to A (assuming B's outbound isn't restricted which it usually isn't). This is a common secure setup (referencing SGs in rules instead of IPs). Exam trick: SGs can reference other SGs even across accounts (if peered and if you use the SG ID, with some restrictions). But mostly within a VPC.
- **Network ACLs (NACL)** – Stateless, subnet-level firewalls. They have numbered rules evaluated in order, separate inbound and outbound rules. Default NACL allows all. You'd use NACLs rarely to block specific IPs or protocols for the whole subnet. Since they're stateless, you must allow both directions of traffic explicitly if using them (e.g., allow inbound port 80 and also allow ephemeral port range outbound for response). If a scenario has one subnet unable to communicate out while others can, check if a NACL is blocking (maybe an overly restrictive NACL). Best practice: leave NACLs default (open) unless you need a specific deny. They can be an added layer of security (deny rules possible on NACL, but not on SG which only allows).
- **NAT Gateway vs NAT Instance** – NAT Gateway is the managed, recommended way for private subnets to have outbound internet (for updates, etc.) without allowing inbound. It auto scales and is highly available within an AZ (and you can have multiple AZs each with a NAT GW). NAT Instances are legacy (you manage an EC2 with a NAT AMI – not recommended). If exam asks how to give instances internet access without making them public, answer: **NAT Gateway in a public subnet + route from private subnet to NAT**.
- **VPC Endpoints** – Access AWS services privately without internet:
- **Interface Endpoints** (powered by AWS PrivateLink) – Creates an ENI in your subnet that provides a private IP for a specific AWS service API (like to connect to S3, SNS, SSM, etc., via AWS backbone). You then use that private IP/endpoint DNS instead of the public endpoint. These are good for calling AWS APIs from private subnets with no IGW/NAT. They cost per hour + data.

- **Gateway Endpoints** – Specifically for **S3** and **DynamoDB** (no cost). These are not ENIs but route table entries. You plop an entry “S3 -> (gateway endpoint)” and now S3 traffic stays in AWS network. Always use these for S3/Dynamo access from a VPC to avoid NAT costs. If a question is about large volumes of S3 traffic from EC2 and reducing cost, answer: use a **VPC Gateway Endpoint for S3** (to avoid NAT data charges).
- **Hybrid Connectivity (On-Prem to AWS)** – Two main methods: **VPN** and **Direct Connect**.
- **Site-to-Site VPN** runs over the internet, typically using IPSec tunnels. It’s relatively quick to set up and cheap, but limited by internet bandwidth and higher latency (~millisecond tens). Good as a backup or for low-volume needs. AWS side is a **Virtual Private Gateway (VGW)** attached to VPC or a **Transit Gateway** can terminate VPNs too.
- **AWS Direct Connect (DX)** is a leased private line into AWS. Offers predictable performance and higher throughput (1 Gbps, 10 Gbps links etc., even aggregated). It’s more expensive and takes time to provision with a provider. Usually used when a consistent, heavy traffic or low latency connection is needed (e.g., data center extension).
- Also, know **Transit Gateway** can act as a hub for multiple VPCs *and* on-prem (VPN or DX) – making a hub-and-spoke. If a company has many VPCs that all need on-prem access, instead of making separate VPNs to each, they can do one or two VPNs into a Transit Gateway, and TGW to all VPCs.
- **AWS PrivateLink (Service Endpoints)** – Not just AWS services, you can expose your own service running in a VPC to other VPCs or accounts via PrivateLink. That creates an endpoint in consumer VPCs. Possibly out of scope for associate, but maybe know PrivateLink provides a secure way to share a service without VPC peering (e.g., AWS services use it, and you can too). If question involves connecting VPCs from different companies securely for a specific service, PrivateLink is likely (because it doesn’t require full peering or opening all traffic, just that service).
- **Network Performance** – For optimization: ensure instances in same AZ if high bandwidth is needed (traffic inside same AZ is highest bandwidth, and if in same placement group (cluster) even better). **Enable Enhanced Networking (ENA)** for >10 Gbps throughput needs. **Use Elastic Fabric Adapter (EFA)** for HPC MPI workloads (special, likely not covered deeply).
- **Cost Optimization in Networking** – Data transfer can be costly:
 - Same AZ traffic using private IPs is free. Across AZ in same region is chargeable (so architecture where one AZ calls another a lot can incur cost; might be better to keep chatter within AZ if possible, or use caching to reduce cross-AZ calls).
 - Data to internet from AWS costs egress \$. Use CloudFront to cache content at edge (cheaper and faster for users, and CloudFront has lower egress rates).
 - Use **Gateway endpoints** to avoid NAT data costs for S3/Dynamo.
 - If many VPCs need to interconnect, **Transit Gateway** might simplify and could be cheaper than many peering links (TGW has per GB cost too though).
 - **Inter-Region:** Data transfer between regions costs, unless using AWS backbone via services like VPC peering (which still costs cross-region) or CloudFront (which can fetch from a origin in another region with some cost but optimized path). There’s no free lunch; minimize cross-region chatter to reduce cost and complexity.

2. DNS and Content Delivery

- **Amazon Route 53** – AWS's DNS service. Key routing policies:
 - **Simple** (single answer, no health check by default),
 - **Weighted** (split % traffic between multiple endpoints – used for A/B testing or gradual cutover),
 - **Latency-based** (sends user to the region/AZ (with Route53 Traffic Flow for local) with lowest latency – needs health checks in place ideally),
 - **Failover** (active-passive setup, with health check on primary; if fails, Route53 returns secondary's IP),
 - **Geo-location** (route based on user's geography to specific endpoints, e.g., direct European users to EU servers always),
- **Multi-value answer** (returns multiple IPs, with health checks, sort of poor-man's load balance at DNS level; client picks one).

Understand health checks: they can monitor an endpoint (HTTP, TCP, or even check other health checks) and will stop including unhealthy endpoints in DNS responses for certain routing types like failover, latency, weighted (if set to do so). Also, Route53 can health-check other AWS resource states, like CloudWatch alarms (less common).
- **Public vs Private Hosted Zones** – Public hosted zones for internet DNS. Private hosted zones for VPC-internal DNS (e.g., you can have a private zone "corp.local" that only resolves inside your VPCs). If instances need to resolve a custom internal name, use a private hosted zone or your own DNS server. The exam might have a scenario: "instances in VPC can't resolve on-prem DNS names" – solution could be to use **Route 53 Resolver** endpoints (inbound/outbound) to bridge DNS between on-prem and AWS, or use private zones if hosting something in AWS that needs a custom domain internally.
- **CloudFront CDN** – Serves content globally from edge locations. Great for **static content (images, videos, CSS)**, and can also proxy dynamic content back to origin. It reduces latency for users and offloads origin servers. CloudFront can pull from an S3 bucket (even private bucket with an Origin Access Identity or Origin Access Control now) or from a custom origin (like an EC2 or on-prem server). It also provides SSL, OAI, caching, and can do geo-restriction or WAF integration. Typical exam use: if users in various geos experience slowness retrieving files from S3 in one region, answer: put CloudFront in front of S3 to cache worldwide. Also, CloudFront + [S3 static website] is a common pattern for a globally distributed website.
- **AWS Global Accelerator** – Often confused with CloudFront. GA is a networking service that provides two anycast static IPs to your app and routes user traffic through AWS's global network to the endpoint closest to your application's region. It's good for non-HTTP use cases or when you need to keep using TCP/UDP at layer 4 (like gaming servers, VoIP, or even HTTP APIs where you want improved performance but can't use a CDN due to dynamic, non-cacheable content). GA doesn't cache; it's about optimizing network path. If a question says "users globally need low-latency access to a *non-HTTP application (say, a gaming server on UDP)* in region X," Global Accelerator is the answer. If it's HTTP and cacheable, CloudFront is the go-to. Also, GA can do smart routing to nearest healthy regional endpoint in an active-active multi-region setup (like two regional ALBs, GA can health-check and route to nearest healthy one).

- **S3 Transfer Acceleration** – Mentioned earlier, but relevant here too: it uses CloudFront’s network to speed uploads to S3. If users are far from the S3 bucket’s region, enabling transfer acceleration on the bucket gives them a url like `bucket.s3-accelerate.amazonaws.com` which will ingest to the nearest CloudFront edge and then hop to S3 region. It’s ideal for improving upload speeds from far locations (and downloads in some cases, but primarily marketed for uploads). If question: remote office complains of slow S3 uploads to us-east-1, solution: enable S3 Transfer Accel.
- **Route 53 Resolver** – A newer feature for hybrid DNS. You can set up **Inbound endpoints** so on-prem DNS can query Route53 private zones (for stuff in VPC) and **Outbound endpoints** with rules so that your VPC DNS can resolve on-prem domain names by forwarding to your on-prem DNS server. For example, if EC2 at AWS needs to resolve “db.corp.local” which is an on-prem domain, you’d configure an Outbound Resolver Rule for `corp.local` to forward to your on-prem DNS IPs over the VPN or DX. This is key for seamless hybrid DNS.

3. Network Security & Troubleshooting

- **VPC Flow Logs** – Can be enabled on VPC, subnet, or ENI level to log accepted/rejected traffic (source/dest, ports, etc.). Great for troubleshooting connectivity or monitoring for anomalous traffic. If an instance isn’t communicating with another and you can’t figure why, flow logs might show “REJECT” indicating security groups or NACLs are blocking. Or if you see no flow logs at all for traffic you expect, maybe it’s not hitting the instance at all (bad route or wrong IP). They can send logs to CloudWatch or S3. In exam, a question on “how to investigate why traffic isn’t flowing or to verify network traffic patterns” suggests enabling **VPC Flow Logs**.
- **ELB Access Logs** – For ALB/NLB/CLB, can be turned on to record requests/connections, which helps in troubleshooting or analysis of requests (source IPs, latency, etc.). If an app behind ALB is getting intermittent failures, ALB logs might show patterns or specific error codes from targets. CloudFront also has access logs (and even real-time logs). Use these logs to pinpoint issues like “user’s requests never reaching the server” (maybe a routing policy issue or ALB misconfig).
- **Reachability Analyzer** – An AWS tool that analyzes network path between two resources (say EC2 in subnet A and EC2 in subnet B or EC2 and an S3 endpoint) and tells you if it’s reachable or where it’s blocked (SG, NACL, route, etc.). It’s extremely handy and likely to be an answer if the question is “How to *quickly* identify what’s blocking connectivity between resource X and Y” – the **Reachability Analyzer** is purpose-built for that. It can save you from manually checking every SG/NACL.
- **AWS Network Firewall** – A managed firewall (layer 4 and 7 rules) that you can put in your VPC (especially in a central VPC) to inspect traffic egress/ingress. You can filter by IP, port, or even application layer patterns (it uses Suricata rules engine). Usually used in an **inspection VPC scenario** (Hub-Spoke with TGW, all traffic goes through the firewall). If a company needs to enforce more advanced rules than SGs (like block URLs, or IDS/IPS), Network Firewall is the AWS-native way. Alternatives: partner appliances from marketplace. If exam scenario is about needing to block specific outgoing URL patterns or do Deep Packet Inspection, **AWS Network Firewall** is likely answer (or possibly **WAF** if it’s web-specific injection patterns at HTTP level).
- **AWS WAF and Shield** – WAF (Web Application Firewall) attaches to CloudFront, ALB, or API Gateway to block/allow HTTP(s) requests based on rules (like block SQL injection, block by country, rate-limit traffic, etc.). If scenario: “mitigate web attacks (OWASP Top 10) on an application”, answer: **use AWS WAF with appropriate rules**. Shield Standard is always on (protects against common DDoS at network layer free by AWS). **Shield Advanced** (paid) gives enhanced DDoS detection/mitigation and access to the SRT (support), plus some WAF integration and insurance against DDoS costs. If a

question mentions DDoS specifically and impact, Shield Advanced might be relevant (for important production facing public endpoints where DDoS risk is high). Otherwise, mention CloudFront and WAF because CloudFront can absorb some volumetric attacks and WAF filters bad requests.

- **Troubleshooting tips:**

- If an instance *can't connect out* (to the internet): check **NAT Gateway** (exists? route?), **DNS resolution** (instances use VPC DNS, so if custom DNS or DHCP option changed, maybe it can't resolve domain names), **SG outbound** (usually open by default, but could be custom blocking), and **NACL**. Also check that instance has internet access – if in public subnet, needs Elastic IP or Public IP. If in private, need NAT for internet.
- If *can't connect in* to an instance: check **SG inbound** (open for that port from source), **NACL inbound/outbound**, **Route table** (public subnet route to IGW if expecting internet traffic), **Windows instance** (check Windows firewall on the OS in addition to SG!).
- If VPC peering connection issues: ensure **route tables** have entries for the peer VPC's CIDR on both sides, ensure SGs allow the traffic (SGs don't automatically know peered traffic vs same VPC, but effectively treat it like within VPC), note that peering doesn't support transitive – if trying to go through one VPC to another third, won't work.
- If overlapping IP ranges between VPCs: peering or TGW won't route those (overlap is not allowed in peering, and in TGW attachments it will likely cause undefined routing). The solution is to renumber one VPC (not trivial) or avoid connecting them. Exam might hint at "can't create peering because CIDR conflict" -> answer: must change VPC CIDR via migration or use a workaround such as NAT instance proxies (rare).
- **Direct Connect troubleshooting:** if DX is up but traffic not flowing, check **advertised routes** (BGP), check VPC route table has the DX routes via VGW, check security (NACL/SG) allowing on both ends, and ensure no ACL on on-prem side. Likely not asked but just in case.

Exam Rules of Thumb – Domain 5

- **Public vs Private** → To be public, resource needs: in public subnet + public IP + route via IGW + SG allows ingress. If any piece missing, it's not reachable from internet. For private instances needing internet, think NAT or endpoints.
- **Connectivity options** → Connect VPCs with **VPC Peering** (if small scale, no transitive, low latency, no additional cost for data in region). Many VPCs or need transitive hub? **Transit Gateway**. Connect VPC to on-prem? **VPN for quick/cheap**, **Direct Connect for stable/high throughput**. Need VPC-to-service? **VPC Endpoints (PrivateLink)**.
- **Use CloudFront for content** → Any static or streaming content distribution scenario, CloudFront is best. Offloads S3 or web servers and improves global latency.
- **Global Accelerator vs CloudFront** → If question is about **improving latency for dynamic, non-cacheable applications** (or protocols like UDP), **Global Accelerator** is likely answer. For standard web content and static files, **CloudFront**.
- **Network cost** → Keep traffic local (same AZ if possible for chatty apps, or at least same region – cross-region should be avoided unless necessary). Use **endpoints** to avoid NAT when accessing AWS services. Use **CloudFront** to reduce internet egress (as edge caches are cheaper for high volume). Consider **compression** or efficient data transfer (like AWS DataSync for moving large data efficiently with checksums, though that's more Data transfer service domain).
- **Diagnose connectivity systematically:**

- Ping/telnet from source (if allowed) to see reachability.
- Use **VPC Reachability Analyzer** for a quick assessment of where traffic might be blocked between two AWS resources.
- Turn on **VPC Flow Logs** on the subnet or ENI in question to see if traffic is hitting and being rejected by NACL/SG.
- Double-check NACL rules: remember order and that a lower rule number with a deny can override later allows.
- Check if any AWS Firewall or Security service is in line (Network Firewall, etc., might be dropping traffic if configured).
- **DNS issues** → If instances can't resolve DNS: ensure the VPC DHCP option set is using a valid DNS (usually AWSProvidedDNS by default). If using custom DNS server on an instance, that instance must be reachable. Commonly, people forget to update DNS in VPC peering (each VPC has its own DNS resolution scope – you can't resolve private records of one VPC from another unless you use Route 53 Resolver endpoints or share DNS). So for cross-VPC name resolution, maybe use Route 53 Private Hosted Zones associated with both VPCs or a centralized DNS server accessible to both.
- **WAF vs Security Group** → SG is network layer allowlist for specific ports/ips, WAF is for HTTP(s) layer 7 filtering (SQLi, XSS, etc.) on CloudFront/ALB/API GW. If threat is an app-layer attack (like someone trying SQL injection on your website), **AWS WAF** is the mitigation. If threat is general port scanning or DDoS, **Security Groups + Shield** cover a lot (SG blocks unsolicited ports by default, Shield standard protects common L3/L4 attacks).
- **Encryption in transit** → For public endpoints, always use HTTPS (ALB, CloudFront with ACM certs). For internal traffic, consider if needed – e.g., encrypt data between services if it goes over untrusted networks. For within VPC, often plaintext is okay (already isolated), but some compliance might require even internal encryption – in which case, set up private PKI or use ACM for internal certs on custom domains in private DNS.

This concludes the comprehensive study notes for AWS CloudOps Engineer – Associate (SOA-C03). Focus on the above key points, practice with scenario questions, and you'll be well-prepared to make the right decisions quickly in the exam. Good luck with your certification!
