

Section 5: EC2 Fundamentals

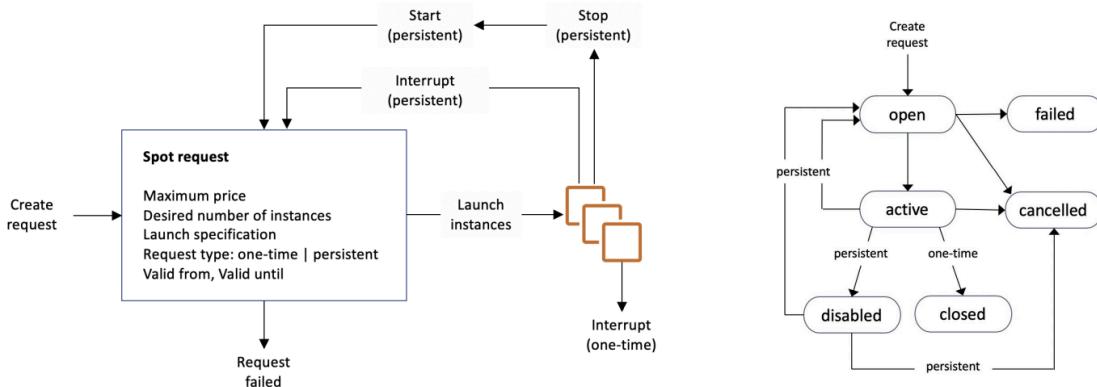
EC2 Spot Instance Requests

- Define max spot price to get instance while current spot price < max
 - Hourly spot price varies
 - If current spot price > max price, stop or terminate instance with 2 minute grace period
- Spot Block
 - “Block” spot instance during a specified time frame (1-6 hours) without interruptions
 - Rare instances where spot instance is reclaimed
- Used for batch jobs, data analysis, or workloads that are resilient to failures, not critical workloads or DB

How to terminate Spot Instances?

- One time or persistent spot requests. If it is one time, the spot instance will start and the spot request is gone. If persistent, the spot request will remain and if an instance is terminated, another will be launched.
- Cancel spot requests in open, active, or disabled state
- Canceling stop request does not terminate instances, must cancel Spot request then stop instances

How to terminate Spot Instances?



You can only cancel Spot Instance requests that are **open, active, or disabled**.
Cancelling a Spot Request does not terminate instances
You must first cancel a Spot Request, and then terminate the associated Spot Instances

Spot Fleets

- Set of spot instances + on demand instances
- Automatically request spot instances with lowest price
- Fleet will try to meet target capacity with price constraints
 - Define possible launch pools: instance type, OS, AZ
 - Can have multiple launch pools for fleet to choose
 - Spot fleet stops launching instances when reaching capacity or max cost
- Strategies to allocate Spot Instances:
 - lowestPrice: from pool with lowest price (cost optimization, short workload)
 - Diversified: distributed across all pools (great for availability, log workloads)
 - capacityOptimized: pool with optimal capacity for number of instances
 - priceCapacityOptimized (recommended): pools with highest capacity available, then select pool with lowest price (best choice for most workloads)

Section 6: EC2 - SAA Level

Private vs Public IP (IPv4)

- Public IP:
 - Machine can be identified on internet
 - Must be unique across whole web, can be geo-located easily
- Private IP:
 - Machine can only be identified on private network only
 - IP must be unique across private network, but 2 different private networks can have same IPs
 - Machines connect to internet via internet gateway as proxy
 - Only specified range of IPs can be used as private IP

Elastic IP

- When you stop and start EC2 instance, it can change its public IP
- If you need to have a fixed public IP for instance, you need Elastic IP
- Elastic IP is public IPv4 IP you own as long as you don't delete it
- Attached to 1 instance at a time
- With Elastic IP, you can mask the failure of an instance by rapidly remapping the address to another instance in your account
- Only 5 Elastic IP in account (can be increased)
- Avoid using, instead use random public IP and register DNS name
 - LB and no public IP
- Default EC2:
 - Private IP for AWS, public IP for internet
- SSH into EC2:
 - Can't use private IP because not on same network, only public IP

- If machine is stopped and started, public IP can change

Placement Groups

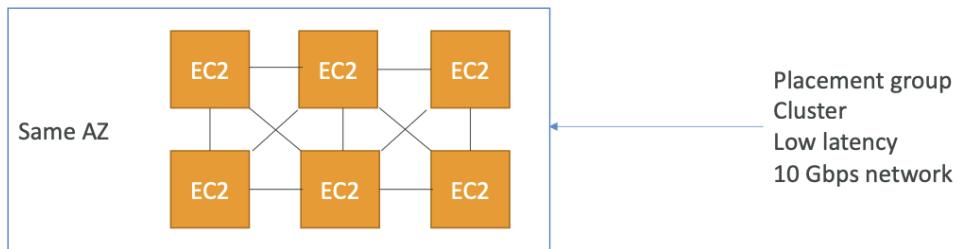
- Control over EC2 instance placement
- Defined using placement groups
- Specify strategy:
 - Cluster: cluster instances in low latency group in single AZ
 - Spread: spreads instances across underlying hardware (7 instances per group per AZ) for critical apps
 - Partition: spreads instances across many different partitions within AZ; scales to 100s of instances per group

Cluster

- Single AZ, great network with enhanced networking enabled
- Con: if AZ fails, all instances fail at same time
- Use cases: big data to complete very fast, app needs low latency and high network throughput

Placement Groups

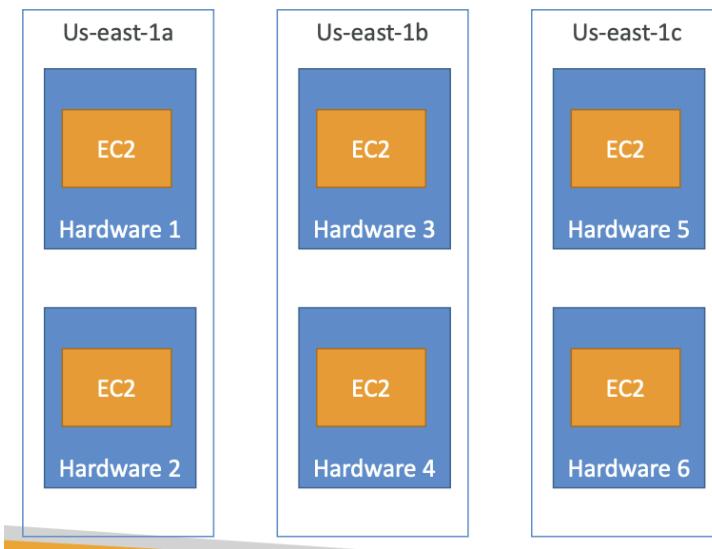
Cluster



- Pros: Great network (10 Gbps bandwidth between instances with Enhanced Networking enabled - recommended)
- Cons: If the AZ fails, all instances fail at the same time
- Use case:
 - Big Data job that needs to complete fast
 - Application that needs extremely low latency and high network throughput

Spread

Placement Groups Spread



- Pros:
 - Can span across Availability Zones (AZ)
 - Reduced risk of simultaneous failure
 - EC2 Instances are on different physical hardware
- Cons:
 - Limited to 7 instances per AZ per placement group
- Use case:
 - Application that needs to maximize high availability
 - Critical Applications where each instance must be isolated from failure from each other

- Span multiple AZ, EC2 instances on different hardware for reduced risk of failure
- Limits 7 instances per AZ per placement group
- Use case: High availability, critical apps for isolated hardware

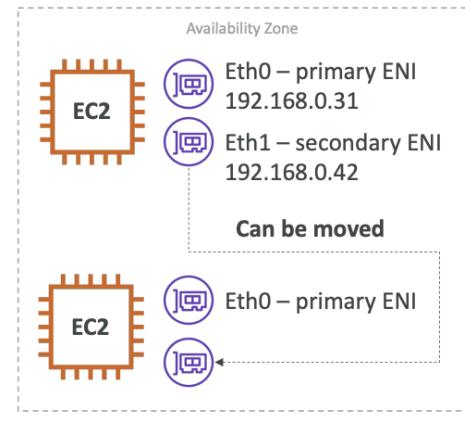
Partition

- Up to 7 partitions per AZ, spanning across multiple AZ in same region - up to 100s of instances
- Instances in partition do not share hardware with other partitions
 - Partition failure does not affect other partitions
- Metadata used by EC2 instances to get partition information
- Use case: Kafka

Elastic Network Interface (ENI)

Elastic Network Interfaces (ENI)

- Logical component in a VPC that represents a virtual network card
- The ENI can have the following attributes:
 - Primary private IPv4, one or more secondary IPv4
 - One Elastic IP (IPv4) per private IPv4
 - One Public IPv4
 - One or more security groups
 - A MAC address
- You can create ENI independently and attach them on the fly (move them) on EC2 instances for failover
- Bound to a specific availability zone (AZ)
 - Component in VPC that represents virtual network card
 - Attributes:
 - Primary private IPv4, 1+ secondary IPv4
 - 1 elastic IPv4 per IPv4
 - 1 public IPv4
 - 1+ SG
 - MAC address
 - Create ENI independently and attach on fly on EC2 instances for failover
 - Bound to AZ

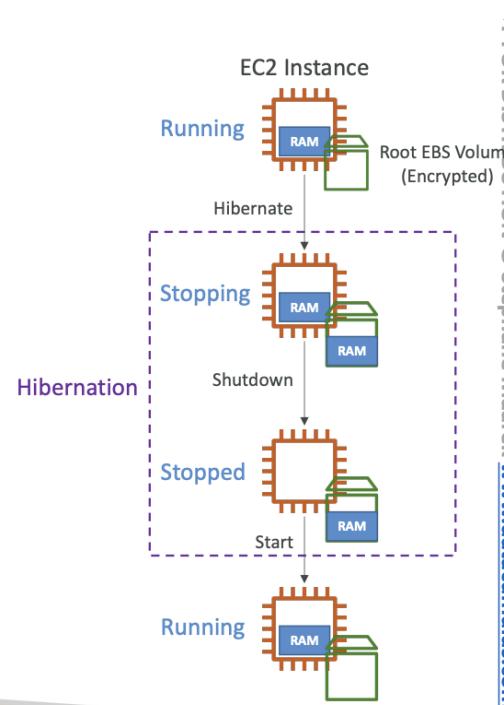


EC2 Hibernate

EC2 Hibernate

- Introducing EC2 Hibernate:
 - The in-memory (RAM) state is preserved
 - The instance boot is much faster! (the OS is not stopped / restarted)
 - Under the hood: the RAM state is written to a file in the root EBS volume
 - The root EBS volume must be encrypted
- Use cases:
 - Long-running processing
 - Saving the RAM state
 - Services that take time to initialize

Stephane Maarek



- Stop: data on disk is kept intact in next start
- Terminate: any EBS volumes (root) also set up to be destroyed is lost
- On start:
 - First start: OS boots and user data script is run
 - App starts, cache warmed
- Hibernate:
 - RAM state is preserved
 - Instance boot much faster (OS not stopped)
 - Under the hood: RAM state written to file at root EBS volume
 - Root EBS volume must be encrypted
- Use case: long running processes, save RAM state, services that take time to initialize

Good to know:

- Many instance families
- RAM size must be < 150 GB
- Instance size - not supported for bare metal instances
- Many AMI's
- Root volume - EBS only, encrypted
- Available for on demand, reserved and spot instances
- Not hibernated for more than 60 days

Section 7: EC2 Instance Storage

EBS Encryption

- When encrypting EBS:
 - Data at rest encryption
 - Data in flight between instance and volume encrypted
 - All snapshots encrypted
 - All volumes created from snapshot
- Encryption and decryption are handled transparently (do nothing)
 - Minimal latency
- KMS keys (AES 256)
- Copying unencrypted snapshot allows encryption
- Snapshots of encrypted volumes are encrypted

Encrypt unencrypted EBS volume

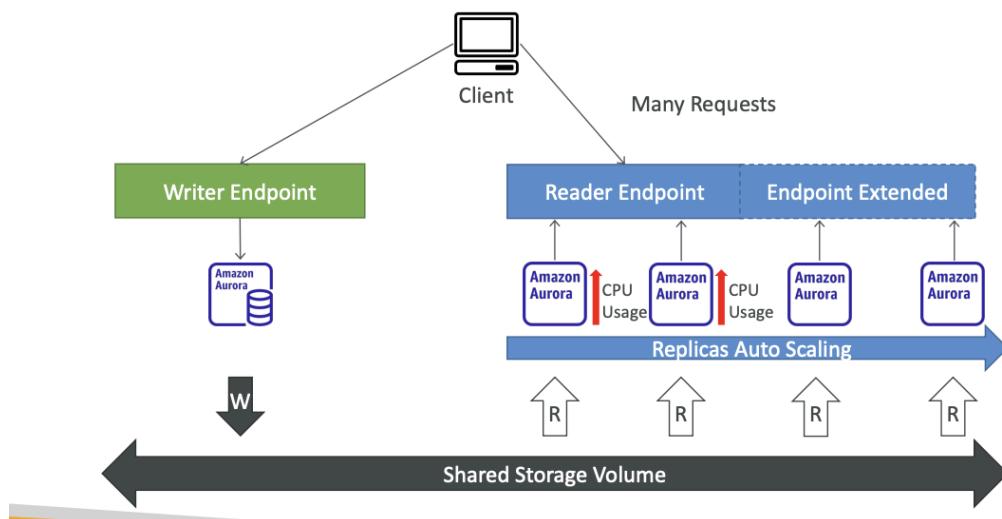
- Create EBS snapshot of volume
- Encrypt EBS snapshot (using copy)
- Create new EBS volume from snapshot (volume will also be encrypted)
- Attach encrypted volume to original instance

Section 9: AWS Fundamentals: RDS + Aurora + ElastiCache

Aurora Advanced Concepts

Replica Auto Scaling

Aurora Replicas - Auto Scaling

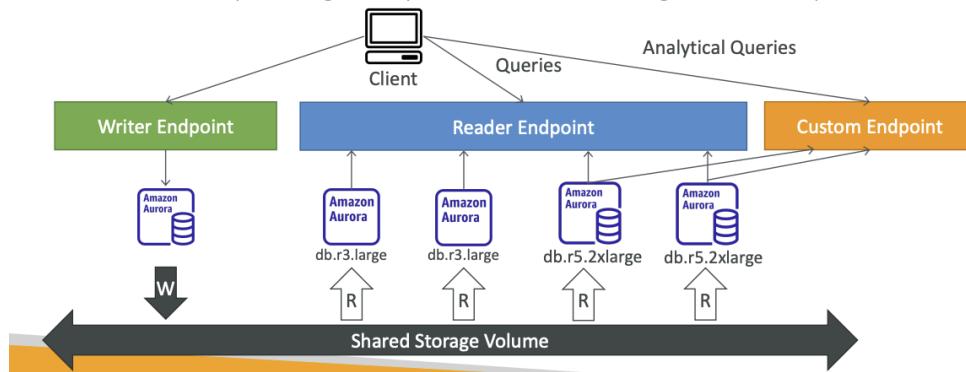


- Add replicas and have extended reader endpoints to cover the new replicas to distribute reads

Custom Endpoints

Aurora – Custom Endpoints

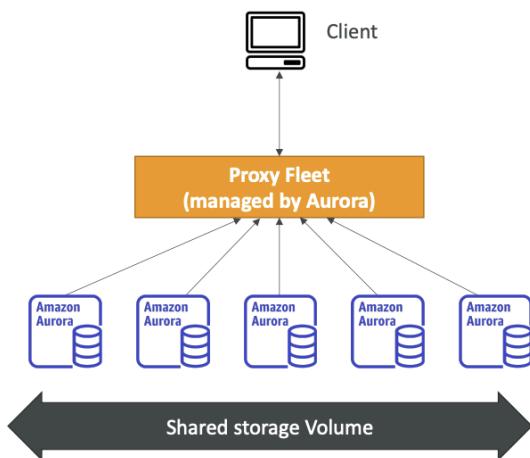
- Define a subset of Aurora Instances as a Custom Endpoint
- Example: Run analytical queries on specific replicas
- The Reader Endpoint is generally not used after defining Custom Endpoints



- With different DB sizes, a subset of reader instances are used as custom endpoints
- Reader endpoint generally not used after custom endpoints

Aurora Serverless

- Automated database instantiation and auto-scaling based on actual usage
- Good for infrequent, intermittent or unpredictable workloads
- No capacity planning needed
- Pay per second, can be more cost-effective



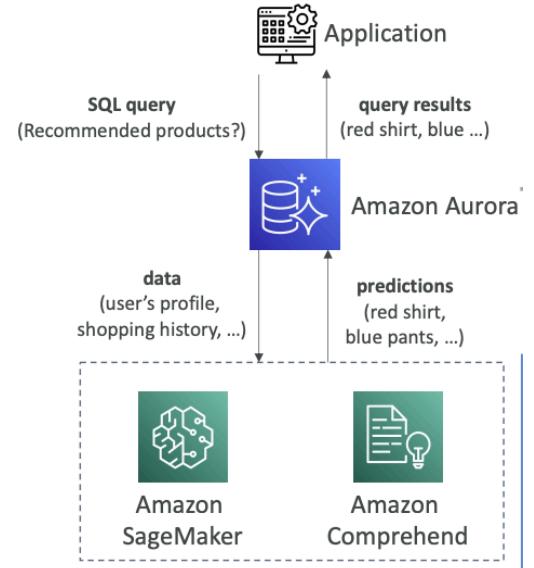
- Automated DB instantiation and auto scaling based on usage
 - Good for infrequent or intermittent / unpredictable workloads
 - No capacity planning, pay per second

Global Aurora

- Cross region read replicas
 - Disaster recovery, simple to put in place
- Global DB (recommended)
 - 1 primary region for read / write
 - Up to 5 secondary (read only) regions, replication lag is < 1 second
 - Up to 16 read replicas per secondary region
 - Helps for decreasing latency
 - Promoting another region < 1 min
- Typical cross region replication takes < 1 second

Aurora Machine Learning

- ML based predictions to apps via SQL
- Simple, optimized integration between Aurora and ML services
- Supports:
 - Sagemaker
 - Comprehend
- Use cases: fraud detection, sentiment analysis, etc...



RDS & Aurora – Backup and Monitoring

RDS Backups

- Automated backups
 - Daily full backup of DB (during backup window)
 - Transaction logs backed up by RDS every 5 min
 - Ability to restore to any point in time
 - 1 to 35 day retention, 0 to disable backup
- Manual DB snapshots
 - Retain backup as long as you want
- Trick: in stopped RDS DB, if you plan on stopping it for a long time, snapshot and delete original, then restore

Aurora Backups

- Automated backups
 - 1 to 35 day retention (cannot be disabled)
 - Point in time recovery in that timeframe
- Manual backups:
 - Manual trigger, retains as long as you want

RDS & Aurora Restore Options

- Restoring a backup or snapshot creates a new DB
- Restore MySQL RDS DB from S3
 - Create a backup from on premise DB
 - Store in S3
 - Restore backup file from S3 onto new RDS instance running MySQL
- Restoring MySQL Aurora cluster from S3
 - Create backup of on premise DB using Percona XtraBackup
 - Store backup file on S3
 - Restore backup file on new Aurora cluster running MySQL

Aurora DB Cloning

- Create new Aurora DB cluster on existing one
- Faster than backup & restore
- Uses copy on write protocol
 - New DB cluster uses same data volume as original DB cluster (fast and efficient - no copying needed)
 - When updates are made to new DB cluster data, additional storage is allocated and data is copied to be separated
- Useful to create a “staging” DB from “prod” DB without impacting the production DB

RDS & Aurora Security

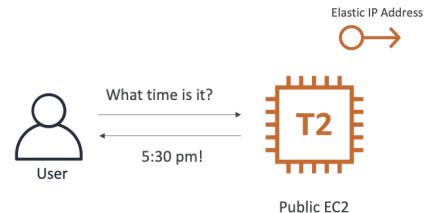
- At rest encryption:
 - DB master & replicas encrypted via KMS defined at launch time
 - If master is not encrypted, read replicas cannot be encrypted
 - To encrypt unencrypted DB, go through DB snapshot & restore as encrypted
- In flight encryption:
 - TLS ready by default, use AWS TLS root certificates client side
- IAM Auth: IAM roles to connect to DB (instead of username and password)
- SG: control network access to RDS / aurora DB
- No SSH available except on RDS custom

- Audit logs can be enabled and sent to CW logs for longer retention

Section 11: Classic Solutions Architecture Discussions

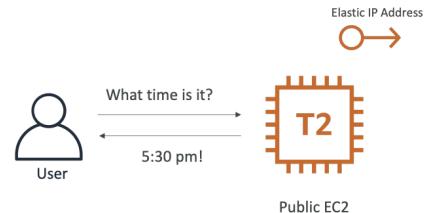
Stateless web app: WhatsTheTime.com

- Starting small with a T2 micro EC2 instance will host website and has an elastic IP to have static IP if any restarts need to happen



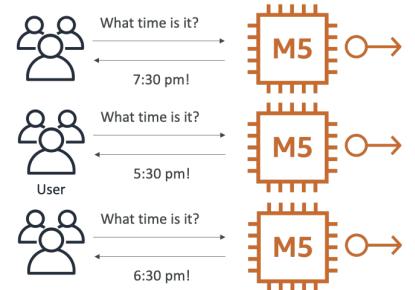
- More uses, scale to M5 instance, but there is downtime while M5 instance is being deployed. Add more M5 instances for more users with more elastic IPs

Stateless web app: What time is it? Starting simple

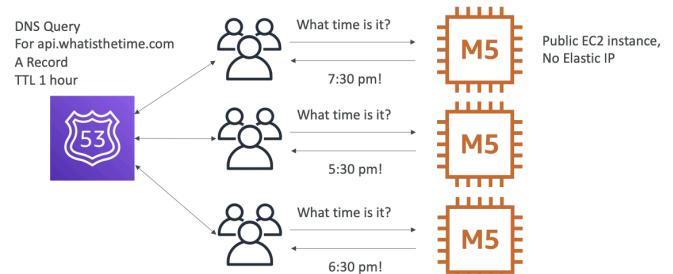


- Users leverage Route 53 with A record of TTL 1 hour for no elastic IPs. Route 53 will keep instances in sync. However if an instance is removed, TTL might be saved to a removed instance.

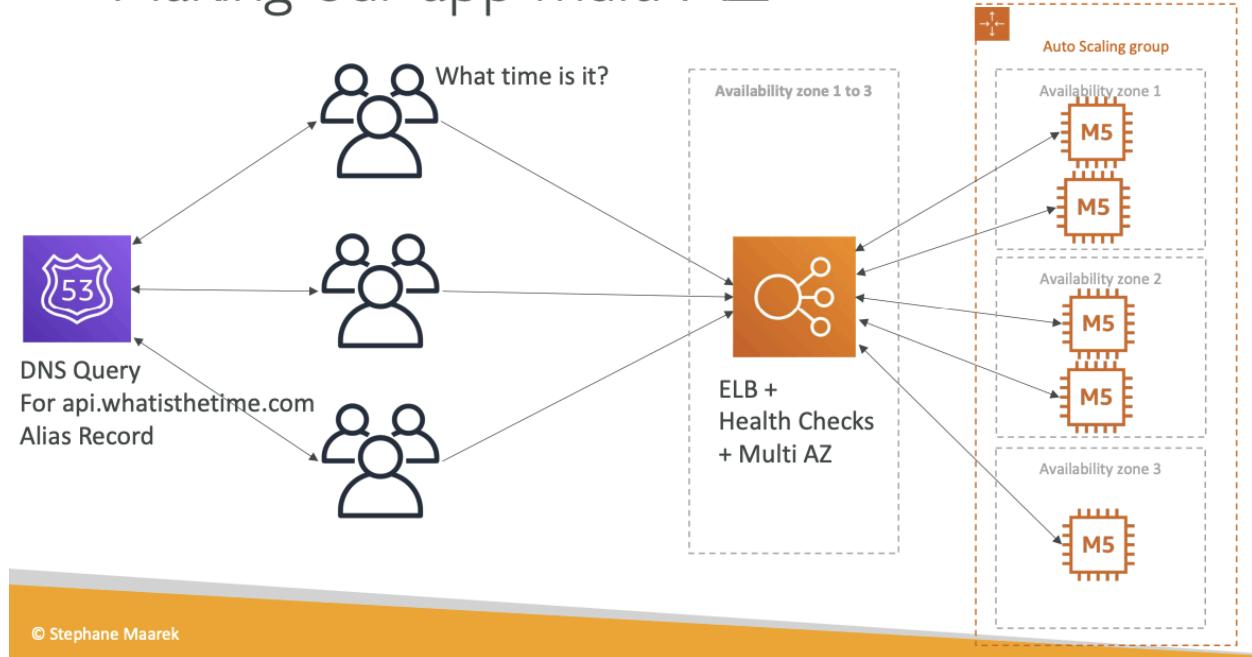
Stateless web app: What time is it? Scaling horizontally



Stateless web app: What time is it? Scaling horizontally



Stateless web app: What time is it? Making our app multi-AZ

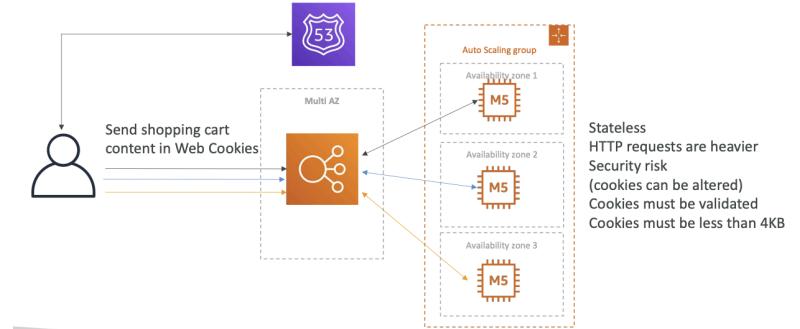


- Auto scaling instances in multiple AZs for disaster recovery. ELB as public facing with SG leading to instances. Alias record to track ELB for users. Additional cost savings can be made by reserving at least 1 instance per AZ.

Stateful Web App

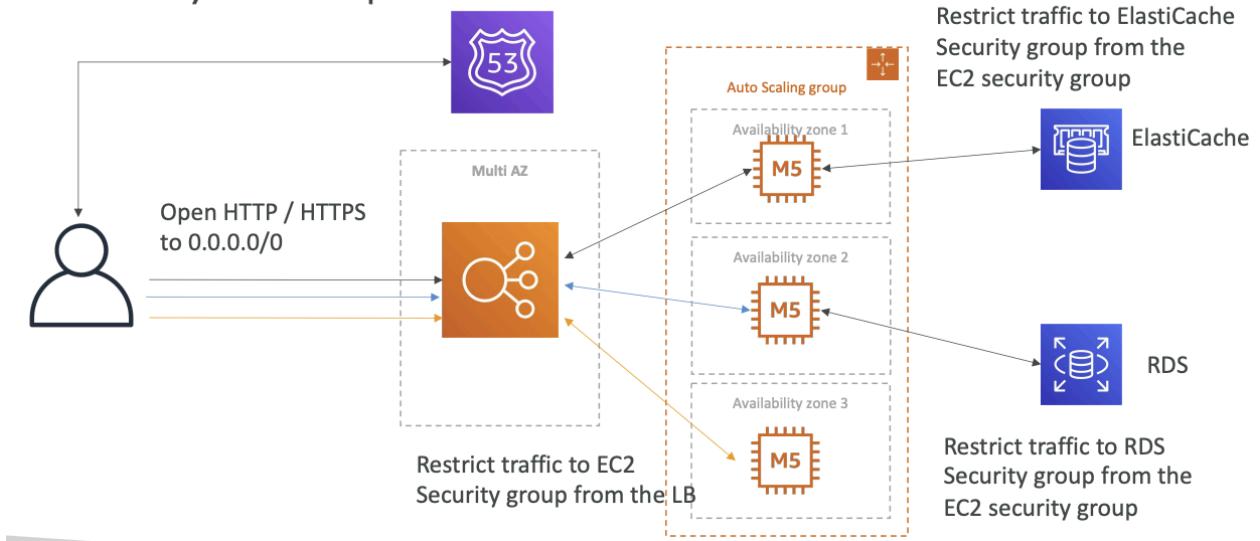
- The ending architecture from the other example poses a risk where if the user moves to another instance, session state is lost. Using stickiness fixes the issue, but if instance terminated, session is lost.
- Using user cookies solves the issue, but HTTP requests are heavier and must be validated if cookies are altered.

Stateful Web App: MyClothes.com Introduce User Cookies



Stateful Web App: MyClothes.com

Security Groups

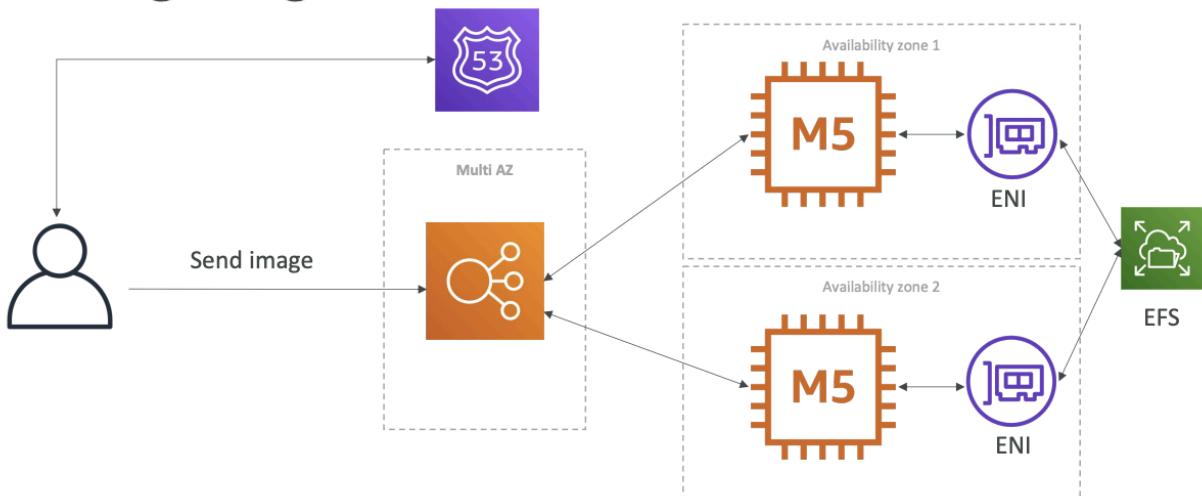


- Using a session ID in web cookies to store and retrieve data from the cache to save session state. RDS can be used with read replicas to have disaster recovery. Add SG to restrict traffic to EC2 instances and ElastiCache + RDS
 - 3 tier architecture

Stateful Web App

Stateful Web App: MyWordPress.com

Storing images with EFS



- Building off the previous example, you can exchange RDS for Aurora for a serverless architecture. To store images, EBS can be used as a block store. However, with increased instances you need additional EBS volumes on each AZ. EFS can be used to solve this issue by creating ENI's in each AZ for AZ's to access single EFS.

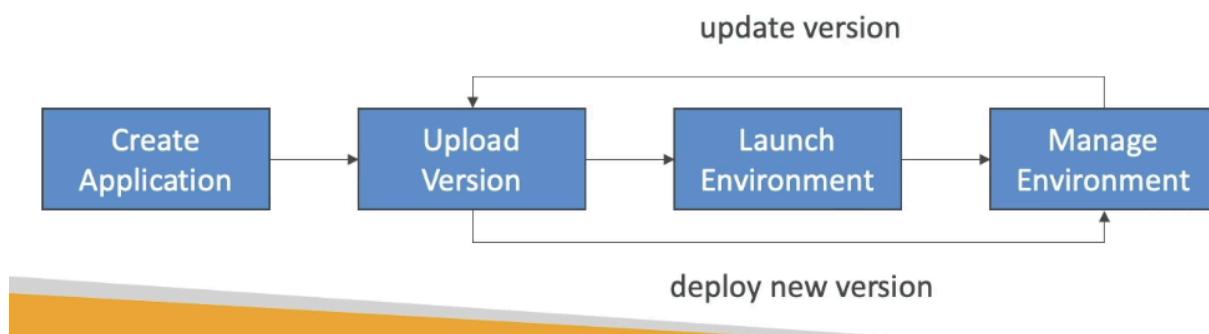
Instantiating Applications Quickly

- EC2 Instances
 - Golden AMI: install applications, OS dependencies... beforehand and launch EC2 instance from golden AMI
 - Bootstrap via user data: for dynamic configuration
 - Hybrid: mix of golden AMI and user data (Elastic Beanstalk)
- RDS
 - Restore from snapshot: DB will have schemas and data ready
- EBS Volumes
 - Restore from snapshot, disk will be read and formatted

Beanstalk Overview

- Developer centric view of deploying apps on AWS
 - Managed service
 - Automatically handles capacity provisioning, LB, scaling, application health monitoring, instance configuration
 - Developer in charge of application code
 - Full control over configuration, paid for underlying services

Beanstalk Components

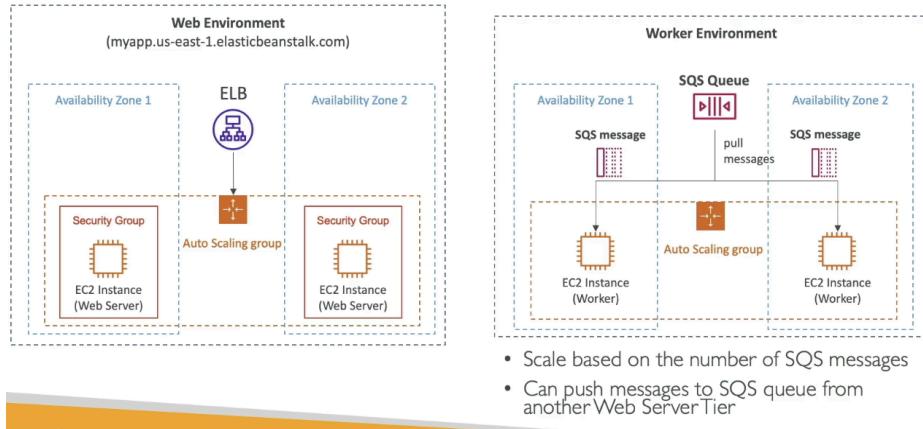


- Application: collection of Elastic Beanstalk components (environments, versions, configurations)
- Application version, environment (tiers)
 - Can create multiple environments and has worker vs web server tier

Supports many platforms

Web Server Tier vs Worker Tier

Web Server Tier vs. Worker Tier



Deployment Modes

Elastic Beanstalk Deployment Modes

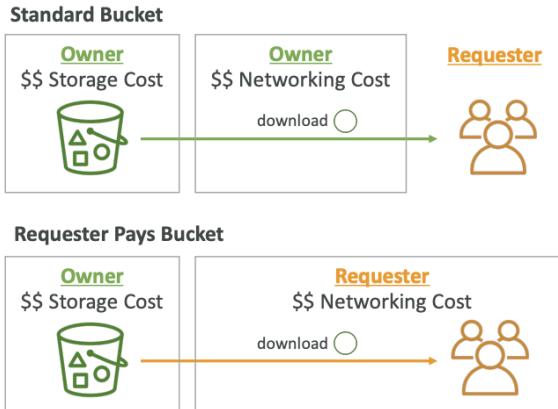


- Single Instance (for dev) or high availability with LB (prod)

Section 13: Advanced S3

S3 Requester Pays

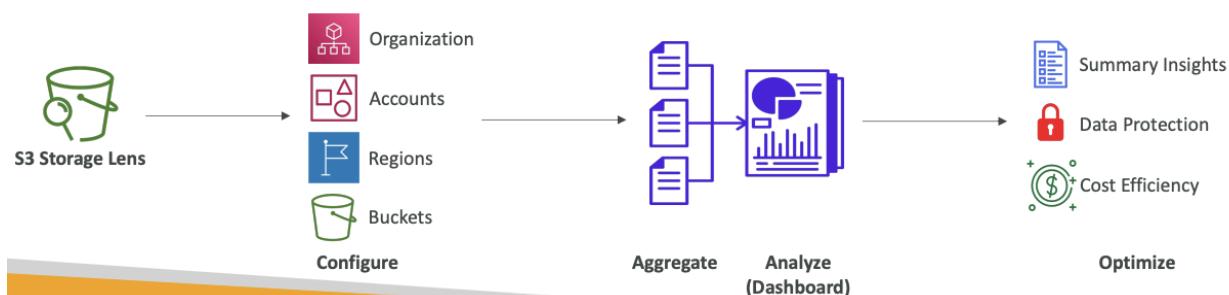
- In general, bucket owners pay for all S3 storage and data transfer costs with their bucket
- With requester pays buckets, requester pays the costs of the request and data download from bucket
 - Helpful when you want to share large datasets with other accounts
 - Requester must be authenticated with AWS



S3 Batch Operations

- Perform bulk operations on existing AS3 objects with a single request
 - Modify object metadata & properties
 - Copy objects between S3 objects
 - Encrypt un-encrypted objects
 - Modify ACLs, tags
 - Restore objects from S3 glacier
 - Invoke lambda function to perform custom action on each object
- Job consists of a list of objects, action to perform, and optional parameters
- Manages retries, tracks progress, sends completion notifications, generate reports
 - S3 inventory to get object list and S3 Select to filter objects

S3 Storage Lens



- Understand, analyze, and optimize storage across entire AWS organization
 - Discover anomalies, identify cost efficiencies, and apply data protection best practices across entire AWS organization (30 days usage & activity metrics)
 - Aggregate data for Organization, specific accounts, regions, buckets, or prefixes

- Default dashboard or create your own
- Can be configured to export metrics daily to S3 bucket (export in CSV, parquet)

Default Dashboard

- Visualize summarized insights and trends for both free and advanced metrics
- Shows multi-region and multi-account data
- Preconfigured by S3
- Can't be deleted, only disabled

Metrics

Storage Lens – Metrics



- Summary Metrics
 - General insights about your S3 storage
 - StorageBytes, ObjectCount...
 - Use cases: identify the fastest-growing (or not used) buckets and prefixes
- Cost-Optimization Metrics
 - Provide insights to manage and optimize your storage costs
 - NonCurrentVersionStorageBytes, IncompleteMultipartUploadStorageBytes...
 - Use cases: identify buckets with incomplete multipart uploaded older than 7 days, Identify which objects could be transitioned to lower-cost storage class
- Summary Metrics
 - General insights of S3 storage
 - StorageBytes, ObjectCount...
 - Use cases: identify fastest growing (or not used) buckets and prefixes
- Cost Optimization Metrics
 - Provide insights to manage and optimize storage costs
 - NonCurrentVersionStorageBytes,
 - IncompleteMultipartUploadStorageBytes...
 - Use cases: identify buckets with incomplete multipart uploaded older than 7 days, identify which object can be transitioned to lower cost storage class
- Data protection Metrics
 - Provide insights to data protection features

- VersioningEnabledBucketCount, MFADeleteEnabledBucketCount, SSEKMSEnabledBucketCount, CrossRegionReplicationRuleCount...
- Use cases: identify buckets that aren't following data protection best practices

Storage Lens – Metrics



- Data-Protection Metrics
 - Provide insights for data protection features
 - VersioningEnabledBucketCount, MFADeleteEnabledBucketCount, SSEKMSEnabledBucketCount, CrossRegionReplicationRuleCount...
 - Use cases: identify buckets that aren't following data-protection best practices
- Access-management Metrics
 - Provide insights for S3 Object Ownership
 - ObjectOwnershipBucketOwnerEnforcedBucketCount...
 - Use cases: identify which Object Ownership settings your buckets use
- Event Metrics
 - Provide insights for S3 Event Notifications
 - EventNotificationEnabledBucketCount (identify which buckets have S3 Event Notifications configured)
- Access Management metrics
 - Provide insights for S3 Object Ownership
 - ObjectOwnershipBucketOwnerEnforcedBucketCount...
 - Use case: identify which Object Ownership settings buckets use
- Event Metrics
 - Provide insights for S3 Event Notifications
 - EventNotificationEnabledBucketCount (identify which buckets have S3 Event Notifications configured)



Storage Lens – Metrics

- Performance Metrics

- Provide insights for S3 Transfer Acceleration
- TransferAccelerationEnabledBucketCount (identify which buckets have S3 Transfer Acceleration enabled)

- Activity Metrics

- Provide insights about how your storage is requested
- AllRequests, GetRequests, PutRequests, ListRequests, BytesDownloaded...

- Detailed Status Code Metrics

- Provide insights for HTTP status codes
- 200OKStatusCount, 403ForbiddenErrorCount, 404NotFoundErrorCode...

Free vs Paid



Storage Lens – Free vs. Paid

- Free Metrics

- Automatically available for all customers
- Contains around 28 usage metrics
- Data is available for queries for 14 days

- Advanced Metrics and Recommendations

- Additional paid metrics and features
- Advanced Metrics – Activity, Advanced Cost Optimization, Advanced Data Protection, Status Code
- CloudWatch Publishing – Access metrics in CloudWatch without additional charges
- Prefix Aggregation – Collect metrics at the prefix level
- Data is available for queries for 15 months

Metrics selection
Choose additional metrics and functionality.

Metrics selection

Free metrics
Includes usage metrics aggregated at the bucket level. Data is available for queries for 14 days.
[Learn more](#)

Advanced metrics and recommendations
Includes options for additional metrics and aggregations and other advanced capabilities. Data is available for queries for 15 months. See [Storage Lens metrics pricing](#) on the Management & Analytics tab.

Advanced metrics and recommendations features [Info](#)

Advanced metrics <input checked="" type="checkbox"/> Choose advanced metrics categories to display in the dashboard. Advanced metrics are not available at the prefix level.	CloudWatch publishing <input type="checkbox"/> Access metrics in CloudWatch without incurring separate CloudWatch metrics publishing charges. See CloudWatch Pricing Prefix-level metrics are not available in CloudWatch.	Prefix aggregation <input type="checkbox"/> Generate insights for usage metrics aggregated by top prefixes.
---	---	---

Advanced metrics categories
Specify which advanced metrics categories to display in the dashboard. [Learn more](#)

Activity
Generate metrics that show details about how your storage is requested, such as requests, bytes uploaded/downloaded, and errors aggregated by bucket.

Detailed status code metrics - new

Section 14: Amazon S3 Security

Glacier Vault Lock

- Lock glacier vault to adopt WORM (write once, read many)
- Create a vault lock policy, lock policy for future edits (can no longer be changed or deleted)
- Helpful for compliance and data retention

S3 Object Lock

- Versioning must be enabled, adopt WORM model to block object version deletion for a specified amount of time
- Retention mode:
 - Compliance
 - Object versions can't be overwritten or deleted by any user; including root user
 - Objects retention modes can't be changed and retention period can't be shortened
 - Governance
 - Most users can't overwrite or delete an object version or alter its lock settings
 - Some users can special permissions to change retention or delete the object
- Retention period: protect the object for a fixed period, can be extended
- Legal hold
 - Protect the object indefinitely, independent from retention period
 - Can be freely placed and removed using s3:PutObjectLegalHold

Section 15: CloudFront

CloudFront Overview

- Content Delivery Network (CDN)
- Improves read performance, content is cached at the edge to improve user experience
 - 216 edge locations (point of presence)
- DDoS protection, integration with Shield and WAF

Origins

- S3 Bucket
 - For distributing files and caching at the edge
 - Enhanced security with CloudFront Origin Access Control (OAC)
 - OAC replaces Origin Access Identity (OAI)
 - CloudFront and be used as an ingress (upload files to S3)
- Custom Origin (HTTP)
 - ALB, EC2, S3 website, any HTTP backend

CloudFront vs S3 Cross Region Replication

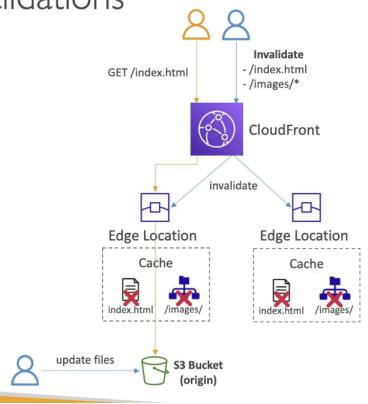
- CloudFront:
 - Global Edge network

- Files are cached for TTL
- Great for static content that must be available everywhere
- S3 Cross Region Replication
 - Must be setup for each region you want replication
 - Files updated near real-time
 - Read only
 - Great for dynamic content that needs to be available at low-latency in a few regions

Cache Invalidations

CloudFront – Cache Invalidations

- In case you update the back-end origin, CloudFront doesn't know about it and will only get the refreshed content after the TTL has expired
- However, you can force an entire or partial cache refresh (thus bypassing the TTL) by performing a CloudFront Invalidation
- You can invalidate all files (*) or a special path (/images/*)

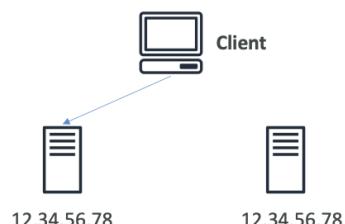
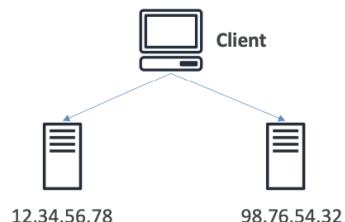


- In case you update the backend origin, CloudFront doesn't know and will only get refreshed content after TTL expired. Force entire or partial cache refresh (bypasses TTL) via CloudFront Invalidation

Unicast IP vs Anycast IP

Unicast IP vs Anycast IP

- Unicast IP: one server holds one IP address
- Anycast IP: all servers hold the same IP address and the client is routed to the nearest one



- Unicast IP: One server holds 1 IP address
- Anycast IP: all servers hold same IP address and client is routed to the nearest one

AWS Global Accelerator Overview

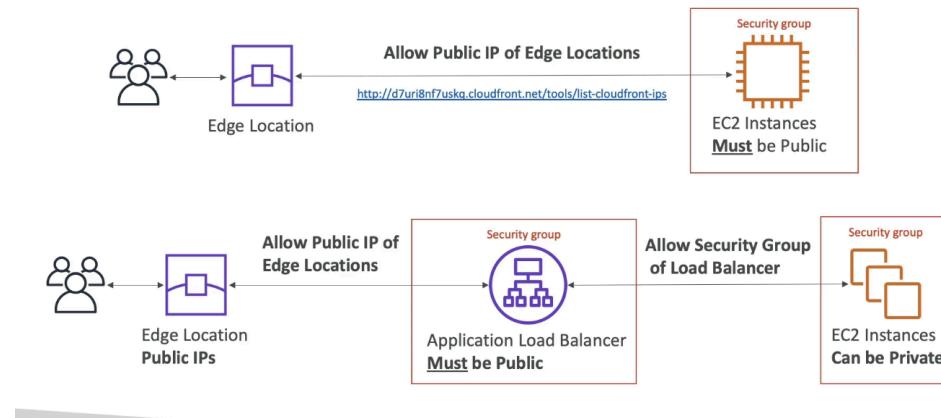
- If you have an application with global users, they go over public internet which adds latency with hops through routers.
- Leverages Anycast IP to use AWS internal network to route to application (edge location)
- 2 Anycast IP are created for application and anycast IP sends traffic directly to edge locations, then to application via AWS network
- Works with Elastic IP, EC2 Instances, ALB, NLB, public or private
- Consistent performance
 - Intelligent routing to lowest latency and fast regional failover
 - No issue with client cache (IP doesn't change)
- Health checks
 - Global accelerator performs health check of applications
 - Helps make application global (failover < 1 min for unhealthy)
 - Great for disaster recovery
- Security
 - Only 2 external IP need whitelisted
 - DDoS protection via AWS Shield

AWS Global Accelerator vs CloudFront

- Both use edge locations and AWS global network, integrations with AWS Shield
- CloudFront
 - Improves performance for both cacheable content (images / videos)
 - Dynamic content, with content served at edge
- Global Accelerator
 - Improves performance for wide range of applications over TCP or UDP
 - Proxying packers at edge to applications running in 1+ regions
 - Good for non HTTP use cases like gaming (UDP), voice over IP
 - Good for HTTP cases that require static IP or require deterministic, fast regional failover

ALB as an Origin

CloudFront – ALB or EC2 as an origin



Geo Restriction

- Restrict who can access distribution
 - Allowlist: allow users to access content only if they're on a list of approved countries
 - Blocklist: prevent users from accessing content if they're on a banned country list
- Country determined via 3rd party Geo-IP database

CloudFront Advanced Concepts

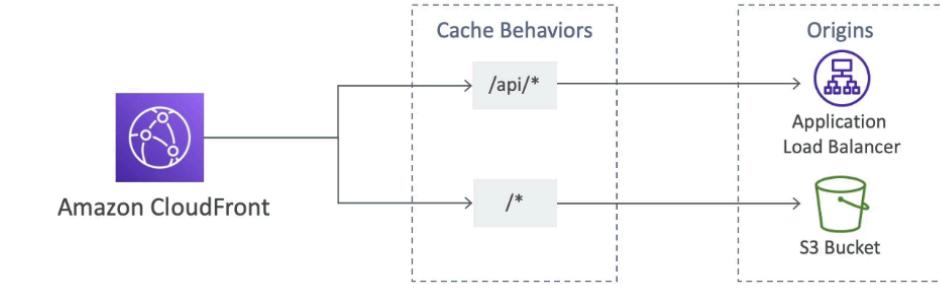
Pricing

- Cost of data out per edge location varies based on location of edge location

Price Classes

- Can reduce the number of edge locations for cost reduction
- 3 price classes:
 1. Price Class All: all regions, best performance
 2. Price Class 200: most regions, excludes most expensive regions
 3. Price Class 300: only the least expensive regions

Multiple Origin

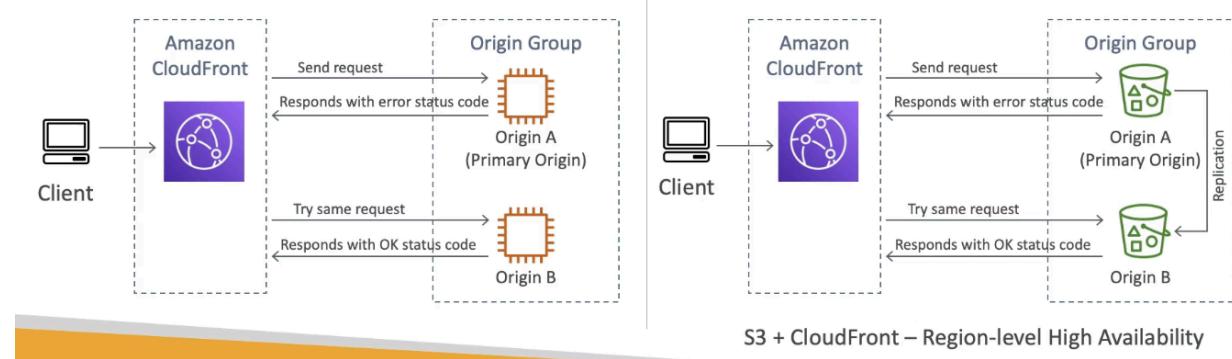


- To route to different kind of origins based on content type based on path pattern

Origin Groups

CloudFront – Origin Groups

- To increase high-availability and do failover
- Origin Group: one primary and one secondary origin
- If the primary origin fails, the second one is used

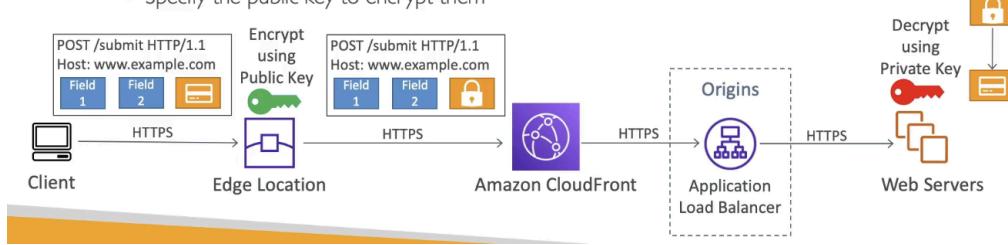


- Increase high availability and do failover
- Origin group: one primary and one secondary
- If the primary origin fails, second is used

Field Level Encryption

CloudFront – Field Level Encryption

- Protect user sensitive information through application stack
- Adds an additional layer of security along with HTTPS
- Sensitive information encrypted at the edge close to user
- Uses asymmetric encryption
- Usage:
 - Specify set of fields in POST requests that you want to be encrypted (up to 10 fields)
 - Specify the public key to encrypt them

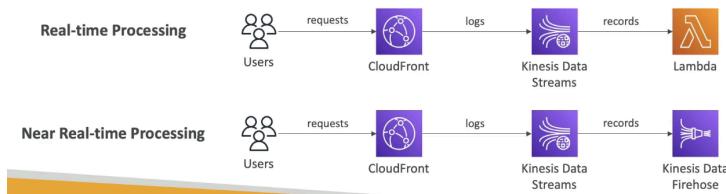


- Protect user sensitive info through application stack
- Additional layer along with HTTPS
- Sensitive information encrypted at the edge close to user
- Uses asymmetric encryption
- Usage:
 - Specify set of fields in POST you want encrypted (up to 10 fields)
 - Specify the public key to encrypt

Real Time Logs

CloudFront – Real Time Logs

- Get real-time requests received by CloudFront sent to Kinesis Data Streams
- Monitor, analyze, and take actions based on content delivery performance
- Allows you to choose:
 - Sampling Rate – percentage of requests for which you want to receive
 - Specific fields and specific Cache Behaviors (path patterns)



- Real time requests received by CloudFront sent to Kinesis Data stream
 - Monitor, analyze based on content delivery performance
- Choose:
 - Sampling rate: % of requests you want to receive
 - Specific fields and specific cache behaviors (path patterns)

Section 16: AWS Storage Extras

AWS Snow Family Overview

- Highly secure portable devices to collect and process data at the edge AND migrate data in / out of AWS
 - Data migration: Snowcone, Snowball Edge, Snowmobile
 - Edge computing: Snowcone, snowball edge

Data Migration with Snow

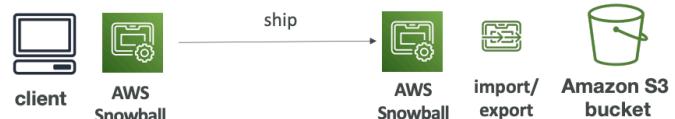
Diagrams

- Offline devices to perform data migrations
 - More than a week to transfer data, use Snowball
 - Done physically

- Direct upload to S3:



- With Snow Family:



AWS Snow Family for Data Migrations



Snowcone



Snowball Edge



Snowmobile

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

Snowball Edge

- Physical data transport to move TB or PB of data in / out of AWS
 - Alternative to moving data over network and paying network fees
- Pay per data transfer job

- Provide block storage and S3 compatible object storage
- Snowball Edge Storage Optimized
 - 80 TB of HDD or 210 TB of NVMe capacity for block volume and S3 compatible object storage
- Snowball Edge Compute Optimized
 - 42 TB or 28 TB NVMe capacity for block volume and S3 compatible object storage
- Use case: large data cloud migrations, disaster recovery

AWS Snow Cone & Snow Cone SSD

- Small portable computing anywhere, rugged and secure to withstand harsh environments
 - Light, device used for edge computing, storage, and data transfer
 - Must provide own battery / cables
 - Used where snowball does not fit (space limited)
- Sent back to AWS offline or connect to internet and use AWS DataSync to send data
- Snow Cone - 8 TB of HDD
- Snow Cone SSD - 14 TB SSD

AWS Snowmobile

- Transfer EBs of data with each has 100 PB capacity (can use multiple in parallel)
 - Better if need > 10 PB
- High security, temp controlled, GPS, 24/7 video surveillance

Usage Process

1. Request Snowball device from AWS for delivery
2. Install snowball client / AWS OpsHub on servers
3. Connect snowball to servers and copy files via client
4. Ship back to AWS
5. Data loaded in S3 bucket
6. Snowball completely wiped

What is Edge Computing?

- Process data while created on edge location
 - Edge location is somewhere that does not have internet access or computing power
- Use Snowball Edge / Snow Cone device
- Use cases: preprocess data, ML at edge, transcoding media streams...
- Ship back to AWS

Snow Family – Edge Computing

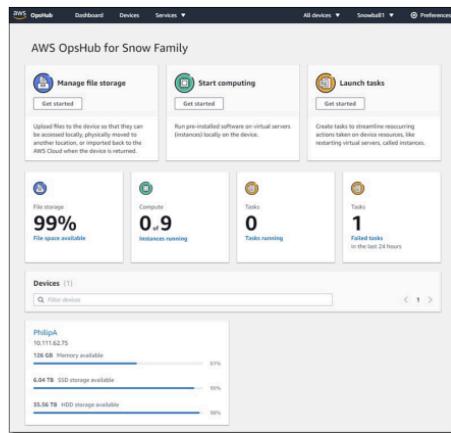
- # Snow Family – Edge Computing
- Snowcone & Snowcone SSD (smaller)
 - 2 CPUs, 4 GB of memory, wired or wireless access
 - USB-C power using a cord or the optional battery
 - Snowball Edge – Compute Optimized
 - 104 vCPUs, 416 GiB of RAM
 - Optional GPU (useful for video processing or machine learning)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



AWS OpsHub

AWS OpsHub

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



<https://aws.amazon.com/blogs/aws/aws-snowball-edge-update/>

- Software installed on computer to manage Snow Family device

Architecture: Snowball into Glacier



- Snowball cannot import to Glacier directly, must use S3 with lifecycle policy

Amazon FSx Overview

- Launch 3rd party high performance file systems in AWS; fully managed service

Amazon FSx for Windows (File Server)

Amazon FSx for Windows (File Server)



- FSx for Windows is a fully managed Windows file system share drive
 - Supports SMB protocol & Windows NTFS
 - Microsoft Active Directory integration, ACLs, user quotas
 - Can be mounted on Linux EC2 instances
 - Supports Microsoft's Distributed File System (DFS) Namespaces (group files across multiple FS)
 - Scale up to 10s of GB/s, millions of IOPS, 100s PB of data
 - Storage Options:
 - SSD – latency sensitive workloads (databases, media processing, data analytics, ...)
 - HDD – broad spectrum of workloads (home directory, CMS, ...)
 - Can be accessed from your on-premises infrastructure (VPN or Direct Connect)
 - Can be configured to be Multi-AZ (high availability)
 - Data is backed-up daily to S3
-
- Fully managed Windows file system share drive
 - Supports SMB protocol & Windows NTFS
 - Microsoft Active Directory integration, ACLs, user quotas
 - Can be mounted on Linux EC2 instances
 - Supports Microsoft's Distributed File System (DFS) Namespaces (groups files across multiple FS)
 - Scale up to 10s of GB/s, millions of IOPS...
 - Storage options:

- SSD – latency sensitive workloads (DB, media processing, data analytics...)
- HDD – broad spectrum of workloads
- Can be accessed from on premise infrastructure (VPN or Direct Connect)
- Can be configured for Multi AZ (high availability)
- Data backed up to S3 daily

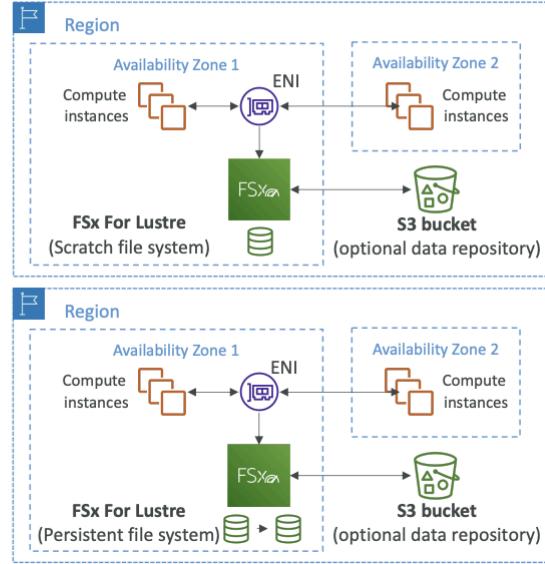
Amazon FSx for Lustre

- ## Amazon FSx for Lustre
- Lustre is a type of parallel distributed file system, for large-scale computing
 - The name Lustre is derived from “Linux” and “cluster
 - Machine Learning, High Performance Computing (HPC)
 - Video Processing, Financial Modeling, Electronic Design Automation
 - Scales up to 100s GB/s, millions of IOPS, sub-ms latencies
 - Storage Options:
 - SSD – low-latency, IOPS intensive workloads, small & random file operations
 - HDD – throughput-intensive workloads, large & sequential file operations
 - Seamless integration with S3
 - Can “read S3” as a file system (through FSx)
 - Can write the output of the computations back to S3 (through FSx)
 - Can be used from on-premises servers (VPN or Direct Connect)
 - Type of parallel distributed file system for large scale computing
 - Name derived from Linux and cluster
 - ML, high performance computing (HPC)
 - Video processing, financial modeling...
 - Scales up to 100s GB, millions of IOPS, low latency
 - Storage options:
 - SSD: low latency, IOPS sensitive workloads, small & random file operations
 - HDD: throughput intensive workloads, large & sequential file operations
 - Seamless integration with S3
 - Can “read S3” as file system (through FSx)
 - Can write the output of computations back to S3 (through FSx)
 - Can be used from on premise servers

FSx File System Deployment Options

FSx Lustre - File System Deployment Options

- Scratch File System
 - Temporary storage
 - Data is not replicated (doesn't persist if file server fails)
 - High burst (6x faster, 200MBps per TiB)
 - Usage: short-term processing, optimize costs
- Persistent File System
 - Long-term storage
 - Data is replicated within same AZ
 - Replace failed files within minutes
 - Usage: long-term processing, sensitive data

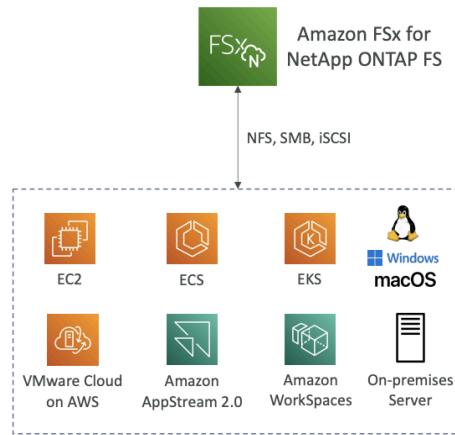


- Scratch File System
 - Temp storage, data not replicated (doesn't persist if file server fails)
 - High burst (6x faster, 200 MBps per TB)
 - Usage: short term processing, optimize cost
- Persistent File System
 - Long term storage, data replicated within same AZ
 - Replace failed files within minutes
 - Use case: long term processing, sensitive data

Amazon FSx for NetApp ONTAP

Amazon FSx for NetApp ONTAP

- Managed NetApp ONTAP on AWS
- File System compatible with NFS, SMB, iSCSI protocol
- Move workloads running on ONTAP or NAS to AWS
- Works with:
 - Linux
 - Windows
 - MacOS
 - VMware Cloud on AWS
 - Amazon Workspaces & AppStream 2.0
 - Amazon EC2, ECS and EKS
- Storage shrinks or grows automatically
- Snapshots, replication, low-cost, compression and data de-duplication
- Point-in-time instantaneous cloning (helpful for testing new workloads)



- Managed NetApp ONTAP on AWS
- Compatible with NFS, SMB, iSCSI protocol
- Move workloads running on ONTAP or NAS to AWS
- Works with many OS and services
- Storage auto scaling
- Snapshots, replication, low-cost, compression and data deduplication
- Point in time instantaneous cloning (for testing new workloads)

Amazon FSx for OpenZFS

Amazon FSx for OpenZFS



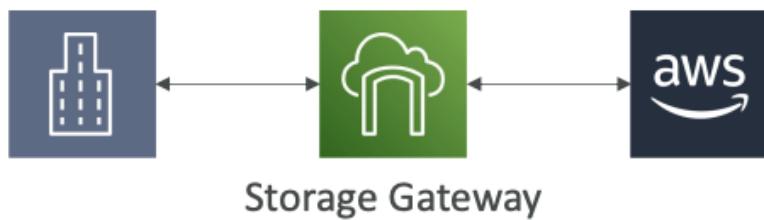
- Managed OpenZFS file system on AWS
- File System compatible with NFS (v3, v4, v4.1, v4.2)
- Move workloads running on ZFS to AWS
- Works with:
 - Linux
 - Windows
 - MacOS
 - VMware Cloud on AWS
 - Amazon Workspaces & AppStream 2.0
 - Amazon EC2, ECS and EKS
- Up to 1,000,000 IOPS with < 0.5ms latency
- Snapshots, compression and low-cost
- Point-in-time instantaneous cloning (helpful for testing new workloads)



- Managed OpenZFS file system in AWS
- File system compatible with NFS
- Move workloads running on ZFS to AWS
- Works with all OS and services
- Up to 1 million IOPS with < 0.5 ms latency
- Snapshots, compression, low cost
- Point in time instantaneous cloning

AWS Storage Gateway

AWS Storage Cloud Native Options

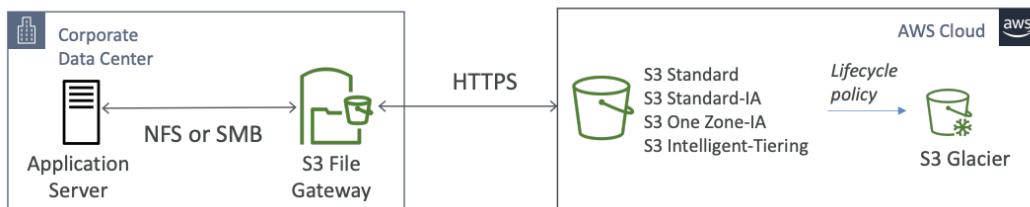


- Bridge between on premise and cloud
 - Use case: disaster recovery, backup and restore, tiered storage, on premise cache & low latency file access
- Types of Storage Gateway
 - S3 File Gateway
 - FSx File Gateway
 - Volume Gateway
 - Tape Gateway

S3 File Gateway

Amazon S3 File Gateway

- Configured S3 buckets are accessible using the NFS and SMB protocol
- Most recently used data is cached in the file gateway
- Supports S3 Standard, S3 Standard IA, S3 One Zone A, S3 Intelligent Tiering
- Transition to S3 Glacier using a Lifecycle Policy
- Bucket access using IAM roles for each File Gateway
- SMB Protocol has integration with Active Directory (AD) for user authentication



- Configured S3 buckets are accessible using NFS and SMB protocol
 - SMB has integration with Active Directory for user authentication
- Most recently used data is cached in file gateway
- Supports all S3 tiers except Glacier
 - Transition to S3 Glacier via lifecycle policy
- Bucket access using IAM roles for each file gateway

FSx File Gateway

- # Amazon FSx File Gateway
- Native access to Amazon FSx for Windows File Server
 - Local cache for frequently accessed data
 - Windows native compatibility (SMB, NTFS, Active Directory...)
 - Useful for group file shares and home directories

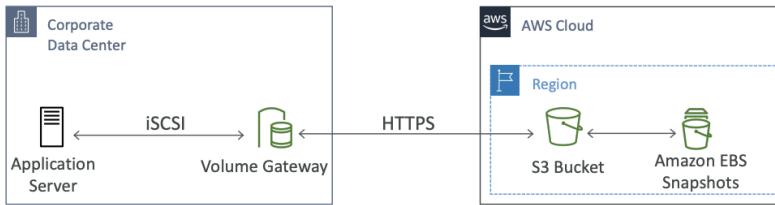


- Native access to Amazon FSx for Windows File Server
 - Windows native compatibility (SMB, NTFS, Active Directory...)
- Local cache for frequently accessed data
- Useful for group file shares and home directories

Volume Gateway

Volume Gateway

- Block storage using iSCSI protocol backed by S3
- Backed by EBS snapshots which can help restore on-premises volumes!
- Cached volumes: low latency access to most recent data
- Stored volumes: entire dataset is on premise, scheduled backups to S3

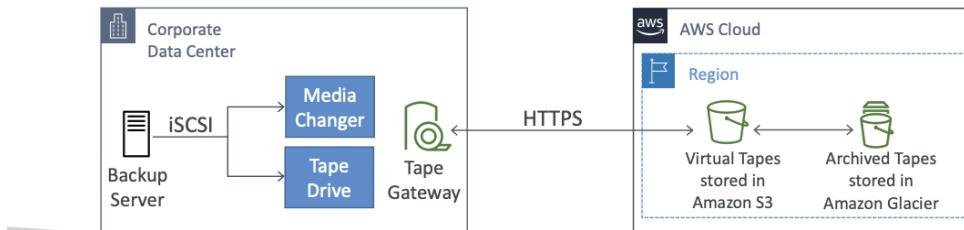


- Block storage using iSCSI protocol backed by S3
- Backed by EBS snapshots to help restore on premise volumes
- Cached Volumes: low latency access to most recent data
- Stored volumes: entire dataset is on premise, scheduled backups to S3

Tape Gateway

Tape Gateway

- Some companies have backup processes using physical tapes (!)
- With Tape Gateway, companies use the same processes but, in the cloud
- Virtual Tape Library (VTL) backed by Amazon S3 and Glacier
- Back up data using existing tape-based processes (and iSCSI interface)
- Works with leading backup software vendors

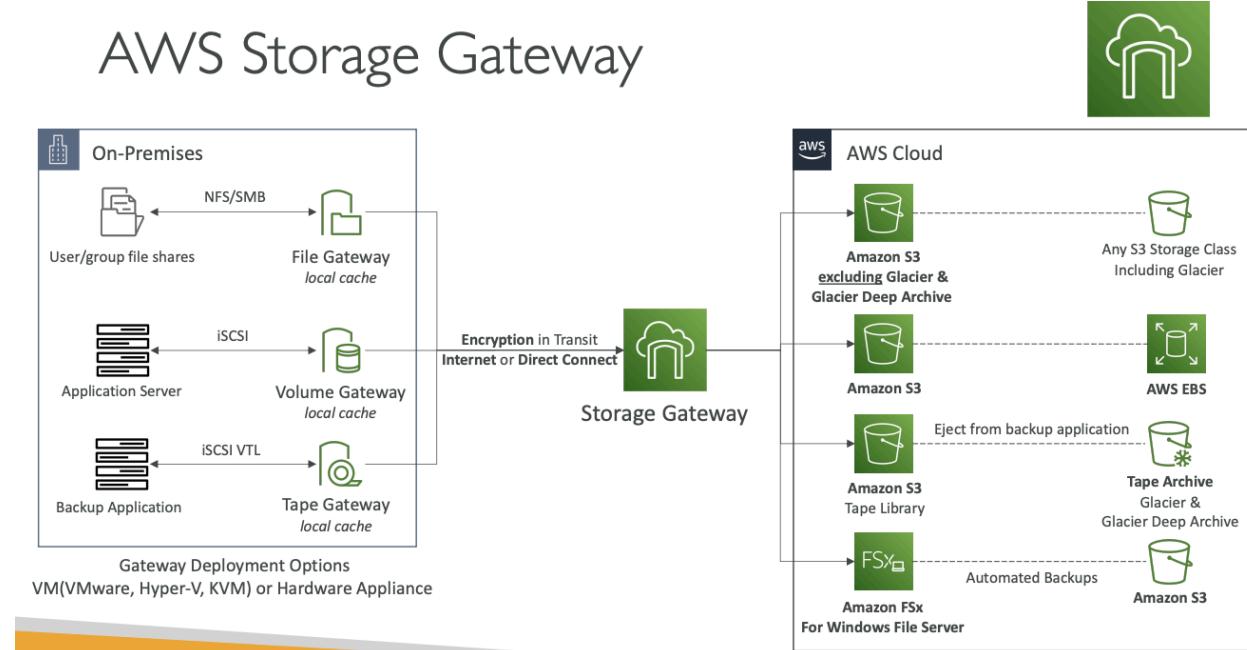


- Backup processes using physical tapes
- Virtual Tape Library (VTL) backed by S3 and Glacier
- Back up data using existing tape based processes
- Works with leading backup software vendors

Storage Gateway – Hardware Appliance

- All gateways require on premise data server for virtualization, but if it isn't there by default, use Storage Gateway Hardware Appliance
 - Works with all Gateways and has required resources already
 - Helps for daily NFS backups in small data centers

Storage Gateway Summary



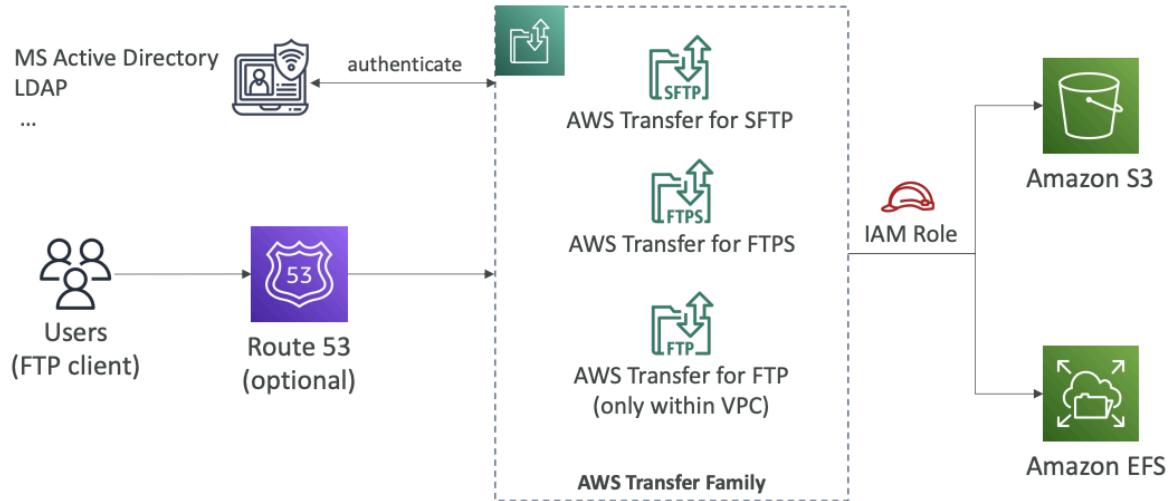
AWS Transfer Family

AWS Transfer Family



- A fully-managed service for file transfers into and out of Amazon S3 or Amazon EFS using the FTP protocol
- Supported Protocols
 - AWS Transfer for FTP (File Transfer Protocol (FTP))
 - AWS Transfer for FTPS (File Transfer Protocol over SSL (FTPS))
 - AWS Transfer for SFTP (Secure File Transfer Protocol (SFTP))
- Managed infrastructure, Scalable, Reliable, Highly Available (multi-AZ)
- Pay per provisioned endpoint per hour + data transfers in GB
- Store and manage users' credentials within the service
- Integrate with existing authentication systems (Microsoft Active Directory, LDAP, Okta, Amazon Cognito, custom)
- Usage: sharing files, public datasets, CRM, ERP, ...
 - Fully managed service for file transfer in / out of S3 or EFS via FTP protocol
 - Supports:
 - AWS Transfer for FTP (File transfer protocol, unencrypted)
 - AWS Transfer for FTPS (encrypted)
 - AWS Transfer for SFTP (secure FTP, encrypted in flight)
 - Managed infrastructure, scalable, reliable, high availability (multi AZ)
 - Pay per provisioned endpoint per hour + data transfers in GB
 - Store and manage users' credentials within the service
 - Integration with existing authentication systems
 - Usage: sharing files, public datasets...

AWS Transfer Family



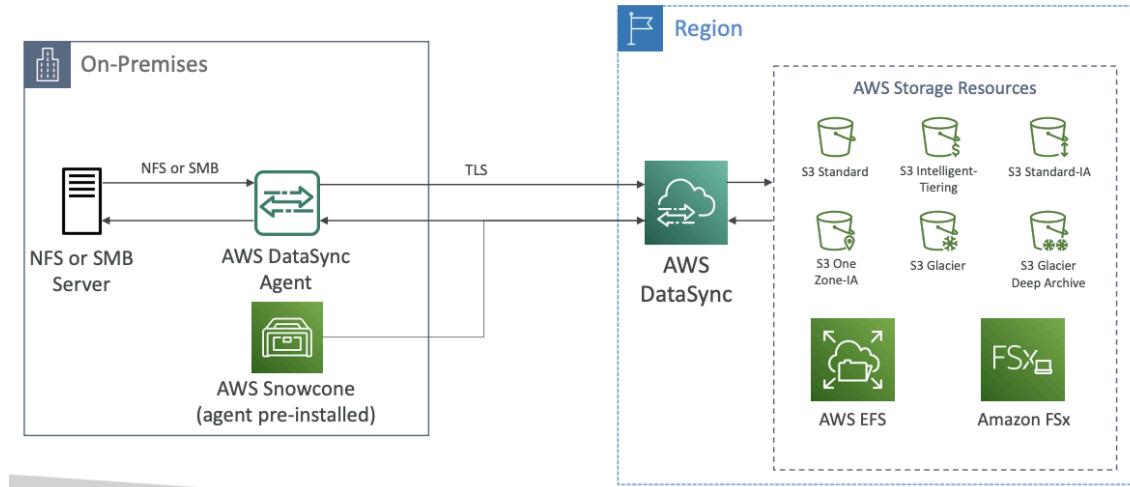
AWS DataSync

AWS DataSync

- Move large amount of data to and from
 - On-premises / other cloud to AWS (NFS, SMB, HDFS, S3 API...) – needs agent
 - AWS to AWS (different storage services) – no agent needed
- Can synchronize to:
 - Amazon S3 (any storage classes – including Glacier)
 - Amazon EFS
 - Amazon FSx (Windows, Lustre, NetApp, OpenZFS...)
- Replication tasks can be scheduled hourly, daily, weekly
- File permissions and metadata are preserved (NFS POSIX, SMB...)
- One agent task can use 10 Gbps, can setup a bandwidth limit
 - Move large amount of data to and from on prem or other cloud to AWS (needs agent), or AWS to AWS (no agent needed)
 - Can synchronize to:
 - S3 (any tier)
 - EFS
 - FSx (all)
 - Replication tasks scheduled hourly, daily, weekly
 - File permissions and metadata are preserved
 - One agent task can use 10 Gbps, can set bandwidth limit

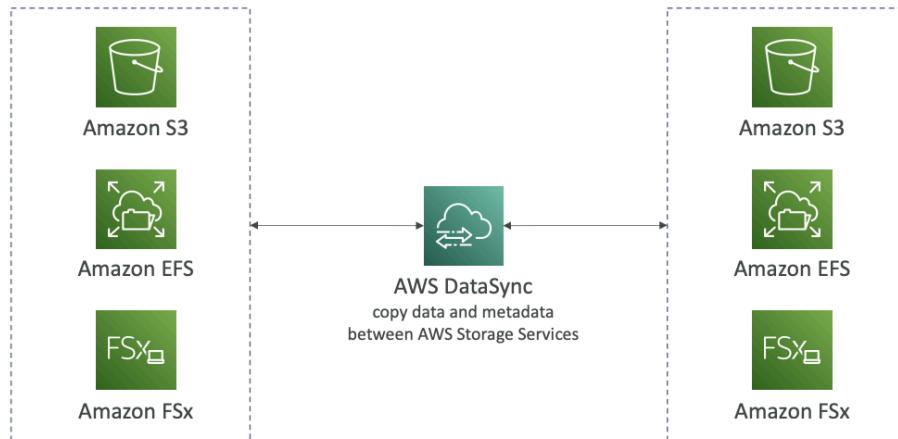
AWS DataSync

NFS / SMB to AWS (S3, EFS, FSx...)



AWS DataSync

Transfer between AWS storage services



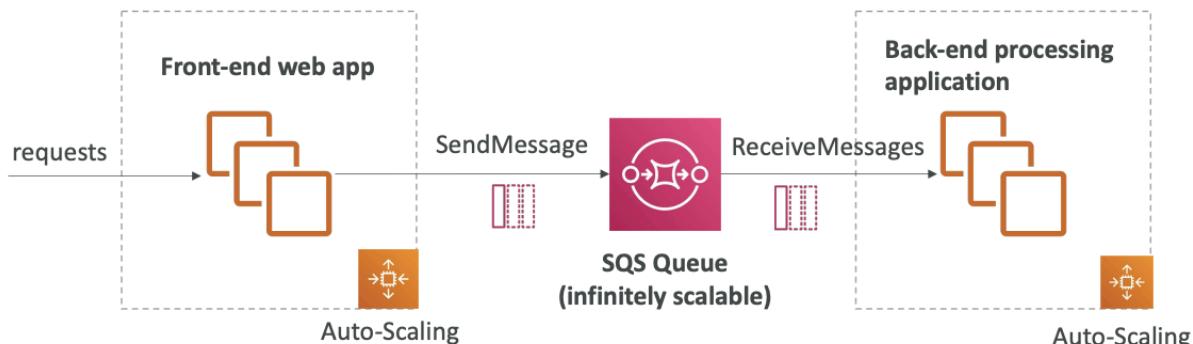
All AWS Storage Options Compared

Storage Comparison

- S3: Object Storage
- S3 Glacier: Object Archival
- EBS volumes: Network storage for one EC2 instance at a time
- Instance Storage: Physical storage for your EC2 instance (high IOPS)
- EFS: Network File System for Linux instances, POSIX filesystem
- FSx for Windows: Network File System for Windows servers
- FSx for Lustre: High Performance Computing Linux file system
- FSx for NetApp ONTAP: High OS Compatibility
- FSx for OpenZFS: Managed ZFS file system
- Storage Gateway: S3 & FSx File Gateway, Volume Gateway (cache & stored), Tape Gateway
- Transfer Family: FTP, FTPS, SFTP interface on top of Amazon S3 or Amazon EFS
- DataSync: Schedule data sync from on-premises to AWS, or AWS to AWS
- Snowcone / Snowball / Snowmobile: to move large amount of data to the cloud, physically
- Database: for specific workloads, usually with indexing and querying

Section 17: Decoupling applications: SQS, SNS, Kinesis, Active MQ

SQS to decouple between application tiers



Section 19: Serverless Overviews from Solution Architect Perspective

DynamoDB Global Tables

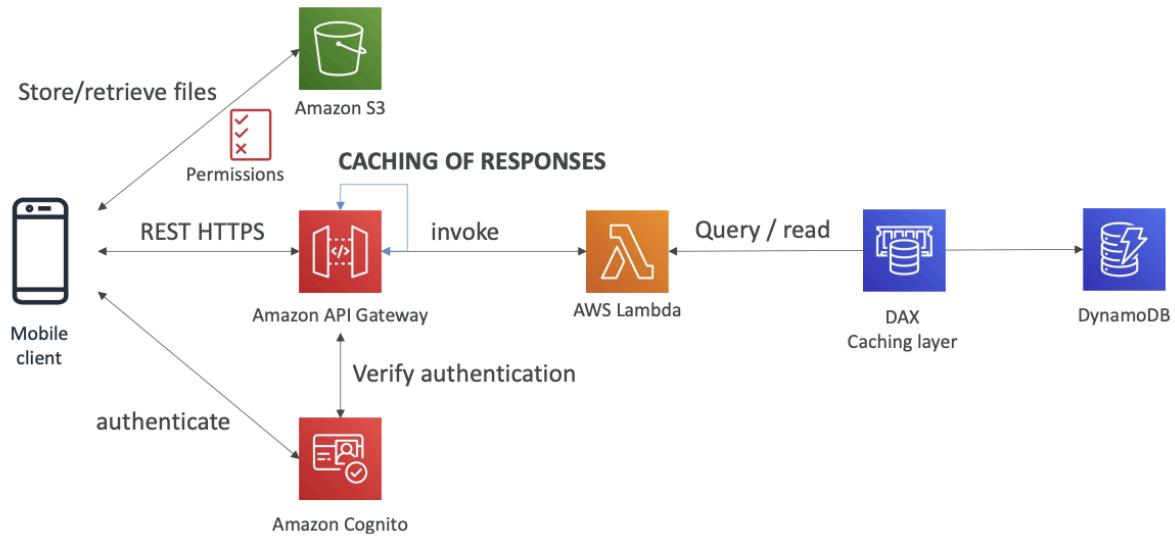
- Make table accessible with low latency in multiple regions
- Active active replication (2 way replication)
- Applications can read / write to table in any region

- Must enable DynamoDB streams as prerequisite

Section 20: Serverless Solution Architecture Discussions

Mobile Application

Mobile app: caching at the API Gateway

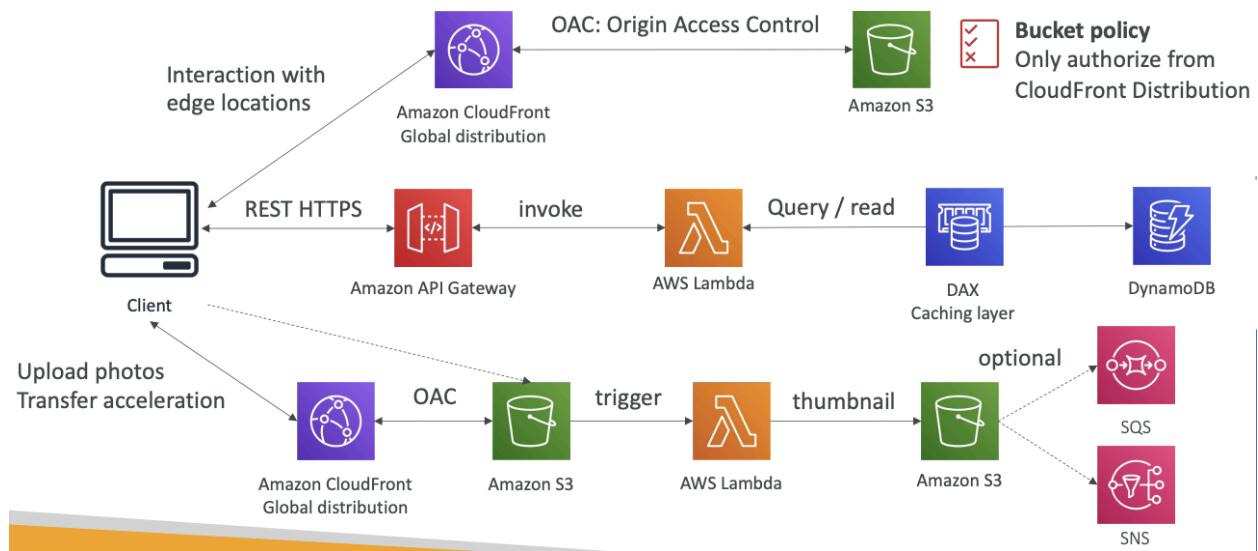


- Serverless REST API: HTTPS, API Gateway, Lambda, DynamoDB
- Using Cognito to generate temporary credentials to access S3 bucket with restricted policy. App users can directly access AWS resources this way. Pattern can be applied to DynamoDB, Lambda...
- Caching the reads on DynamoDB using DAX
- Caching the REST requests at the API Gateway level
- Security for authentication and authorization with Cognito

- For REST API, use API Gateway + Lambda for a serverless architecture. To authenticate to access S3, use Cognito Identity Pool to retrieve files from S3.
 - With increased high read throughput of static data, DAX can be used alongside DynamoDB as a serverless DB. API GW can also cache some responses

Serverless Website

Thumbnail Generation flow

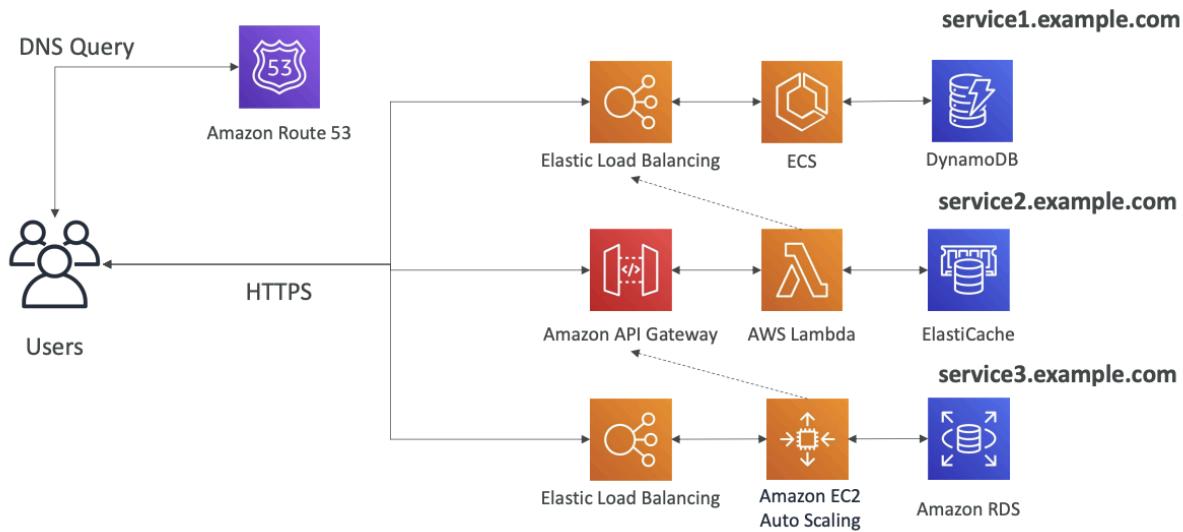


AWS Hosted Website Summary

- We've seen static content being distributed using CloudFront with S3
 - The REST API was serverless, didn't need Cognito because public
 - We leveraged a Global DynamoDB table to serve the data globally
(we could have used Aurora Global Database)
 - We enabled DynamoDB streams to trigger a Lambda function
 - The lambda function had an IAM role which could use SES
 - SES (Simple Email Service) was used to send emails in a serverless way
 - S3 can trigger SQS / SNS / Lambda to notify of events
-
- To serve clients globally, use CloudFront to access S3. To securely access, use Origin Access Control (OAC) to only allow CloudFront to access S3. REST API needs API GW to invoke lambda function and a serverless DB like DynamoDB. To send welcome emails, enable Dynamo streams to invoke a Lambda function to send emails via SES.
 - If users upload images for thumbnails, use CloudFront + OAC to S3 again and have a Lambda function triggered to create the thumbnail to S3.

Microservices Architecture

Micro Services Environment

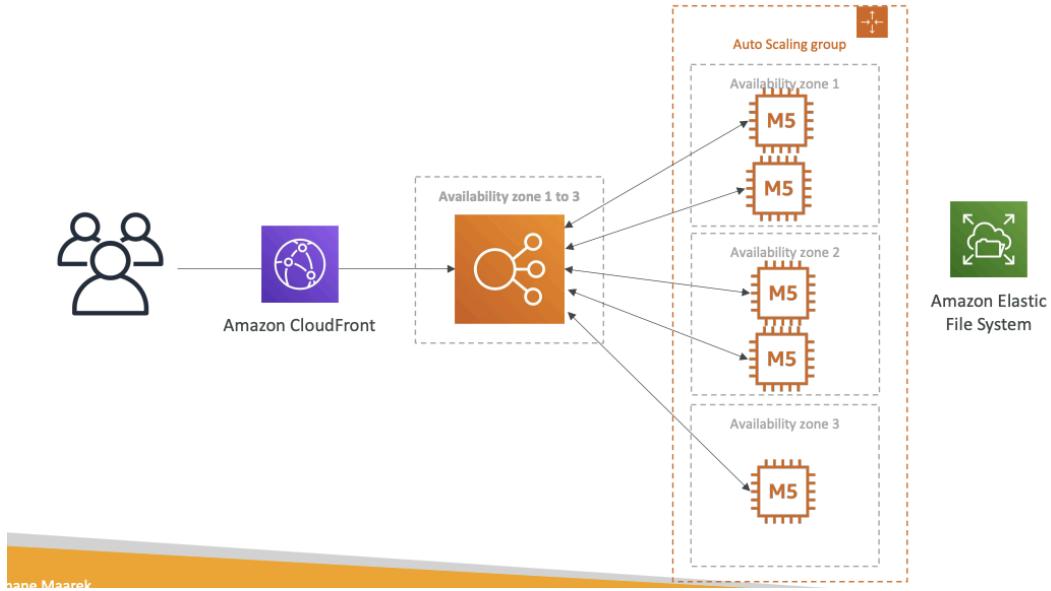


- Design microservices as you want
- Synchronous patterns: API GW, LB
- Asynchronous pattern: SQS, Kinesis, SNS, Lambda triggers (S3)
- Challenges:
 - Repeated overhead for creating each microservices
 - Issues with optimizing server utilization
 - Complexity of running multiple versions of multiple microservices simultaneously
 - Proliferation of client side code requirements to integrate with many separate services
- Can be solved by serverless patterns:
 - API GW, Lambda scale automatically and pay per usage
 - Easily clone API, reproduce environments

Software Updates Offloading

- When a new software update is out, lots of requests are made to update and content is distributed in mass over the network, thus very costly. How to optimize cost?

Easy way to fix things!



- Use CloudFront because there are no changes to architecture and will cache software update files at the edge
 - Software files not dynamic, CloudFront auto scales

Section 21: Databases in AWS

Choosing the right DB

Choosing the Right Database

- We have a lot of managed databases on AWS to choose from
- Questions to choose the right database based on your architecture:
 - Read-heavy, write-heavy, or balanced workload? Throughput needs? Will it change, does it need to scale or fluctuate during the day?
 - How much data to store and for how long? Will it grow? Average object size? How are they accessed?
 - Data durability? Source of truth for the data ?
 - Latency requirements? Concurrent users?
 - Data model? How will you query the data? Joins? Structured? Semi-Structured?
 - Strong schema? More flexibility? Reporting? Search? RDBMS / NoSQL?
 - License costs? Switch to Cloud Native DB such as Aurora?

Database Types

- RDBMS (SQL): RDS / Aurora, great for joins
- NoSQL: no joins; Dynamo (JSON), ElastiCache (key value), Neptune (graphs), DocumentDB (Mongo), Keyspaces (Apache Cassandra)
- Object Store: S3, Glacier
- Data Warehouse: SQL Analytics; Redshift, Athena, EMR
- Search: OpenSearch (JSON), free text, unstructured searches
- Graphs: Neptune - displays relationships between data
- Ledger: Amazon Quantum Ledger DB
- Time series: Amazon Timestream

RDS Summary

- Managed PostgresQL / MySQL / Oracle / SQL Server / DB2 / MariaDB / Custom
- Provisioned RDS instance size and EBS volume type & size
 - Auto scaling for storage
- Support read replicas and multi AZ
- Security via IAM, SG, KMS, SSL in transit
 - Support for IAM Authentication, Secrets Manager
- Automated backup with point in time (35 days)
 - Manual snapshot for long term storage
- Managed and scheduled maintenance (downtime)
- RDS Custom for access to and customize underlying instance (Oracle & SQL server)
- Use case: relational DB, perform SQL queries, transactions

Aurora Summary

- Compatible API for PostgreSQL / MySQL, separation of storage and compute
- Storage: data is stored in 6 replicas, across 3 AZ – highly available, self-healing, auto-scaling
- Compute:
 - Cluster of DB Instance across multiple AZ, auto-scaling of Read Replicas
- Cluster:
 - Custom endpoints for writer and reader DB instances
- Same security / monitoring / maintenance features as RDS
- Aurora Serverless
 - For unpredictable / intermittent workloads, no capacity planning
- Aurora Global:
 - Up to 16 DB Read Instances in each region, < 1 second storage replication
- Aurora Machine Learning:
 - ML using SageMaker & Comprehend on Aurora
- Aurora Database Cloning:
 - New cluster from existing one, faster than restoring a snapshot

- Use case: same as RDS, but with less maintenance / more flexibility / more performance / more features

ElastiCache

- Managed Redis / Memcached (similar as RDS, but for cache)
 - In memory data store with sub ms latency
- Redis clustering and multi AZ, read replicas (sharding)
- Security via IAM, SG, KMS, Redis Auth
- Backup / snapshot / point in time restore
- Managed and scheduled maintenance
- Some application code changes needed
- Use case: key / value store, frequent read, cache DB queries, session data, cannot use SQL

DynamoDB

- NoSQL AWS managed, serverless, millisecond latency
- Capacity mode:
 - Provisioned with optional auto scaling
 - On demand
- Can replace ElastiCache as key value store (storing session data, TTL)
- Highly available, multi AZ by default, read / write decoupled, transaction capability
- DAX cluster for read cache, low read latency
- Security, authentication, authorization done via IAM
- Event processing: Dynamo Streams with Lambda or Kinesis Data Stream
- Global tables: active active setup
- Automated backup up to 35 days with point in time (restore to new table) or on demand backup
- Export to S3 without using RCU within point in time window or import to S3 without using WCU
- Great to rapidly evolve schemas
- Use case: serverless DB, distributed serverless cache

S3

- Key value store for objects
 - Large objects, not small
- Serverless, infinite scaling, max size 5 TB, versioning
- Tiers of S3 with lifecycle policies
- Features:
 - Versioning, encryption, replication, MFA delete, access logs...
- Security:
 - IAM, bucket policy, ACL, access points, object lambda, CORS, object / vault lock

- Encryption:
 - SSE S3, SSE KMS, SSE C, client side, TLS, default encryption
- Batch operations via S3 Batch, list files using S3 Inventory
- Performance
 - Multi part upload, S3 transfer acceleration, S3 select
- Automation:
 - S3 Event Notifications
- Use case: static files, key value store for big files, website hosting

DocumentDB

- AWS version of MongoDB (NoSQL)
 - Used to store, query, index JSON data
 - Auto grows in increments of 10GB
 - Auto scales to workloads
 - Fully managed, highly available with replication across 3 AZ
- Similar deployment concepts as Aurora

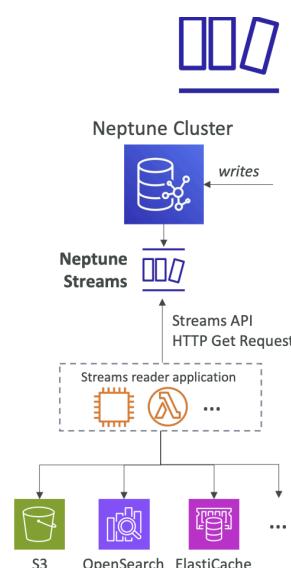
Amazon Neptune

- Fully managed graph DB (like social network)
 - Highly available across 3 AZ, up to 15 read replicas
- Build and run apps working with highly connected datasets
 - Optimized for complex and hard queries
 - Store billions of relations and query the graph with ms latency
- Use cases: knowledge graphs, fraud detection, recommendations, social networking

Neptune Streams

Amazon Neptune – Streams

- Real-time ordered sequence of every change to your graph data
- Changes are available immediately after writing
- No duplicates, strict order
- Streams data is accessible in an HTTP REST API
- Use cases:
 - Send notifications when certain changes are made
 - Maintain your graph data synchronized in another data store (e.g., S3, OpenSearch, ElastiCache)
 - Replicate data across regions in Neptune



- Real time ordered sequence of every change to graph data
 - Changes immediately after writing
- No duplicates, strict ordering
- Stream data accessible via HTTP REST API
- Use cases:
 - Send notifications when certain changes made, maintain graph data synchronized in another data store
 - Replicate data across regions in Neptune

Amazon Keyspaces (for Apache Cassandra)

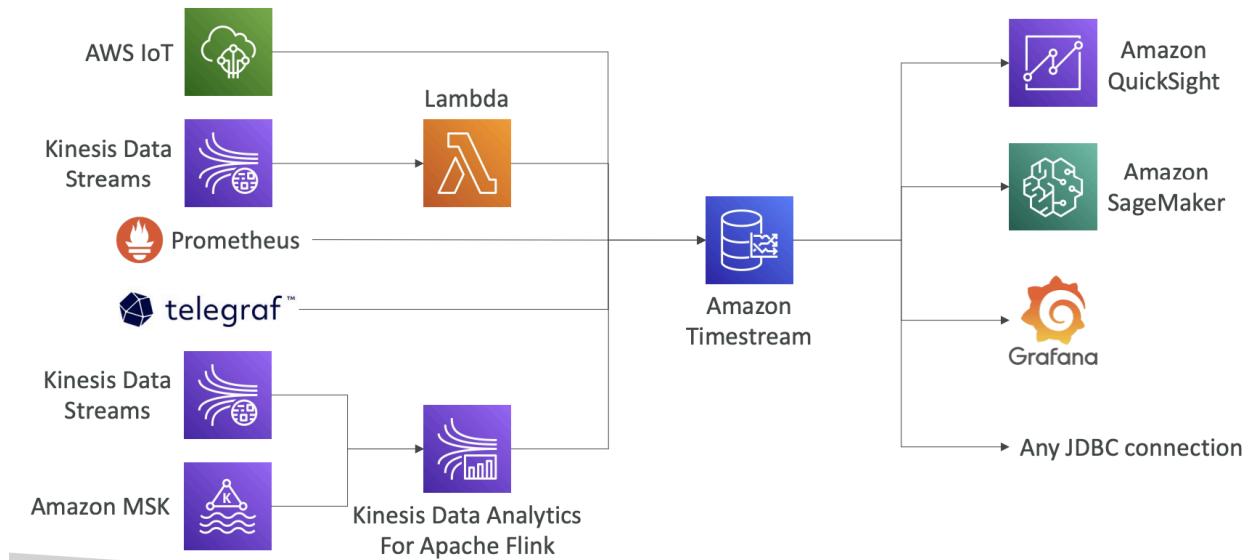
- Cassandra is NoSQL, Keyspaces is AWS managed, serverless, scalable, highly available
- Automatically scales tables up / down based on traffic
 - Tables replicated 3 times across multiple AZ
- Use Cassandra Query Language (CQL)
- Single digit ms latency at any scale, 1000s of requests / second
- Capacity: on demand or provisioned with auto scaling
- Encryption, backup, point in time up to 35 days
- Use cases: Apache Cassandra

Amazon QLDB

- Quantum Ledger DB
 - Ledge is a book for recording financial transaction
 - Fully managed, serverless, highly available, replication across 3 AZ
- Used to review history of all changes made to application data over time
- Immutable system: no entry can be removed or modified, cryptographically verifiable
 - Journal behind the scenes to have a sequence number any time revision made
- 2-3x better performance than others, use SQL
- Difference with Amazon Managed Blockchain: no decentralization component

Amazon Timestream

Amazon Timestream – Architecture



- Fully managed, fast, scalable, serverless time series database
 - Automatically scales up/down to adjust capacity
 - 1000s times faster & 1/10th the cost of relational databases
- Scheduled queries, multi-measure records, SQL compatibility
- Data storage tiering: recent data kept in memory and historical data kept in a cost-optimized storage
- Built-in time series analytics functions (helps you identify patterns in your data in near real-time)
- Encryption in transit and at rest
- Use case: data in relation to time

Section 22: Data and Analytics

Amazon Athena

- Serverless query service to analyze data stored in S3 via SQL to query files (business analytics, reporting, etc...)
 - Supports CSV, JSON...
 - Pricing: \$5 per TB of data scanned
- Commonly used with QuickSight for reporting / dashboards
- Analyze data in S3 using serverless SQL, use Athena



Performance Improvement

- Columnar data for cost savings (less scan) to only scan the columns you need
 - Apache Parquet or ORC is recommended
 - Huge performance improvement
 - Use Glue to convert data to Parquet or ORC
- Compress data for smaller retrievals
- Partition datasets in S3 for easy querying on virtual columns
 - s3://bucket//pathToTable for direct access

```
s3://yourBucket/pathToTable
    /<PARTITION_COLUMN_NAME>=<VALUE>
    /<PARTITION_COLUMN_NAME>=<VALUE>
    /<PARTITION_COLUMN_NAME>=<VALUE>
    /etc...
```

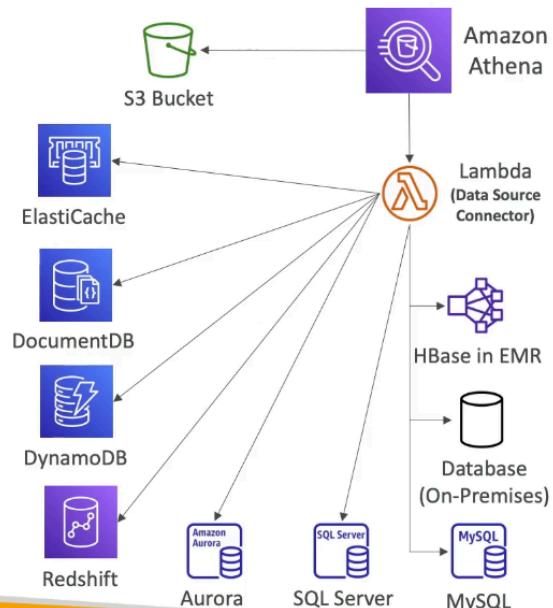
Example: s3://athena-examples/flight/parquet/year=1991/month=1/day=1/

- Use larger files (> 128 MB) to minimize overhead

Federated Query

Amazon Athena – Federated Query

- Allows you to run SQL queries across data stored in relational, non-relational, object, and custom data sources (AWS or on-premises)
- Uses Data Source Connectors that run on AWS Lambda to run Federated Queries (e.g., CloudWatch Logs, DynamoDB, RDS, ...)
- Store the results back in Amazon S3



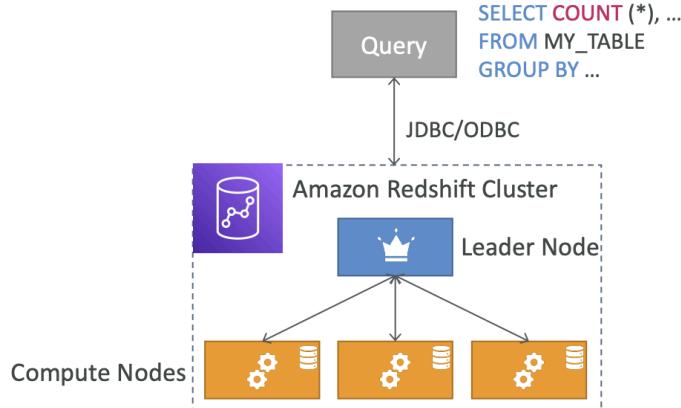
- Allows to run SQL queries across data stored in relational, non relational, object, and custom data sources (AWS or on premise)
- Uses data source connectors that run on Lambda to run Federated Queries in other services (CW logs, DynamoDB, RDS...) and stores results back to S3

Amazon Redshift

- Online analytical processing for analytics and data warehousing
 - SQL interface for queries
 - BI tools like Quicksight or Tableau integrate with it
- Columnar storage of data (instead of row) & parallel query engine
 - Pay as you go based on instances provisioned
- vs Athena: faster queries / joins / aggregations thanks to indexes

Redshift Cluster

Redshift Cluster

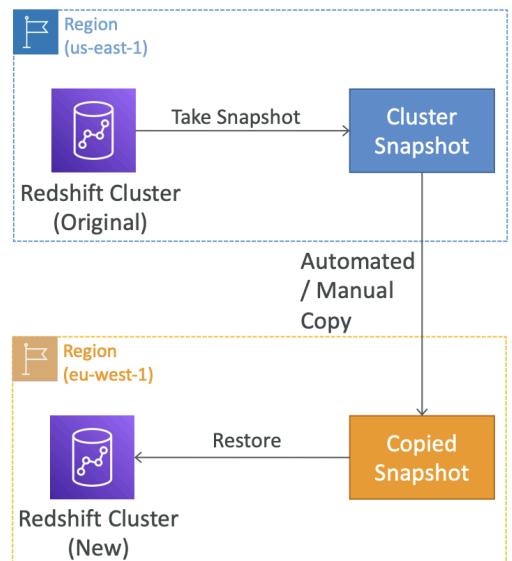


- Leader node: for query planning, results aggregation
- Compute node: for performing the queries, send results to leader
- You provision the node size in advance
- You can use Reserved Instances for cost savings

- Leader node: for query planning, results aggregation
- Compute node: perform queries, send results to leader
- Must provision node size in advance, can use reserved instances for cost savings

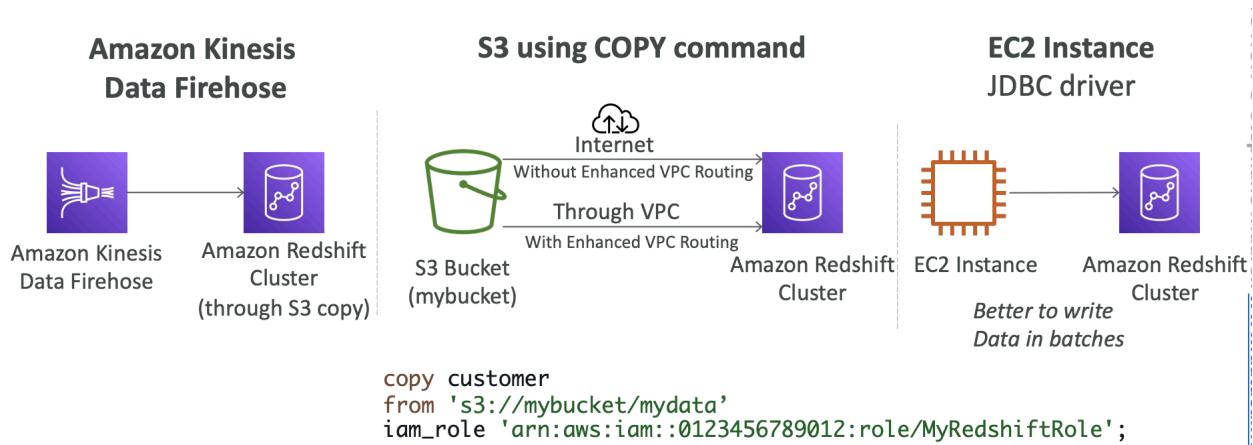
Snapshots & DR

- Redshift has multi AZ for some cluster types
 - Snapshots are point in time backups of a cluster, stored internally in S3
 - Snapshots are incremental (only what has changed is saved)
 - Restore snapshot into new cluster



- Automated: every 8 hours, every 5 GB, or on schedule. Set retention
- Manual: snapshot retained until deleted
- Can configure Redshift to auto copy snapshots of a cluster to another region

Loading Data into Redshift

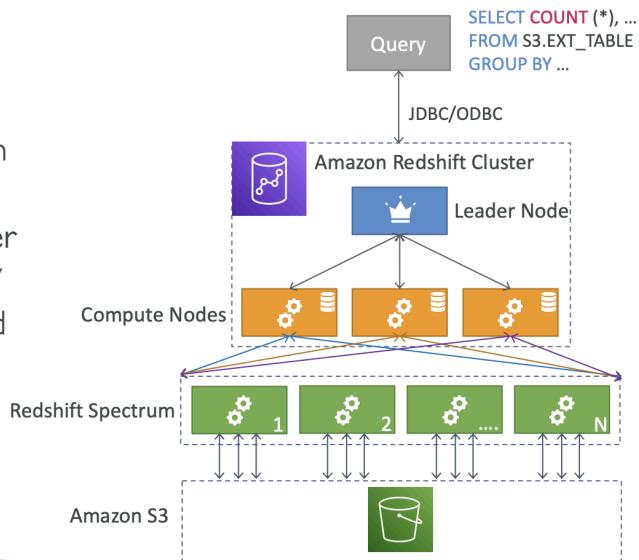


- Large inserts better

Redshift Spectrum

Redshift Spectrum

- Query data that is already in S3 without loading it
- Must have a Redshift cluster available to start the query
- The query is then submitted to thousands of Redshift Spectrum nodes



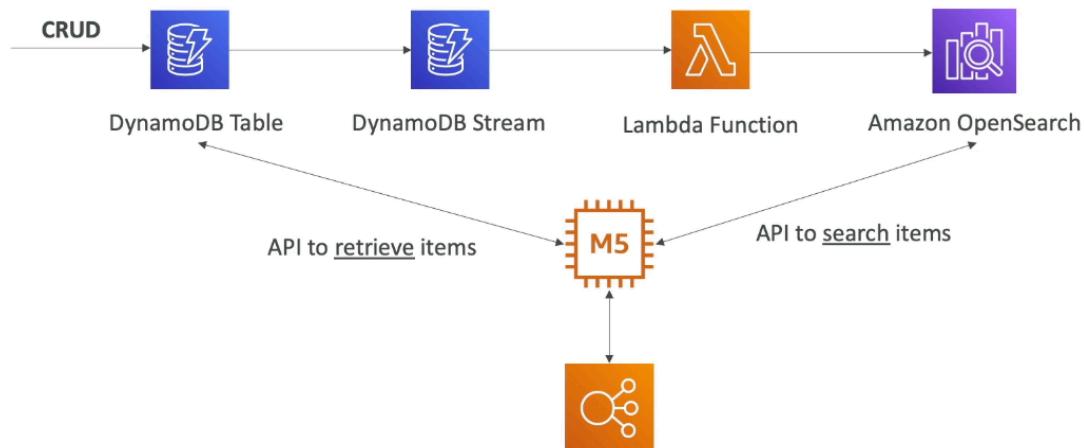
- Query data in S3 without loading
- Must have Redshift cluster available to start query
 - Query is submitted to thousands of Redshift Spectrum nodes

Amazon OpenSearch

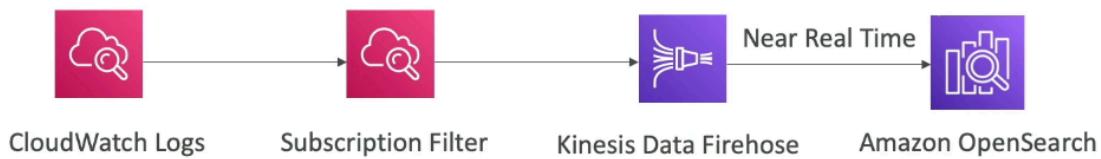
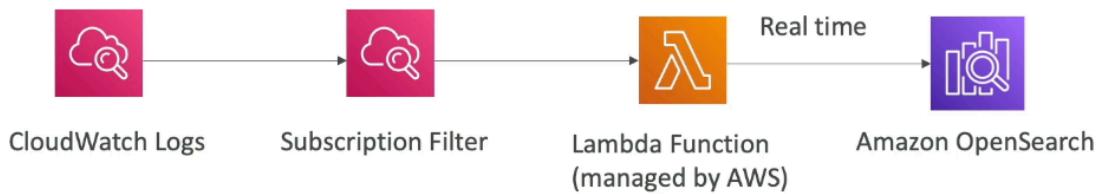
- In DynamoDB, queries only exist by primary key or indexes, but with OpenSearch you can search any field, even partial matches and comes with OpenSearch Dashboards for visualization
 - Common to use OpenSearch as a complement to another DB
 - Ingestion of data from Kinesis Firehose, AWS IoT, CloudWatch logs
- 2 modes: managed cluster or serverless
- No native support for SQL
- Security through Cognito, IAM, KMS, TLS

Architecture Patterns

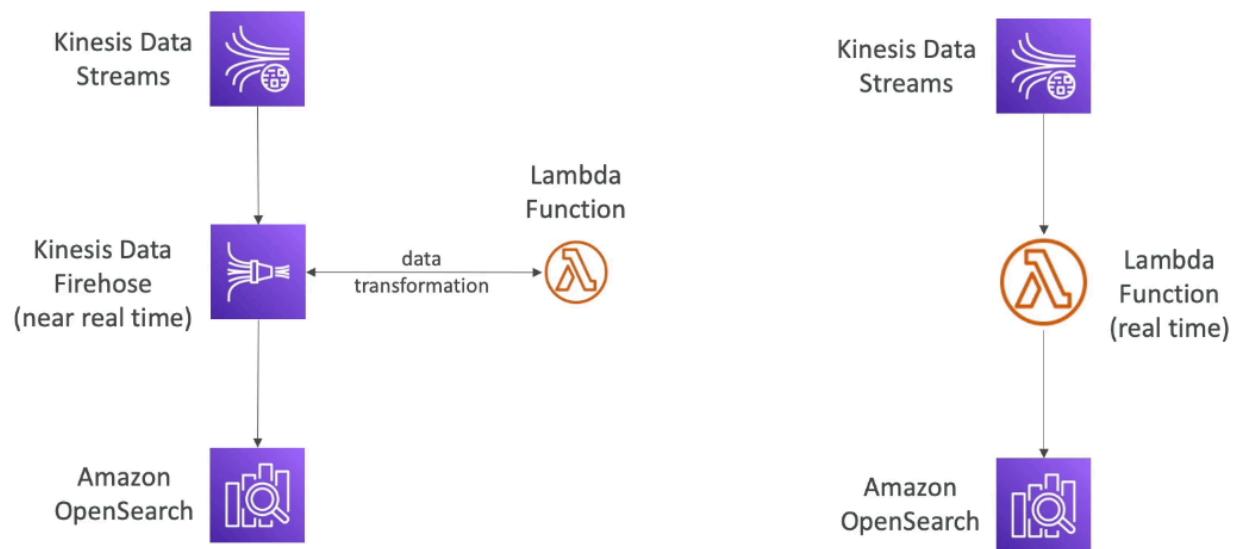
OpenSearch patterns DynamoDB



OpenSearch patterns CloudWatch Logs



OpenSearch patterns Kinesis Data Streams & Kinesis Data Firehose



Amazon EMR

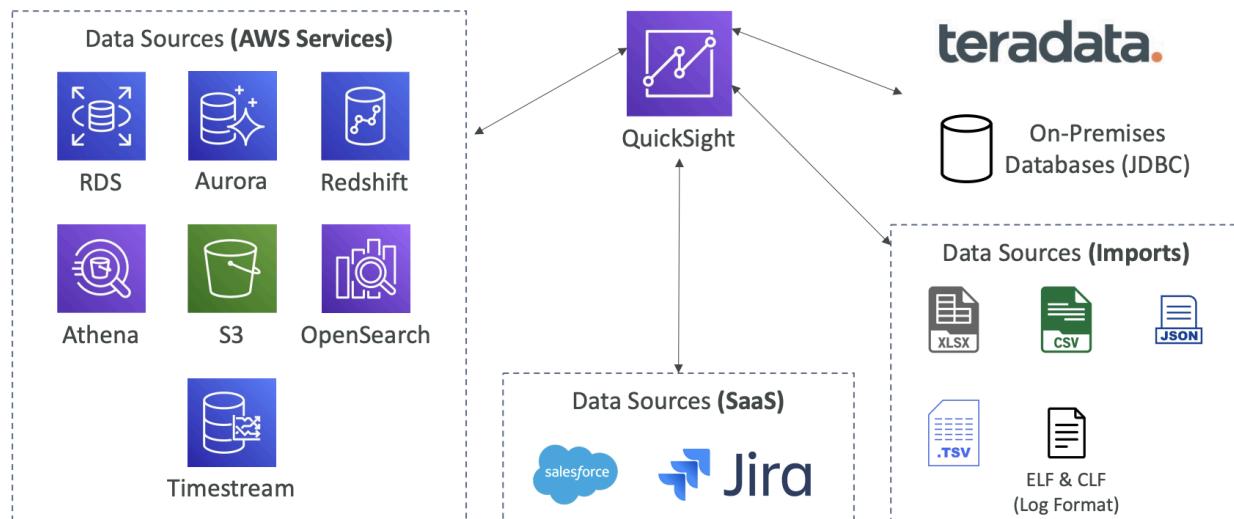
- Create Hadoop clusters (big data) to analyze and process vast amount of data
- Clusters made with hundreds of EC2 instances
 - EMR bundles with Apache Spark, Flink...
- Takes care of all provisioning and configuration, auto scaling and integrated with Spot instances
- Use cases: data processing, ML, web indexing, big data

EMR Node Types & Purchasing

- Master node: manage cluster, coordinate, manage health – long running
- Core node: run task and store data – long running
- Task node (optional): just to run tasks – usually spot
- Purchasing:
 - On demand: reliable, predictable, no termination
 - Reserved (1 year min): cost savings (EMR will automatically use if available)
 - Master and core nodes good options for this
 - Spot Instance: cheaper, can be terminated, less reliable
- Can have long running cluster or transient (temp) cluster deployment modes

Amazon QuickSight

QuickSight Integrations



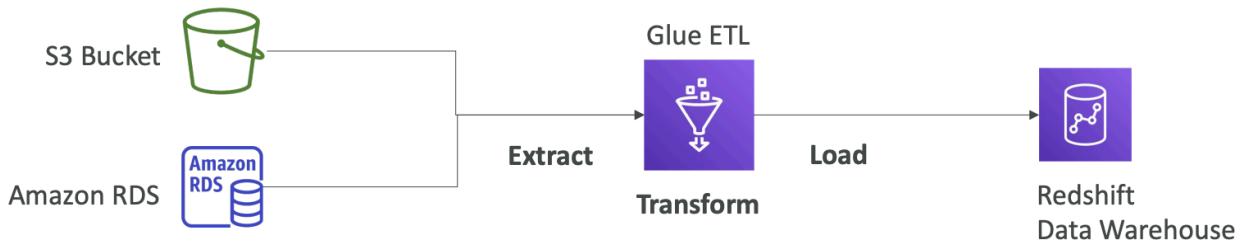
- Serverless ML powered BI to create dashboards
 - Fast, auto scale, embeddable with per session pricing
- Use case: business analytics, visualizations, insights...

- Integrated with RDS, Aurora, Athena, Redshift, S3...
- In memory computation using SPICE engine if data imported into QuickSight
- Enterprise edition:
 - Can set up column level security (CLS) to prevent some columns to be displayed to some users

Dashboards & Analytics

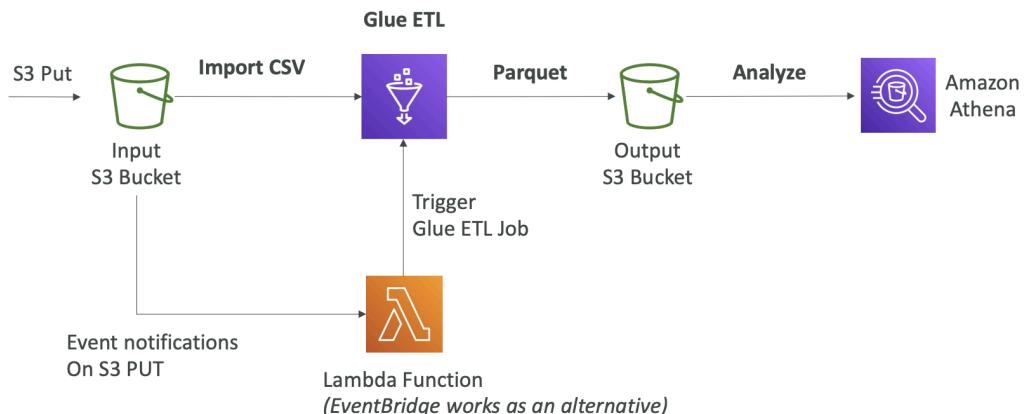
- Define users (standard version) and Groups (enterprise)
 - Users & groups only exist within QuickSight
- Dashboard:
 - Read only snapshot of analysis that can be shared
 - Preserves the configuration of analysis (filtering, parameters, controls...)
- Share analysis or dashboard with users or groups by publishing
 - Users who see the dashboard can also see underlying data

AWS Glue



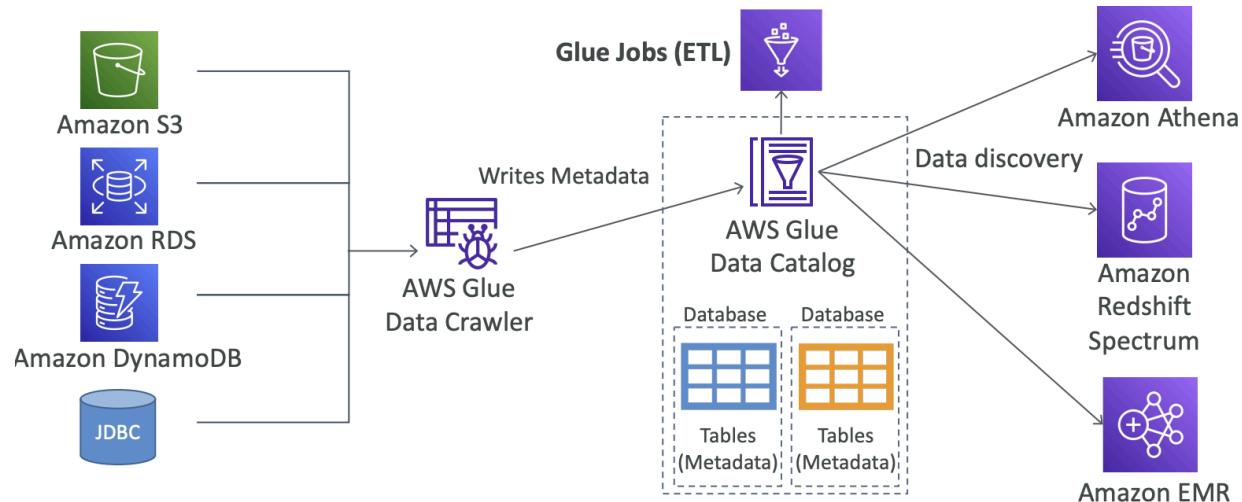
- Fully serverless, managed extract, transform, and load (ETL) service
 - Prepare and transform data for analytics

AWS Glue – Convert data into Parquet format



Glue Data Catalog

Glue Data Catalog: catalog of datasets



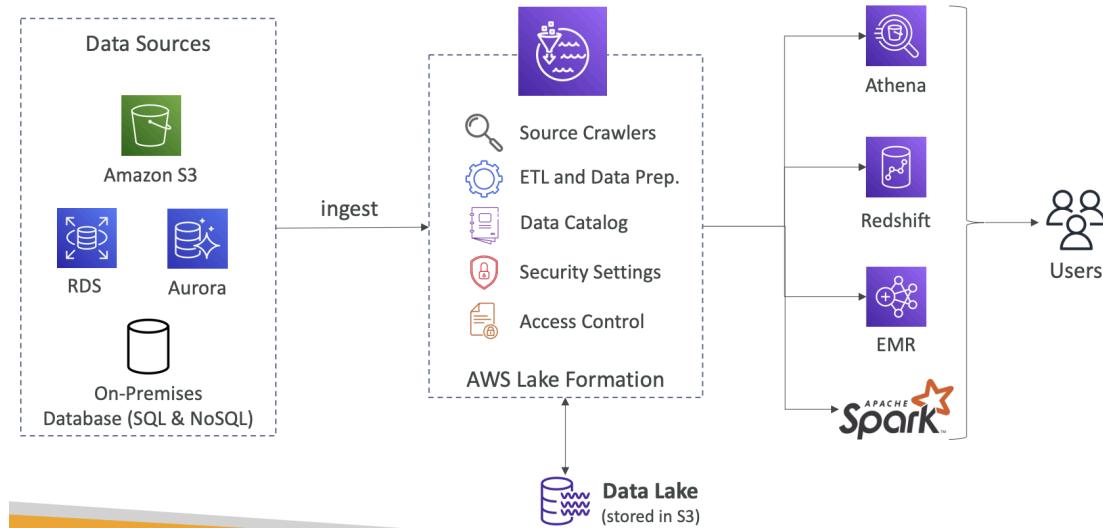
Glue High Level

Glue – things to know at a high-level

- Glue Job Bookmarks: prevent re-processing old data
- Glue Elastic Views:
 - Combine and replicate data across multiple data stores using SQL
 - No custom code, Glue monitors for changes in the source data, serverless
 - Leverages a “virtual table” (materialized view)
- Glue DataBrew: clean and normalize data using pre-built transformation
- Glue Studio: new GUI to create, run and monitor ETL jobs in Glue
- Glue Streaming ETL (built on Apache Spark Structured Streaming): compatible with Kinesis Data Streaming, Kafka, MSK (managed Kafka)

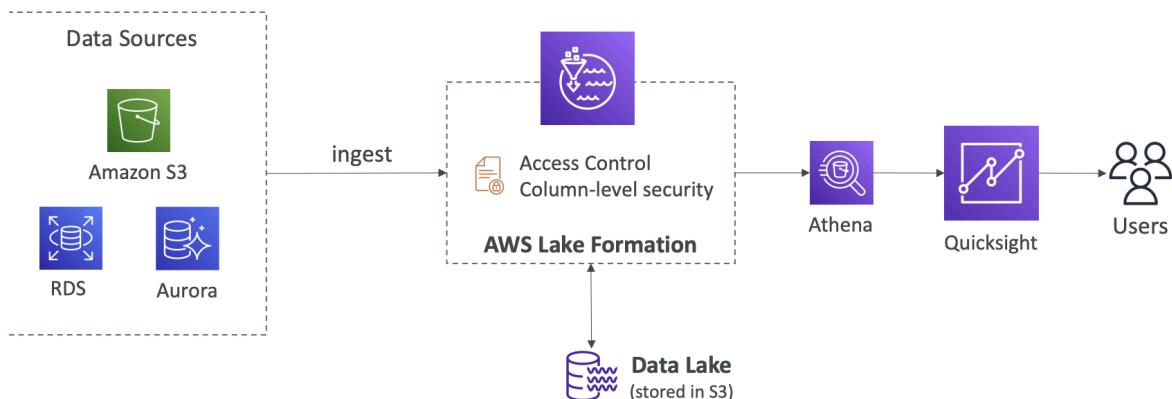
AWS LakeFormation

AWS Lake Formation



- Data lake = central place to have all your data for analytics purposes
- Fully managed service that makes it easy to setup a data lake in days
 - Discover, cleanse, transform, and ingest data into Data lake
 - Automates complex manual steps and deduplicate
- Combine structured and unstructured data in data lake
- Out of box source blueprints: S3, RDS, relational & NoSQL DB...
- Fine grained access control for applications (row and column level)
- Built on top of AWS Glue

AWS Lake Formation Centralized Permissions Example



Section 23: Machine Learning

Rekognition Overview

- Find objects, people, text. Scenes in images and videos using ML
 - Facial analysis and search for user verification
 - Create DB of faces or compare against
- Use case: labeling, content moderation, text / face detection / verification, pathing



Content Moderation

- Detecting unwanted content in social media, online...
 - Flag sensitive content for review in Amazon Augmented AI
- Set minimum confidence threshold for items that are flagged
 - Lower means more content shown

Amazon Transcribe

- Speech to text, using deep learning process called automatic speech recognition (ASR) to convert speech to text
 - Auto removes personally identifiable information using Redaction
 - Supports Automatic Language Detection for multi lingual audio
- Use cases: transcribe, closed captioning, generate metadata for media assets to create a fully searchable archive

Amazon Polly

- Text to speech via deep learning

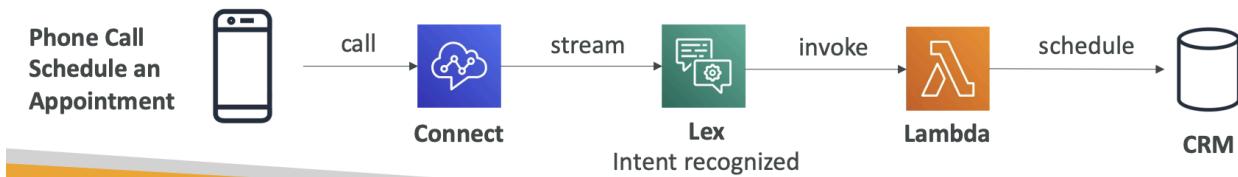
Lexicon & SSML

- Customize pronunciation of words with pronunciation lexicons
 - Upload lexicons and use them in SynthesizeSpeech operation
- Generate speech from plain text or from documents marked up with Speech Synthesis Markup Language (SSML) – more customization
 - Emphasize words or phrases, phonetic pronunciation...

Amazon Translate

- Language translation
 - Can localize content - such as websites and applications - for international users, and to easily translate large volumes of text efficiently

Amazon Lex + Connect



- Lex:
 - Automatic Speech Recognition (ASR) to convert speech to text
 - Natural language understanding
 - Chatbots, call center bots...
- Connect
 - Receive calls, create contact flows, cloud-based virtual contact center
 - Can integrate with other CRM systems or AWS
 - No upfront payments, very cheap

Amazon Comprehend

- NLP, fully managed and serverless
 - ML to find insights and relationships in text
 - Sentiment analysis, extract phrases...
 - Automatically organizes a collection of text files by topic
 - Use cases: sentiment analysis, create and group articles by topic

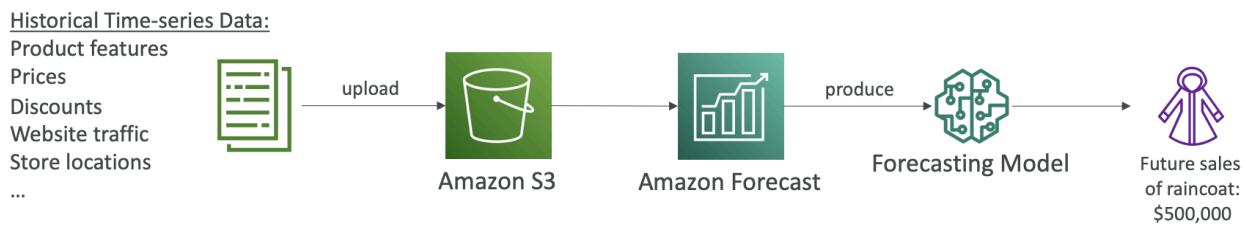
Comprehend Medical

- Detects and returns useful information in unstructured clinical text: doctor notes, test results, etc...
 - Uses NLP to detect protected health info – DetectPHI API
- Store documents in S3, analyze real-time data with Kinesis Data Firehose, or use Amazon Transcribe to transcribe patient narratives into text that can be analyzed by Amazon Comprehend Medical

Amazon SageMaker

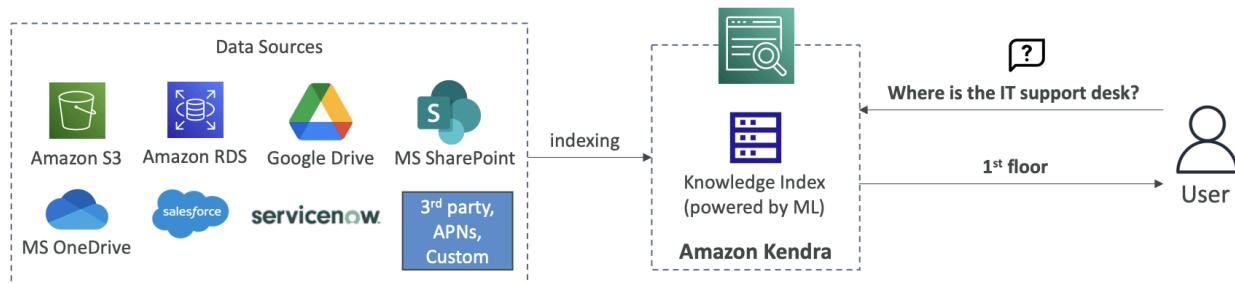
- Fully managed service to build ML models

Amazon Forecast



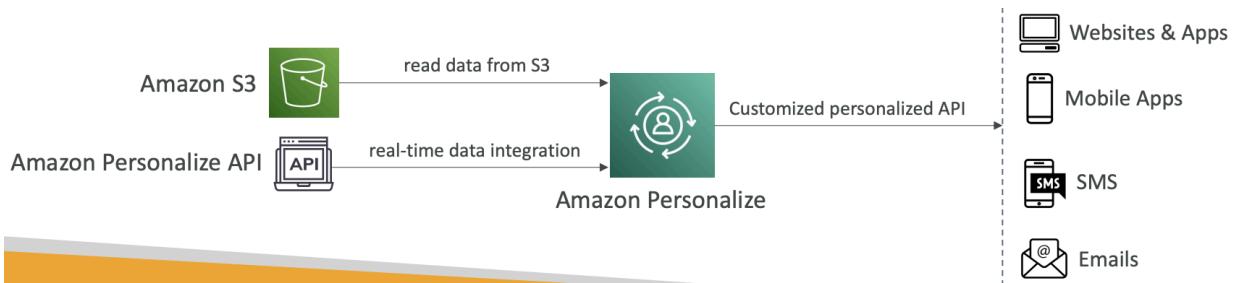
- Fully managed ML to deliver accurate forecasts
 - Reduce forecasting time from months to hours

Amazon Kendra



- Fully managed document search service via ML
 - Extract answers from documents (text, pdf, HTML, PPT, Word...)
 - NLP search capabilities
- Incremental learning: learn from user interactions / feedback to promote preferred results
 - Can manually fine tune search results

Amazon Personalize



- Fully managed ML service to build apps with real time personalized recommendations
 - Example: personalized product recommendations / re-ranking, customized direct marketing

- Integrates into existing websites, applications, SMS, email marketing systems...
 - Implement in days (no need to build, train, deploy ML solutions)

Amazon Textract



- Automatically extract text, handwriting, and data from scanned documents via AI and ML
 - Extra data from forms and tables, read and process any type of documents
 - Use case: financial services, healthcare...

AWS ML Summary

AWS Machine Learning - Summary

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

Section 24: AWS Monitoring & Audit: CloudWatch, CloudTrail & Config

CloudWatch Container Insights

- Collect, aggregate, summarize metrics and logs from containers
 - ECS, EKS, K8 on EC2, Fargate
- In EKS and Kubernetes, CloudWatch Insights uses containerized version of CW Agent to discover containers

CloudWatch Lambda Insights

- Monitoring and troubleshooting for Lambda
 - Collects, aggregates, and summarizes system level metrics including CPU time, memory, disk, and network
 - Collects, aggregates, and summarizes diagnostic information such as cold starts and Lambda worker shutdowns
- Lambda insights provided as lambda layer

CloudWatch Contributor Insights

- Analyze log data and create time series that display contributor data
 - See metrics about top N contributors
 - Total number of unique contributors and usage
- Find top talkers and understand who or what is impacting system performance
- Works for any AWS generated logs
- Build rules from scratch or use sample rules that AWS has created – leverages your CloudWatch Logs
 - CloudWatch also provides built-in rules that you can use to analyze metrics from other AWS services

CloudWatch Application Insights

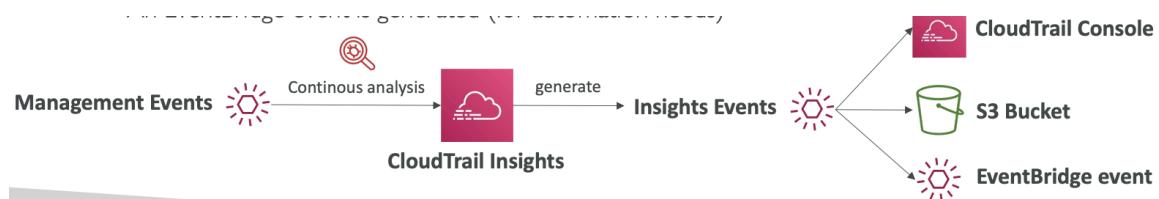
- Provide automated dashboard to show potential problems with monitored applications to help isolate ongoing issues
 - Powered by SageMaker
- Enhanced visibility into your application health to reduce the time it will take you to troubleshoot and repair your applications
- Findings and alerts are sent to EventBridge and SSM OpsCenter

CloudWatch Insights and Operational Visibility Summary

CloudWatch Insights and Operational Visibility

- CloudWatch Container Insights
 - ECS, EKS, Kubernetes on EC2, Fargate, needs agent for Kubernetes
 - Metrics and logs
- CloudWatch Lambda Insights
 - Detailed metrics to troubleshoot serverless applications
- CloudWatch Contributors Insights
 - Find “Top-N” Contributors through CloudWatch Logs
- CloudWatch Application Insights
 - Automatic dashboard to troubleshoot your application and related AWS services

CloudTrail Insights



- Enable to detect unusual activity in account
 - Inaccurate resource provisioning, hitting service limits, bursts of IAM actions...
- Analyze normal management events as a baseline and analyze write events to detect unusual patterns
 - Anomalies appear in CloudTrail console
 - Event sent to S3 and EventBridge event is generated

AWS Config

- Helps with auditing and recording compliance of AWS resources
 - Helps record configuration changes over time
 - Any buckets with public access? Has ALB configuration changed over time?
- Receive alerts (SNS) for any changes
- Per region service, can be aggregated across regions and accounts
- Can store data in S3

Config Rules

AWS Config Resource

- View compliance of a resource over time

sg-077b425b1649da83e	EC2 SecurityGroup	Compliant
sg-0831434f1876c0c74	EC2 SecurityGroup	Noncompliant
sg-09f10ed254d464f30	EC2 SecurityGroup	Compliant

- View configuration of a resource over time

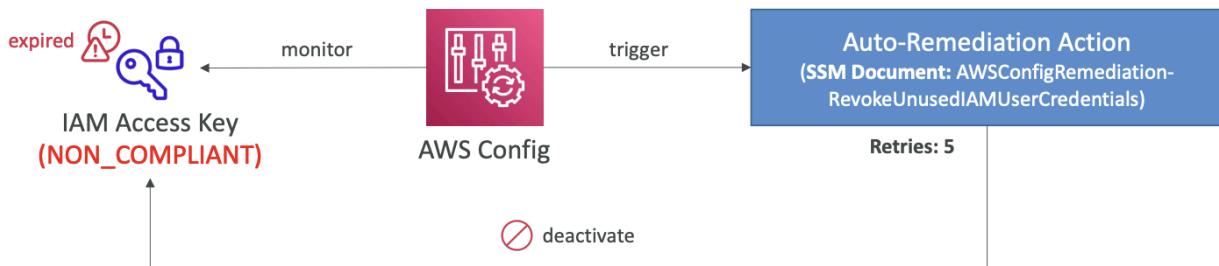


- View CloudTrail API calls of a resource over time



- AWS Managed rules
- Custom rules (defined in Lambda)
- Rules can be evaluated / triggered
 - For each config change
 - And / or at regular time intervals
- Config rules do not prevent actions from happening (no deny)
 - Just gives overview and compliance of resources
- Pricing: no free tier

Remediations



- Automate remediation of non-compliant resources using SSM Automation Documents
 - Use AWS managed Automation Documents or custom
 - Can create custom that invokes Lambda function
 - Remediation retries if resource still non-compliant after auto remediation

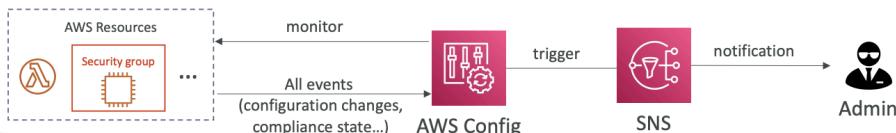
Notifications

Config Rules – Notifications

- Use EventBridge to trigger notifications when AWS resources are non-compliant



- Ability to send configuration changes and compliance state notifications to SNS (all events – use SNS Filtering or filter at client-side)



- Use EventBridge for notifications and can use SNS

CloudWatch vs CloudTrail vs Config

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

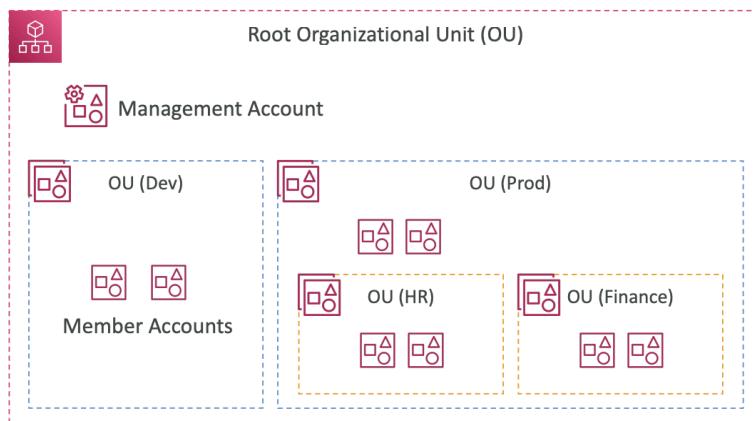
For an Elastic Load Balancer

- CloudWatch:
 - Monitoring Incoming connections metric
 - Visualize error codes as % over time
 - Make a dashboard to get an idea of your load balancer performance
- Config:
 - Track security group rules for the Load Balancer
 - Track configuration changes for the Load Balancer
 - Ensure an SSL certificate is always assigned to the Load Balancer (compliance)
- CloudTrail:
 - Track who made any changes to the Load Balancer with API calls

Section 25: IAM – Advanced

AWS Organizations

AWS Organizations



- Global service to allow management of multiple AWS accounts
 - Main account is management account, other accounts are member accounts
 - Members can be part of only 1 organization
- Consolidated Billing across all accounts – single payment with pricing benefits from aggregated usage
- Shared reserve instances and savings plans discounts across accounts

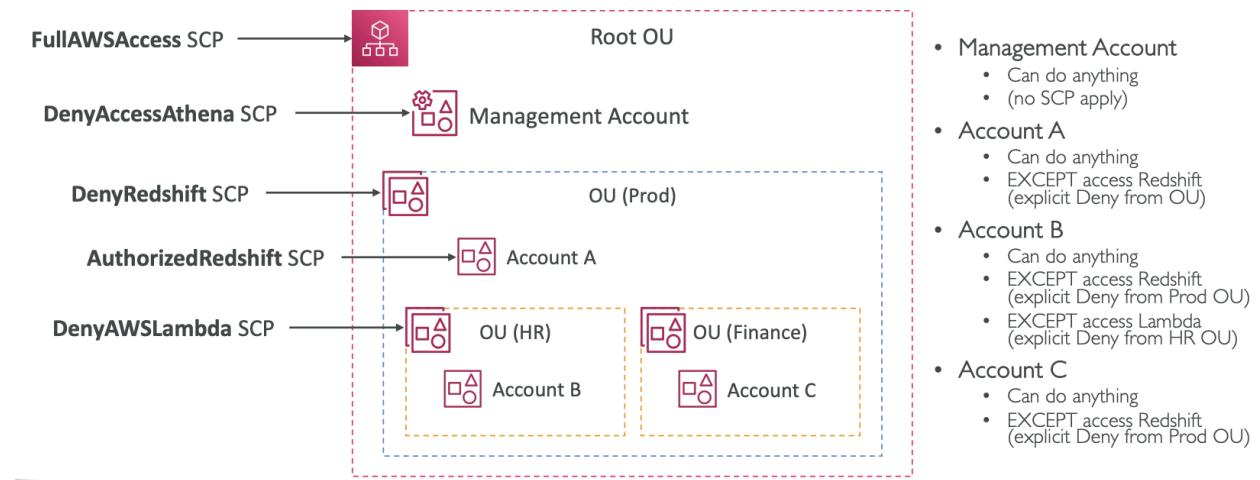
- API available to automate account creation

Advantages

- Multi account vs one account multi VPC
- Tagging standards for billing
- Enable CloudTrail on all accounts, send logs to central S3 account
- Send CloudWatch logs to central logging account
- Establish Cross Account Roles for admin purposes

Security

SCP Hierarchy



SCP Examples

Blocklist and Allowlist strategies

```
[{"Version": "2012-10-17", "Statement": [{"Sid": "AllowsAllActions", "Effect": "Allow", "Action": "*", "Resource": "*"}, {"Sid": "DenyDynamoDB", "Effect": "Deny", "Action": "dynamodb:*", "Resource": "*"}]}, {"Version": "2012-10-17", "Statement": [{"Effect": "Allow", "Action": ["ec2:*", "cloudwatch:*"], "Resource": "*"}]}]
```

- Service Control Policies (SCP)
- IAM policies applied to OU or accounts to restrict users and roles
 - Do not apply to management account
- Must have explicit allow (deny all by default)

IAM Conditions

IAM Conditions

aws:SourceIp

restrict the client IP from
which the API calls are being made

```
{ "Version": "2012-10-17", "Statement": [ { "Effect": "Deny", "Action": "*", "Resource": "*", "Condition": { "NotIpAddress": { "aws:SourceIp": ["192.0.2.0/24", "203.0.113.0/24"] } } } ] }
```

aws:RequestedRegion

restrict the region the
API calls are made to

```
{ "Version": "2012-10-17", "Statement": [ { "Effect": "Deny", "Action": ["ec2:", "rds:", "dynamodb:"], "Resource": "*", "Condition": { "StringEquals": { "aws:RequestedRegion": ["eu-central-1", "eu-west-1"] } } } ] }
```

IAM Conditions

<p><u>ec2:ResourceTag</u> restrict based on tags</p> <pre>{ "Version": "2012-10-17", "Statement": [{ "Effect": "Allow", "Action": ["ec2:StartInstances", "ec2:StopInstances"], "Resource": "arn:aws:ec2:us-east-1:123456789012:instance/*", "Condition": { "StringEquals": { "ec2:ResourceTag/Project": "DataAnalytics", "aws:PrincipalTag/Department": "Data" } } }] }</pre>	<p><u>aws:MultiFactorAuthPresent</u> to force MFA</p> <pre>{ "Version": "2012-10-17", "Statement": [{ "Effect": "Allow", "Action": "ec2:*", "Resource": "*" }, { "Effect": "Deny", "Action": ["ec2:StopInstances", "ec2:TerminateInstances"], "Resource": "*", "Condition": { "BoolIfExists": { "aws:MultiFactorAuthPresent": false } } }] }</pre>
---	---

IAM for S3

- s3>ListBucket permission applies to
`arn:aws:s3:::test`
- => bucket level permission

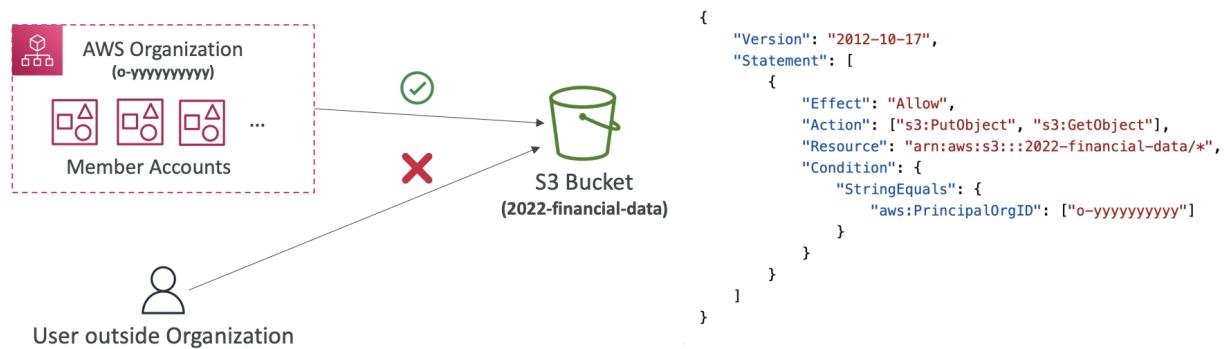
- s3GetObject, s3PutObject,
s3DeleteObject applies to
`arn:aws:s3:::test/*`
- => object level permission

```
{  
    "Version": "2012-10-17",  
    "Statement": [  
        {  
            "Effect": "Allow",  
            "Action": ["s3>ListBucket"],  
            "Resource": "arn:aws:s3:::test"  
        },  
        {  
            "Effect": "Allow",  
            "Action": [  
                "s3:PutObject",  
                "s3:GetObject",  
                "s3>DeleteObject"  
            ],  
            "Resource": "arn:aws:s3:::test/*"  
        }  
    ]  
}
```

Resource Policies & aws:PrincipalOrgID

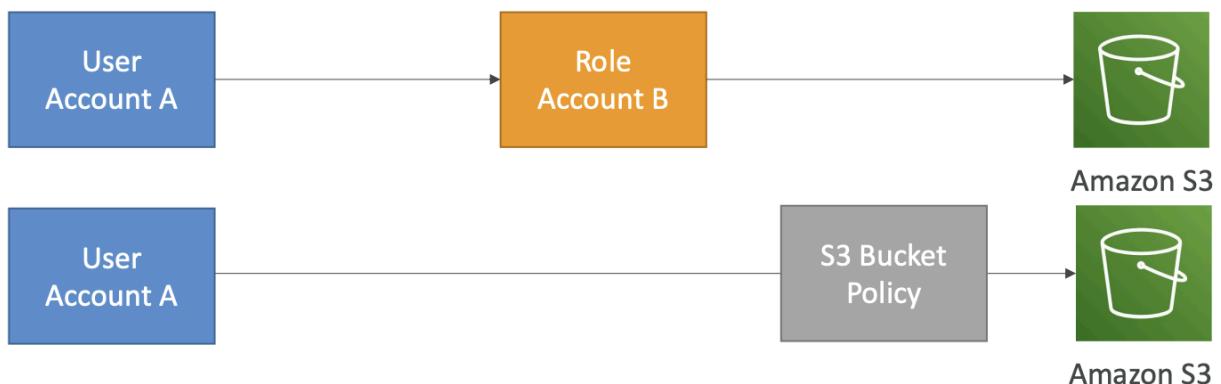
Resource Policies & aws:PrincipalOrgID

- aws:PrincipalOrgID can be used in any resource policies to restrict access to accounts that are member of an AWS Organization



- aws:PrincipalOrgID can be used in any resource policies to restrict access to accounts that are member of AWS org

IAM Roles vs Resource Based Policies

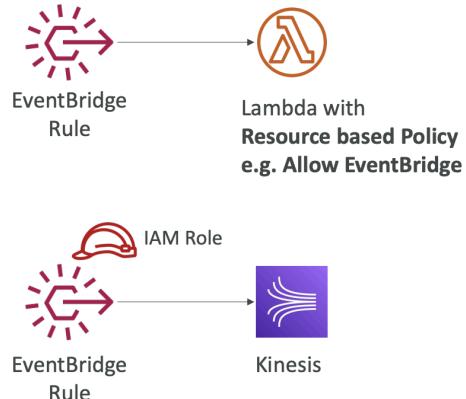


- Cross account:
 - Attach resource based policy to resource or use a role as proxy
- When you assume a role (user, application or service), you give up original permissions and take permissions assigned to the role
- When use a resource based policy, the principal does not give up permissions
 - Example: User in account A needs to scan a DynamoDB table in Account A and dump it in an S3 bucket in Account B

EventBridge – Security

Amazon EventBridge – Security

- When a rule runs, it needs permissions on the target
- Resource-based policy: Lambda, SNS, SQS, S3 buckets, API Gateway...
- IAM role: Kinesis stream, Systems Manager Run Command, ECS task...



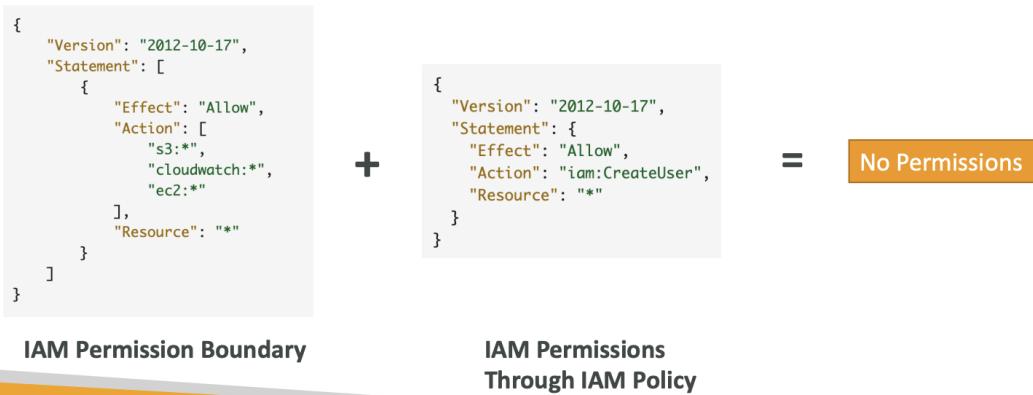
- When a rule runs, it needs permissions on target
- Resource based policy: lambda, S3, API GW...
- IAM Role: Kinesis, ECS task...

IAM Permission Boundaries

IAM Permission Boundaries

- IAM Permission Boundaries are supported for users and roles (not groups)
- Advanced feature to use a managed policy to set the maximum permissions an IAM entity can get.

Example:

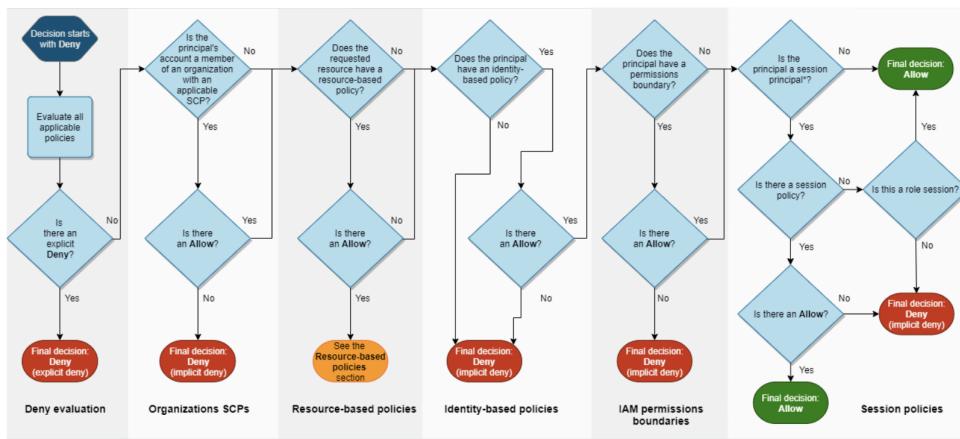


- Supported for users and roles (not groups)

- Can be used in combinations of AWS Org SCP
- Advanced feature to use a managed policy to set the maximum permissions an IAM entity can get
 - In example, user cannot create user because the max permissions are in the permissions boundary
- Use cases:
 - Delegate responsibilities to non administrators within their permission boundaries, for example create new IAM users
 - Allow developers to self-assign policies and manage their own permissions, while making sure they can't "escalate" their privileges (= make themselves admin)
 - Useful to restrict one specific user (instead of a whole account using Organizations & SCP)

IAM Policy Evaluation Logic

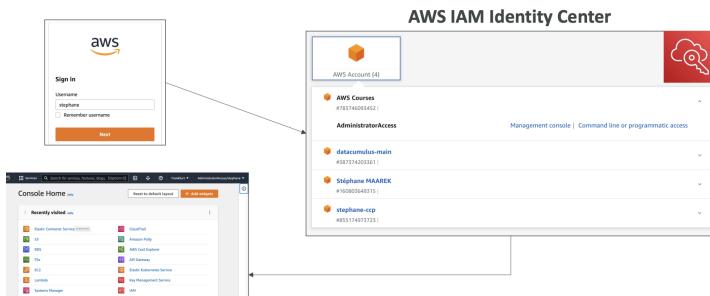
IAM Policy Evaluation Logic



*A session principal is either a role session or an IAM federated user session.

AWS IAM Identity Center

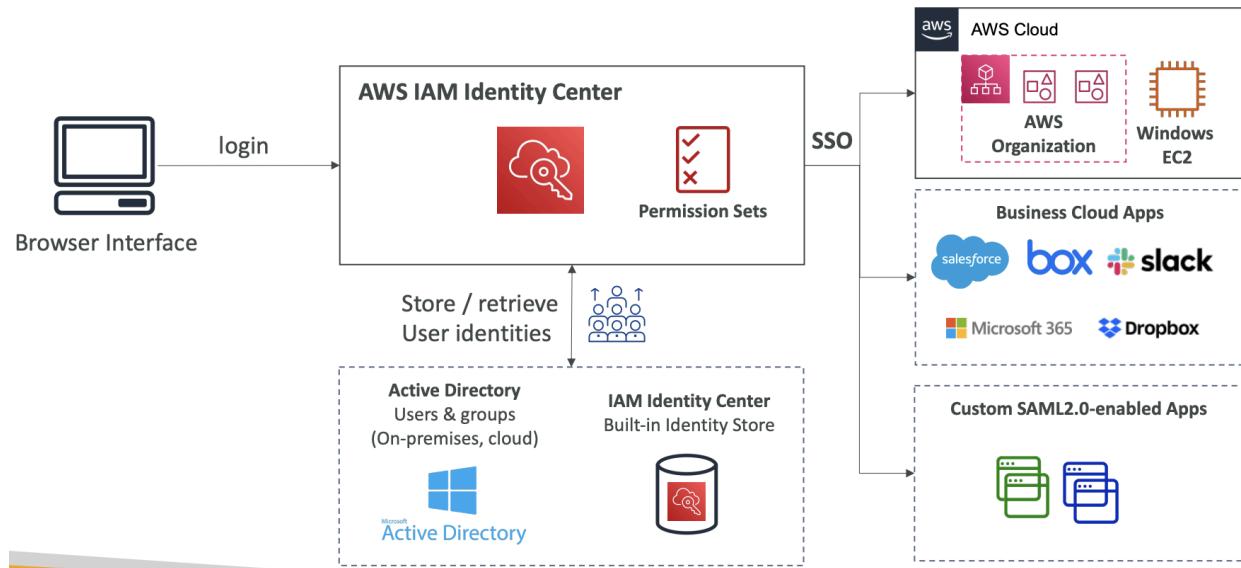
AWS IAM Identity Center – Login Flow



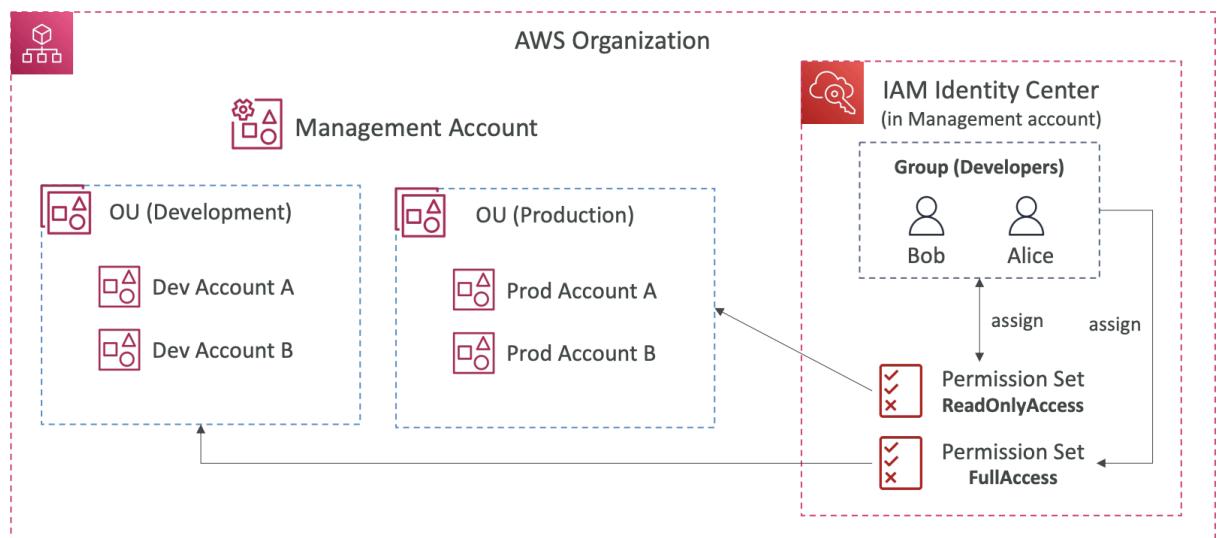
- 1 login SSO for:

- AWS accounts in AWS Orgs
- Business cloud applications
- SAML 2.0 applications
- EC2 Windows instances
- Identity Provider
 - Built in identity store in IAM Identity Center
 - 3rd party: Active Directory, Okta...

AWS IAM Identity Center



IAM Identity Center



Fine Grained Permissions and Assignments

AWS IAM Identity Center Fine-grained Permissions and Assignments



- Multi-Account Permissions

- Manage access across AWS accounts in your AWS Organization
- Permission Sets – a collection of one or more IAM Policies assigned to users and groups to define AWS access

- Application Assignments

- SSO access to many SAML 2.0 business applications (Salesforce, Box, Microsoft 365, ...)
- Provide required URLs, certificates, and metadata

- Attribute-Based Access Control (ABAC)

- Fine-grained permissions based on users' attributes stored in IAM Identity Center Identity Store
- Example: cost center, title, locale, ...
- Use case: Define permissions once, then modify AWS access by changing the attributes

- Multi Account Permissions

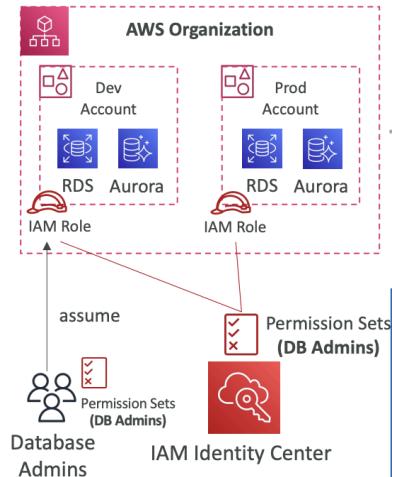
- Managed access across AWS accounts in AWS Org
- Permission sets – collection of 1+ IAM policies assigned to users and groups to define AWS access
 - Automatically create IAM role for users

- Application Assignments

- SSO access to many SAML 2.0 business apps
- Provide required URLs, certificates, metadata

- Attribute based access control (ABAC)

- Fine grained permissions based on users' attributes stored in IAM Identity Center Identity Store
 - Ex: cost center, title...
- Use case: define permissions once, then modify AWS access by changing the attributes



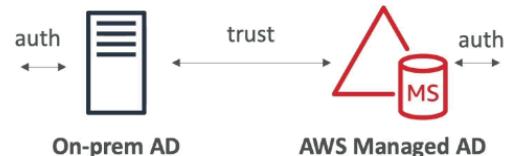
AWS Directory Services

Microsoft Active Directory (AD)

- DB of objects: user accounts, computers, printers, etc... with centralized security management
 - Objects are organized in trees and a group of trees is a forest

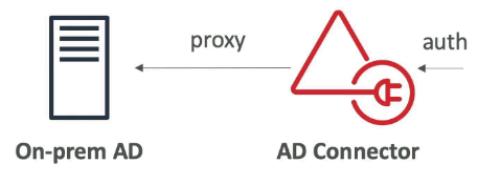
AWS Managed Microsoft AD

- Create own AD in AWS, manage users locally, supports MFA
- Establish trust connections with on premise AD



AD Connector

- Directory Gateway (proxy) to redirect to on premise AD, supports MFA
 - Users are managed on on premise AD



Simple AD

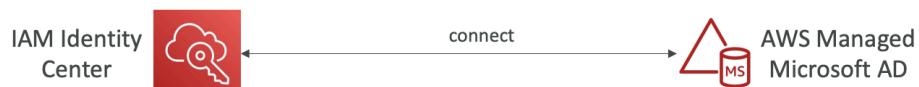
- AD compatible managed directory on AWS, cannot be joined with on premise AD



Active Directory Setup

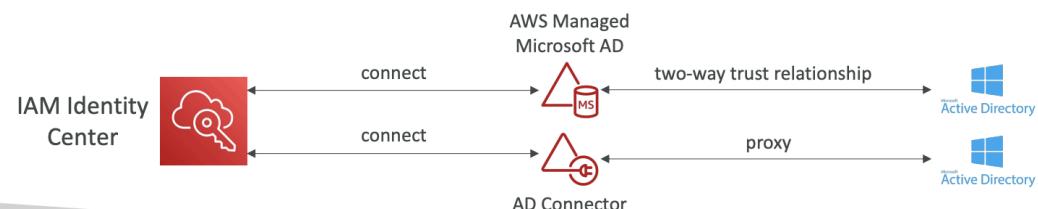
IAM Identity Center – Active Directory Setup

- Connect to an AWS Managed Microsoft AD (Directory Service)
 - Integration is out of the box



- Connect to a Self-Managed Directory

- Create Two-way Trust Relationship using AWS Managed Microsoft AD
- Create an AD Connector

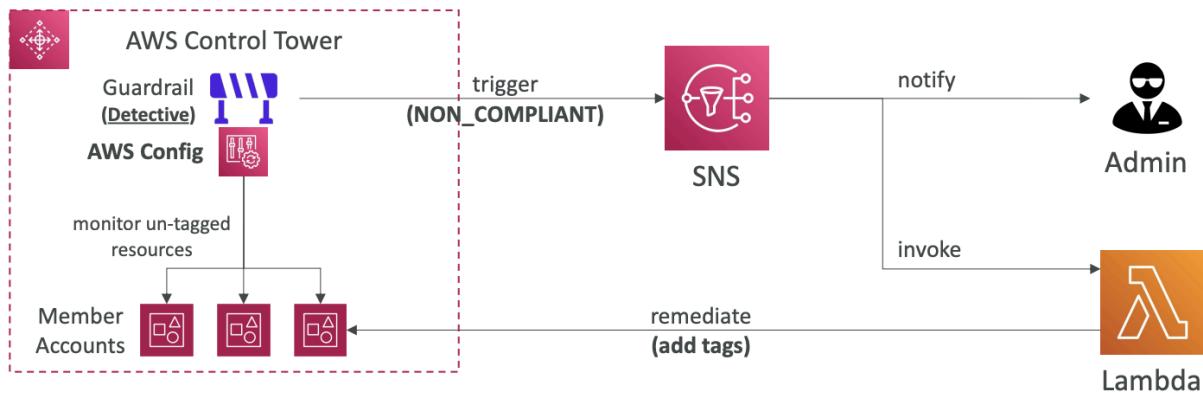


- Connect to AWS Managed Microsoft AD
 - Integration out of the box
- Connect to self managed directory
 - Create 2 way trust relationship using AWS Managed Microsoft AD
 - Create AD connector

AWS Control Tower

- Easy way to set up and govern a secure and compliant multi account AWS environment
 - Use AWS Organizations to create accounts
- Benefits:
 - Automate environment set up
 - Automate ongoing policy management using guardrails
 - Detect policy violations and remediate
 - Monitor compliance through interactive dashboard

Guardrails

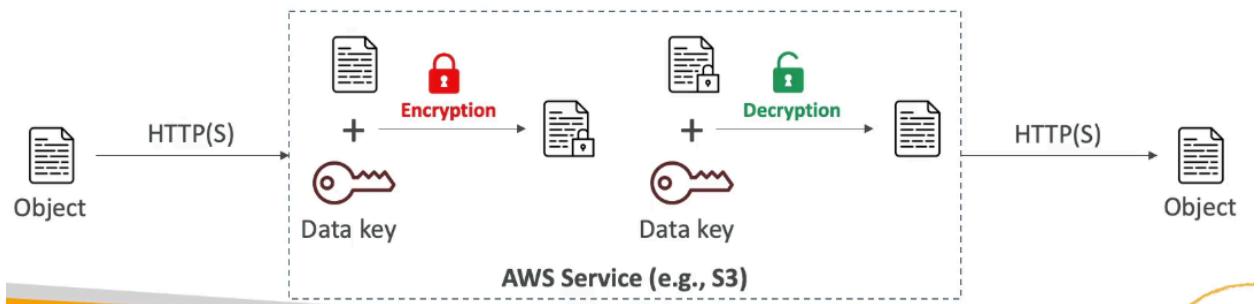


- Provides ongoing governance for your Control Tower environment (AWS Accounts)
- Preventive Guardrail – using SCPs
 - Ex: restrict regions across all accounts
- Detective Guardrail – using AWS Config
 - Ex: identify untagged resources

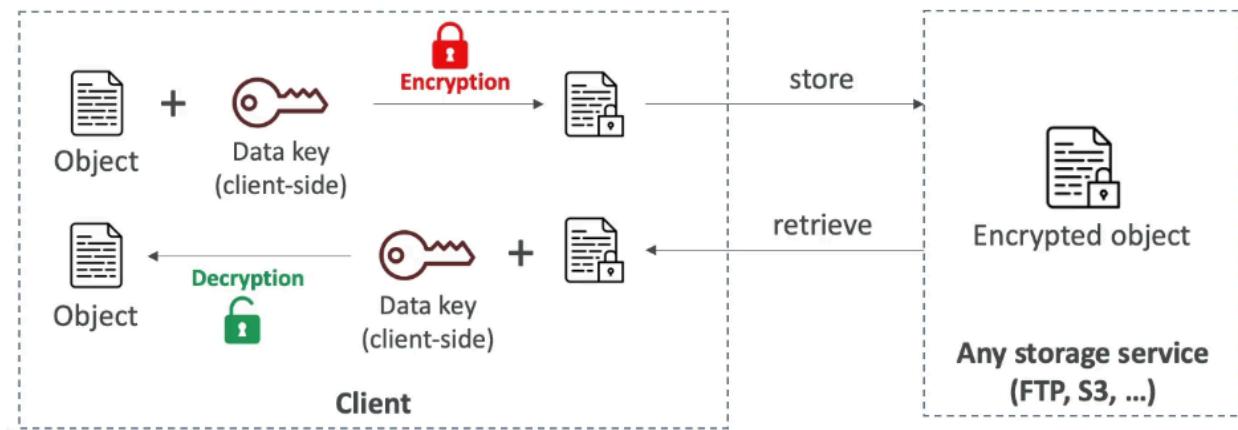
Section 26: AWS Security & Encryption: KMS, SSM Parameter Store, Shield, WAF

Encryption 101

- Encryption in flight (TLS / SSL for HTTPS encryption)
 - Data encrypted before sending and decrypted after receiving → only target server can receive it; ensures no middle man attack occurs



- Server Side encryption at rest
 - Data is encrypted after received by the server and decrypted before being sent
 - Stored in encrypted form via (data) key where the encryption / decryption keys must be managed somewhere and the server needs access to it



- Client side encryption
 - Data is encrypted by the client and never decrypted by the server (server not trusted and cannot decrypt data)
 - Data will be decrypted by a receiving client
 - Could leverage envelope encryption

KMS

- AWS managed encryption keys, fully integrated with IAM for authorization and easy ways to control access to data; can audit with CloudTrail and integration with AWS services
 - Scoped per region, the same KMS key cannot be in the same region
- KMS key types:
 - AWS owned → default free key for SSE-S3, SSE-SQS, SSE-DDB
 - AWS managed key → free for aws/service name
 - Customer managed key created in KMS → \$1 /month
 - Customer managed imported → \$1 /month
 - + pay for API call to KMS

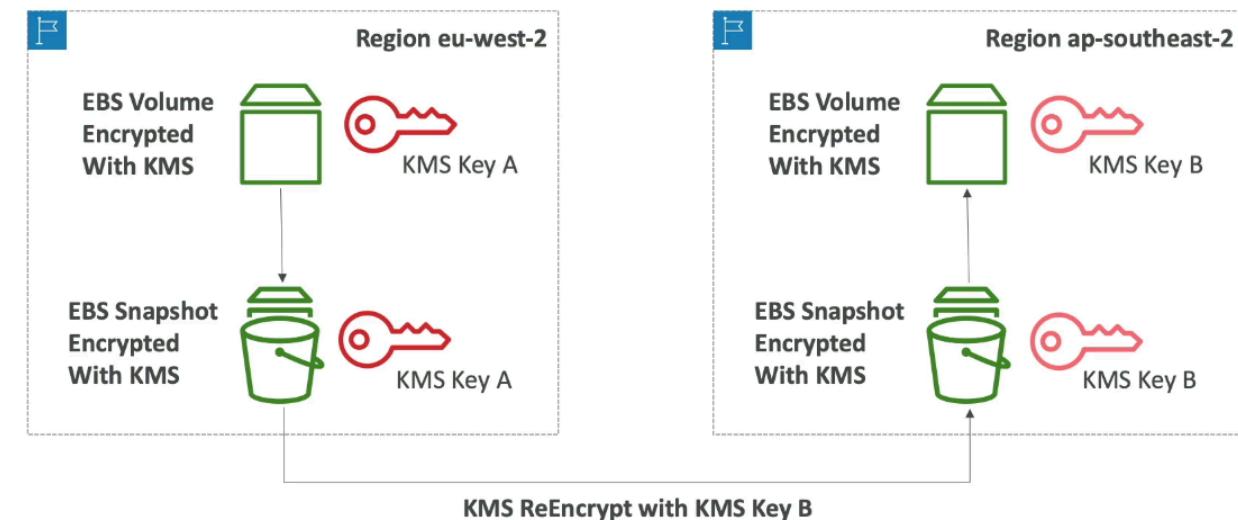
- Automatic Key rotation
 - AWS managed: auto rotate every 1 year
 - Customer managed KMS: enabled feature for automatic rotation 1 year
 - Imported: only manual rotation using alias

Key Types

- Symmetric (AES-256)
 - Single encryption key for encryption / decryption that all AWS services integrate with
 - Never get access to KMS key unencrypted
- Asymmetric (RSA & ECC Key Pairs)
 - Public (encrypt) and private (decrypt) key pair for sign / verify
 - Public key is downloadable, but can't access private key unencrypted
 - Used for encryption outside of AWS by users who can't use KMS

Copying Snapshots across Regions

Copying Snapshots across regions



Key Policies

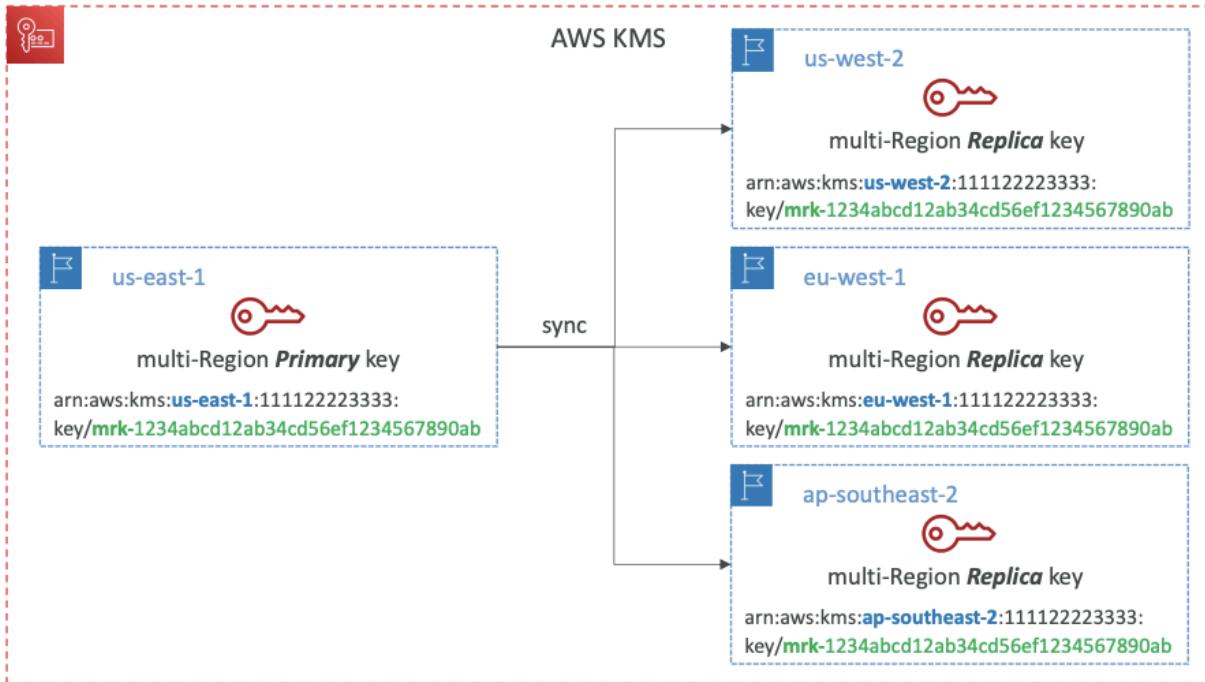
- Control access to KMS keys similar to S3 bucket policy
 - Difference: cannot control access without them
 - If there is no policy, no one can access the KMS key
- Default KMS policy:
 - Created if you don't provide a key policy
 - Complete access to the key to root user = entire AWS account

- Custom Key Policy:
 - Define users, roles access to key and define who can administer the key
 - Useful for cross account access of KMS key

Copying Snapshots across Accounts

1. Create snapshot, encrypted with own KMS key (customer managed)
2. Attach KMS key policy to authorize cross account access
3. Share encrypted snapshot
4. (in target) create a copy of the snapshot, encrypt it with a different customer managed key in account
5. Create a volume from snapshot

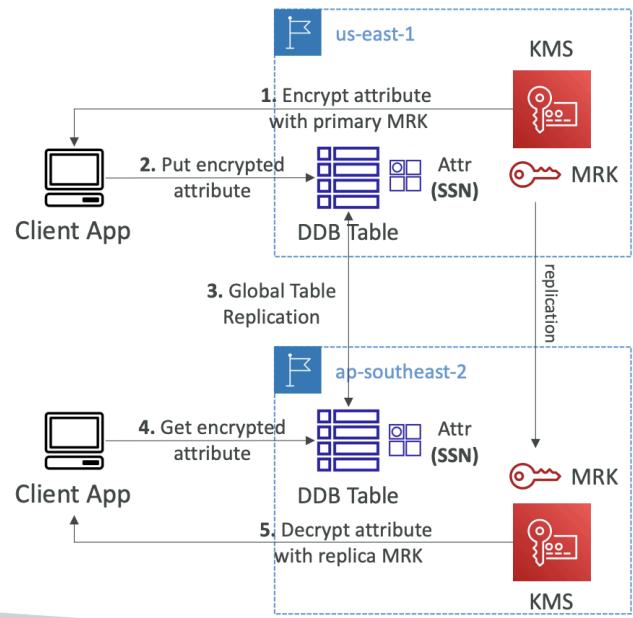
Multi Region Keys



- Identical KMS keys in different regions that can be used interchangeably
- Multi region keys have same key ID, key material, automatic rotation
 - Encrypt in 1 region and decrypt in another
 - No need to re-encrypt or make cross region API calls
- KMS multi region are NOT global (primary + replicas)
 - Each key is managed independently
- Use case: global client side encryption, encryption on global DynamoDB, global Aurora

DynamoDB Global Tables and KMS Multi-Region Keys Client side encryption

- Can encrypt specific attributes client side in DynamoDB using [Amazon DynamoDB Encryption Client](#)
 - Combined with Global Tables, the client-side encrypted data is replicated to other regions
 - If we use a multi-region key, replicated in the same region as the DynamoDB Global table, then clients in these regions can use low latency API calls to KMS in their region to decrypt the data client-side
 - Using client-side encryption we can protect specific fields and guarantee only decryption if the client has access to an API key



Global Aurora and KMS Multi-Region Keys Client Side Encryption

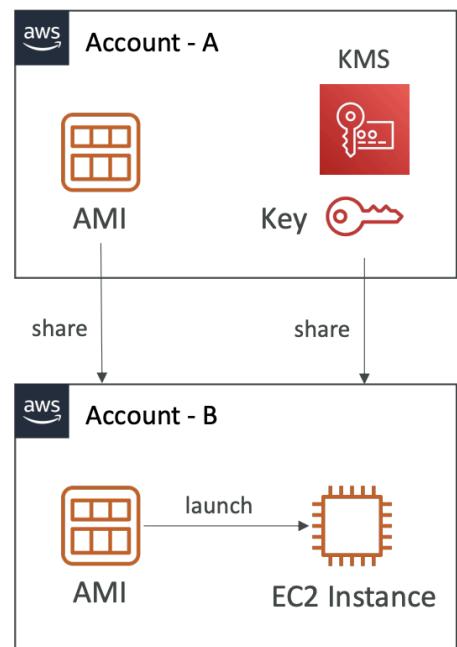
- Encrypt specific attributes client side with [AWS Encryption SDK](#)
 - Combine with Aurora Global to have client side encrypted data replicated to other regions
 - If using multi region key, replicated in the same region as Global Aurora DB, then clients can use API calls to KMS in their region to decrypt data client side

S3 Replication Encryption Considerations

- Unencrypted objects and object encrypted with SSE-S3 are replicated by default
- Objects with SSE-C can be replicated
- With SSE-KMS, must enable option
 - Specify which KMS key
 - Adapt KMS key policy for target key
 - IAM role with kms:Decrypt for source KMS key and kms:Encrypt for the target KMS key
 - May get KMS throttling errors
- Can use multi region KMS key, but treated as independent keys by S3
 - Object will still be decrypted and encrypted by the same key even though the key is multi region

AMI Sharing Process Encrypted via KMS

- AMI in source account encrypted with KMS key from source account
- Must modify the image attribute to add launch permission which corresponds to the specified target AWS account
- Share KMS keys used to encrypt snapshot that the AMI references with the target account / IAM role
- IAM role / user in target account must have permissions to DescribeKey, ReEncrypted, CreateGrant, Decrypt
- When launching EC2 instance from AMI, optionally the target account can re-encrypt AMI with a key of its own

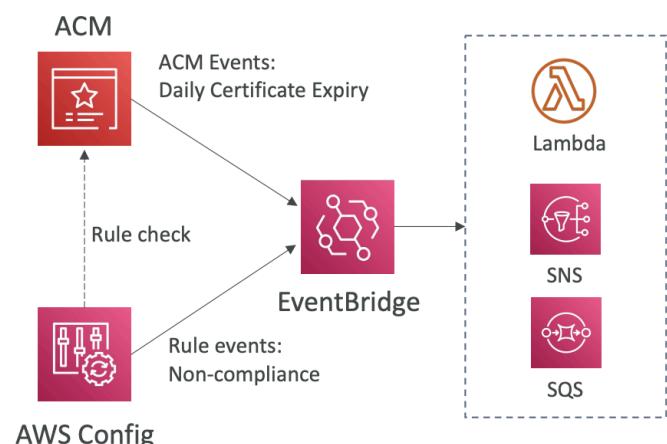


Requesting Public Certificates

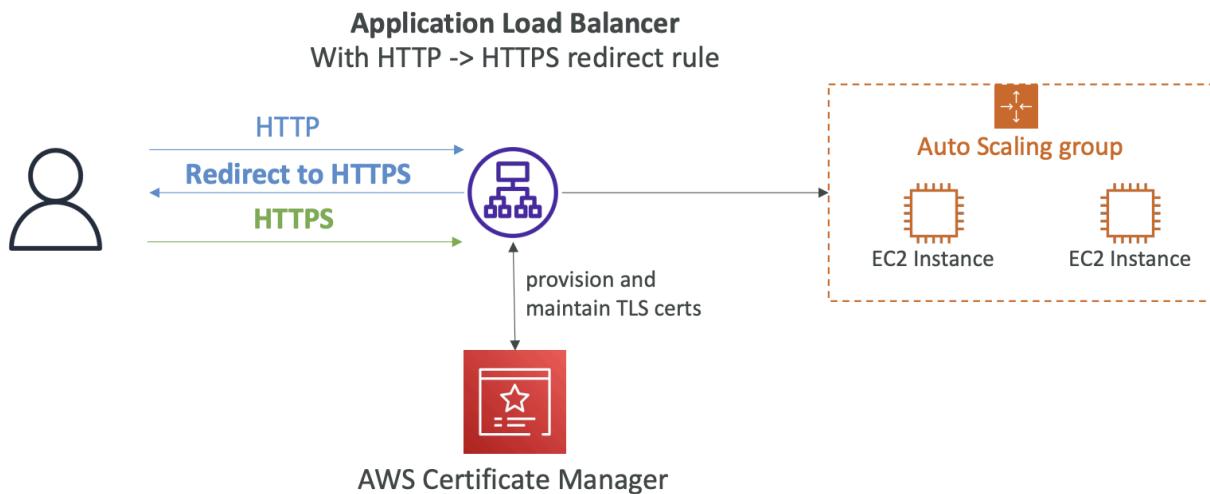
1. List domain names to be included in certificate
 - a. Fully qualified domain name (FQDN)
 - b. Wildcard domain
2. Select validation method: DNS validation, Email validation
 - a. DNS validation preferred for automation purposes
 - b. Email validation will send emails to contact addresses in the WHOIS database
 - c. DNS validation will leverage a CNAME record to DNS config
3. Takes a few hours to get verified and public certificate will be enrolled for automatic renewal
 - a. ACM auto renews ACM generated certificates 60 days before expire

Importing Public Certificates

- Option to generate certificate outside of ACM and import
- No automatic renewal, must import a new certificate before expiry
- ACM sends daily expiration events starting 45 days prior to expiration
 - # of days can be configured
 - Events appear in EventBridge
- AWS Config has managed rule named acm-certificate-expiration-check to check for expiring certificates (configurable days)



Integration with ALB



AWS Shield

- Protect from DDoS (many requests at same time)
- AWS Shield Standard
 - Free, protection from attacks like SYN / UDP Floods, Reflection attacks and other layer 3-4 attacks
- AWS Shield Advanced
 - DDoS mitigation service
 - Protect against more sophisticated attack on EC2, ELB, CloudFront, Global Accelerator, Route 53
 - 24/7 access to AWS DDoS response team
 - Protect against higher fees during usage spikes due to DDoS
 - Automatic application layer DDoS mitigation automatically created and deploys WAF rules to mitigate layer 7 attacks

AWS Firewall Manager

- Manage firewall rules in all accounts of AWS Org
- Security policy: common set of security rules
 - WAF rules (ALB, API GW, CloudFront)
 - Shield Advanced rules (ALB, CLB, NLB, Elastic IP, CloudFront)
 - SG for EC2, ALB, ENI in VPC
 - AWS Network Firewall (VPC level)
 - Route 53 resolver DNS firewall
- Policies created at region level

- Rules applied to new resources as they are created across all and future accounts in organization

WAF vs Firewall Manager vs Shield

WAF vs. Firewall Manager vs. Shield



AWS WAF



AWS Firewall Manager



AWS Shield

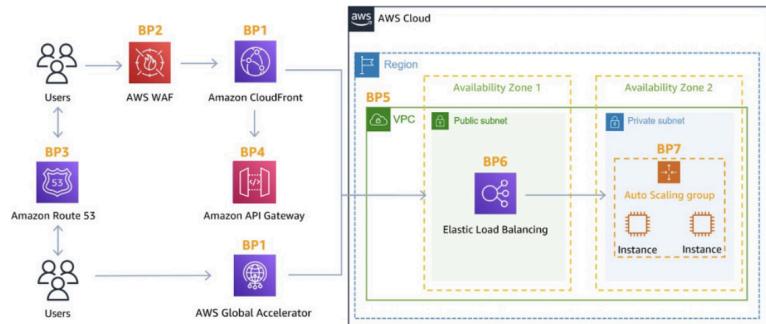
- WAF, Shield and Firewall Manager are used together for comprehensive protection
 - Define your Web ACL rules in WAF
 - For granular protection of your resources, WAF alone is the correct choice
 - If you want to use AWS WAF across accounts, accelerate WAF configuration, automate the protection of new resources, use Firewall Manager with AWS WAF
 - Shield Advanced adds additional features on top of AWS WAF, such as dedicated support from the Shield Response Team (SRT) and advanced reporting.
 - If you're prone to frequent DDoS attacks, consider purchasing Shield Advanced
-
- Used together for comprehensive protection
 - Define Web ACL rules in WAF
 - For granular protection of resources, WAF alone
 - Using WAF across accounts, accelerate WAF configuration, automate protection of new resources, use Firewall Manager with WAF
 - Shield Advanced adds features on top of WAF, protects against DDoS

Best Practices for DDoS Resiliency

Edge Location Mitigation (BP1, BP3)

AWS Best Practices for DDoS Resiliency Edge Location Mitigation (BP1, BP3)

- BP1 – CloudFront
 - Web Application delivery at the edge
 - Protect from DDoS Common Attacks (SYN floods, UDP reflection...)
- BP1 – Global Accelerator
 - Access your application from the edge
 - Integration with Shield for DDoS protection
 - Helpful if your backend is not compatible with CloudFront
- BP3 – Route 53
 - Domain Name Resolution at the edge
 - DDoS Protection mechanism

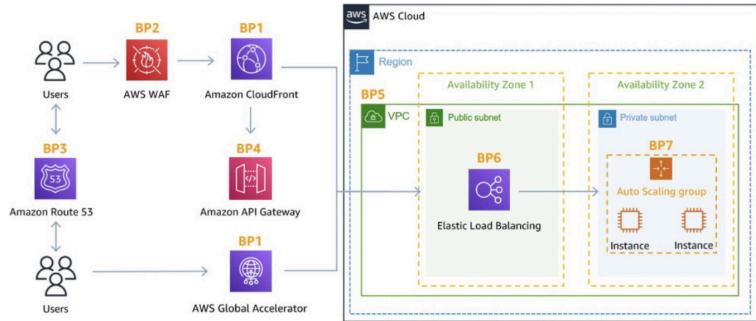


- Best Practice (BP) 1 – CloudFront
 - Web application delivery at edge
 - Protect from DDoS common attacks
- BP1 – Global Accelerator
 - Access app from edge
 - Integration with Shield for DDoS protection
 - Helpful if backend is not compatible with CloudFront
- BP3 – Route 53
 - Domain name resolution at the edge
 - DDoS protection

AWS Best Practices for DDoS Resiliency

Best practices for DDoS mitigation

- Infrastructure layer defense (BP1, BP3, BP6)
 - Protect Amazon EC2 against high traffic
 - That includes using Global Accelerator, Route 53, CloudFront, Elastic Load Balancing
- Amazon EC2 with Auto Scaling (BP7)
 - Helps scale in case of sudden traffic surges including a flash crowd or a DDoS attack
- Elastic Load Balancing (BP6)
 - Elastic Load Balancing scales with the traffic increases and will distribute the traffic to many EC2 instances

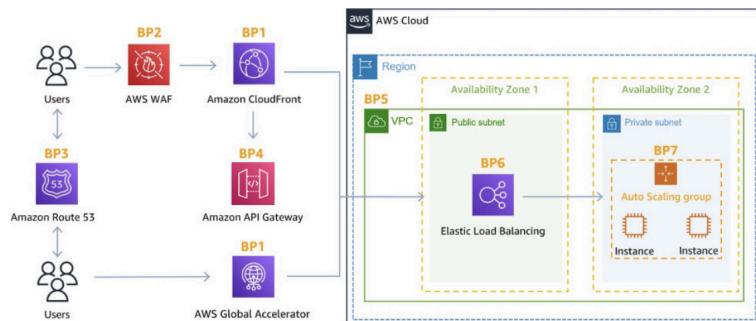


- Infrastructure Layer Defense (BP1, 3, 6)
 - Protect EC2 instances against high traffic
 - Using Global Accelerator, Route 53, CloudFront, ELB
- BP7 – EC2 + Auto scaling
 - Helps scale in traffic surges like DDoS
- BP6 – ELB
 - Scales traffic increases and distribute across instances

AWS Best Practices for DDoS Resiliency

Application Layer Defense

- Detect and filter malicious web requests (BP1, BP2)
 - CloudFront cache static content and serve it from edge locations, protecting your backend
 - AWS WAF is used on top of CloudFront and Application Load Balancer to filter and block requests based on request signatures
 - WAF rate-based rules can automatically block the IPs of bad actors
 - Use managed rules on WAF to block attacks based on IP reputation, or block anonymous IPs
 - CloudFront can block specific geographies
- Shield Advanced (BP1, BP2, BP6)
 - Shield Advanced automatic application layer DDoS mitigation automatically creates, evaluates and deploys AWS WAF rules to mitigate layer 7 attacks



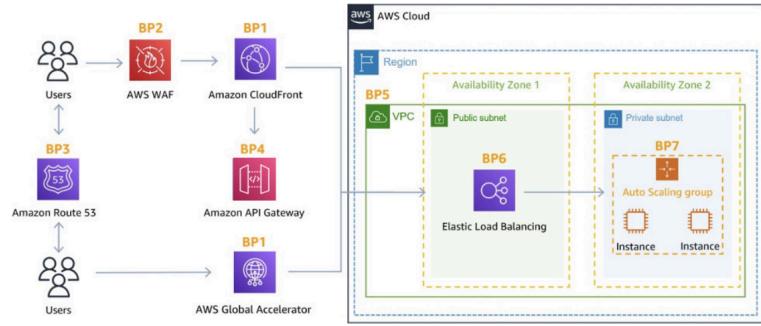
- BP 1, 2 – Detect and filter malicious web requests

- CloudFront cache static content and serve from edge locations, protecting backend
- WAF on top of CloudFront and ALB to filter and block requests based on request signatures
- WAF rate based rules to auto block IPs of bad actors
- WAF managed rules to block attacks based on IP reputation or anonymous IP
- CloudFront can block specific geographies
- Shield Advanced (BP1, 2, 6)
 - Auto application layer DDoS mitigation automatically creates, evaluates and deploys WAF rules to mitigate layer 7 attacks

AWS Best Practices for DDoS Resiliency

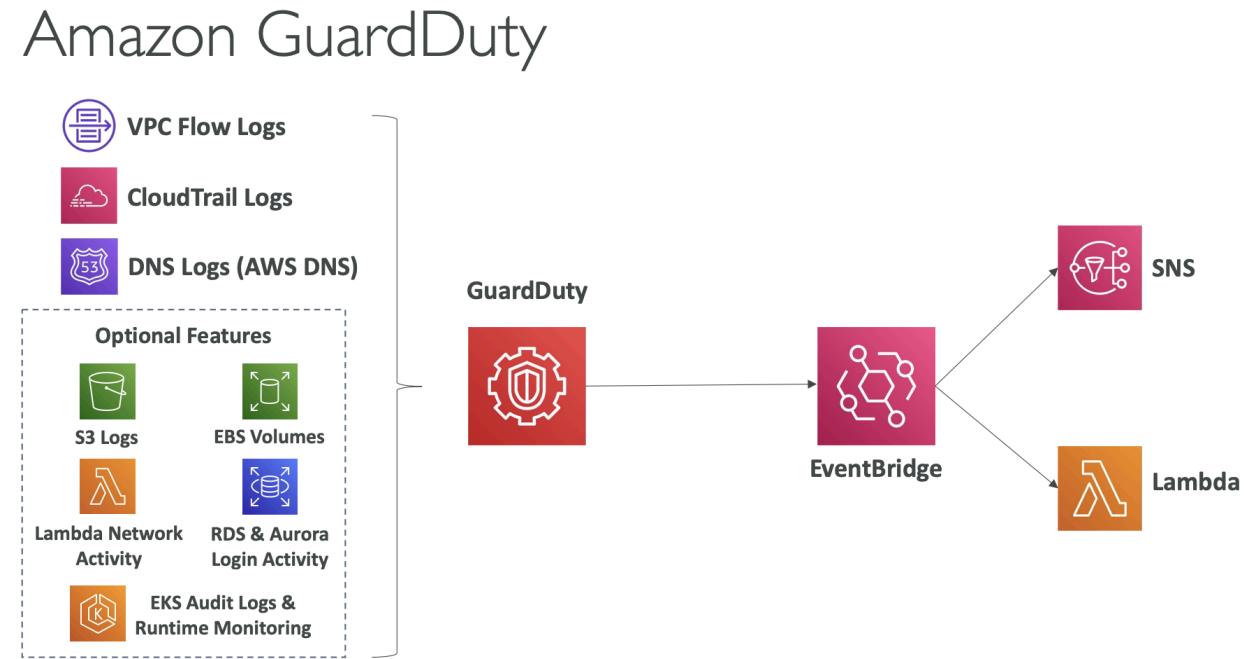
Attack surface reduction

- Obfuscating AWS resources (BP1, BP4, BP6)
 - Using CloudFront, API Gateway, Elastic Load Balancing to hide your backend resources (Lambda functions, EC2 instances)
- Security groups and Network ACLs (BP5)
 - Use security groups and NACLs to filter traffic based on specific IP at the subnet or ENI-level
 - Elastic IP are protected by AWS Shield Advanced
- Protecting API endpoints (BP4)
 - Hide EC2, Lambda, elsewhere
 - Edge-optimized mode, or CloudFront + regional mode (more control for DDoS)
 - WAF + API Gateway: burst limits, headers filtering, use API keys



- Obfuscating AWS resources (BP1, 4, 6)
 - CloudFront, API GW, ELB to hide backend resources
- SG and Network ACLs (BP5)
 - SG and NACLs to filter traffic based on IP at the subnet or ENI level
 - Elastic IP protected by Shield Advanced
- Protecting API endpoints (BP4)
 - Hide EC2, lambda
 - Edge optimized mode or CloudFront + regional mode (more control for DDoS)
 - WAF + API GW: burst limits, header filtering, API keys

Amazon GuardDuty

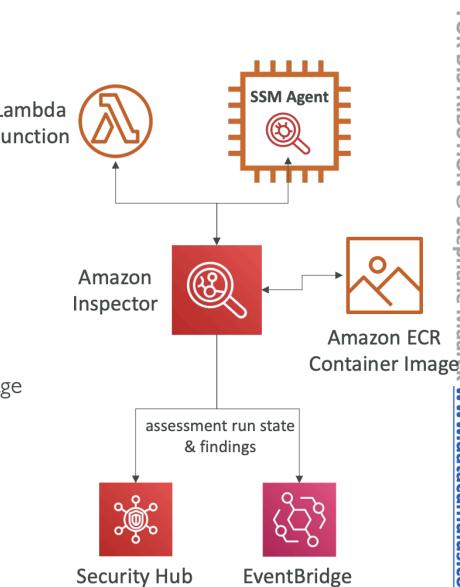


- Intelligent threat discovery via ML for anomaly detection and 3rd party data
 - 1 click, no software install
- Input data includes:
 - CloudTrail Event Logs – unusual API calls, unauthorized deployments...
 - Management events: create VPC subnet, create trail...
 - S3 Data events: get object, delete object...
 - VPC Flow Logs – unusual internal traffic, IP address
 - DNS logs – compromised EC2 instances sending encoded data within DNS queries
 - Optional features: EKS Audit logs, RDS / Aurora, EBS, Lambda, S3 data events...
- EventBridge rules to be notified via Lambda or SNS
- Can protect against cryptocurrency attacks

Amazon Inspector

Amazon Inspector

- Automated Security Assessments
- For EC2 instances
 - Leveraging the AWS System Manager (SSM) agent
 - Analyze against unintended network accessibility
 - Analyze the running OS against known vulnerabilities
- For Container Images push to Amazon ECR
 - Assessment of Container Images as they are pushed
- For Lambda Functions
 - Identifies software vulnerabilities in function code and package dependencies
 - Assessment of functions as they are deployed
- Reporting & integration with AWS Security Hub
- Send findings to Amazon Event Bridge



- Automated security assessments
 - Only for running EC2 instances, container images, and Lambda
- For EC2 instances:
 - Leverage AWS System Manager (SSM) agent
 - Analyze against unintended network accessibility
 - Analyze running OS against known vulnerabilities
- For container images to ECR:
 - Assess container images as they are pushed
- For Lambda:
 - Identified software vulnerabilities in function code and package dependencies
 - Assessment of functions as they are deployed
- Reporting & integration with AWS Security Hub
- Sending findings to EventBridge

What does Amazon Inspector evaluate?

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

Section 27: Networking – VPC

CIDR – IPv4

- Classless Inter-Domain Routing – method for allocating IP
 - Define IP ranges
 - xxx/32 → 1 IP
 - 0.0.0/0 → all IPs
 - 192.168.0.0/26 → 192.168.0.0 – 192.168.0.63 (64 IPs)
- Used in SG rules and AWS Networking
- 2 components
 - Base IP
 - Represents an IP contained in range
 - Subnet mask

Understanding CIDR – Subnet Mask

- The Subnet Mask basically allows part of the underlying IP to get additional next values from the base IP

192	.	168	.	0	.	0	/32 => allows for 1 IP (2^0)	→ 192.168.0.0
192	.	168	.	0	.	0	/31 => allows for 2 IP (2^1)	→ 192.168.0.0 -> 192.168.0.1
192	.	168	.	0	.	0	/30 => allows for 4 IP (2^2)	→ 192.168.0.0 -> 192.168.0.3
192	.	168	.	0	.	0	/29 => allows for 8 IP (2^3)	→ 192.168.0.0 -> 192.168.0.7
192	.	168	.	0	.	0	/28 => allows for 16 IP (2^4)	→ 192.168.0.0 -> 192.168.0.15
192	.	168	.	0	.	0	/27 => allows for 32 IP (2^5)	→ 192.168.0.0 -> 192.168.0.31
192	.	168	.	0	.	0	/26 => allows for 64 IP (2^6)	→ 192.168.0.0 -> 192.168.0.63
192	.	168	.	0	.	0	/25 => allows for 128 IP (2^7)	→ 192.168.0.0 -> 192.168.0.127
192	.	168	.	0	.	0	/24 => allows for 256 IP (2^8)	→ 192.168.0.0 -> 192.168.0.255

192	.	168	.	0	.	0	/16 => allows for 65,536 IP (2^{16})	→ 192.168.0.0 -> 192.168.255.255

192	.	168	.	0	.	0	/0 => allows for All IPs	→ 0.0.0.0 -> 255.255.255.255

Quick Memo			
Octets			
1 st	2 nd	3 rd	4 th
• /32 – no octet can change			
• /24 – last octet can change			
• /16 – last 2 octets can change			
• /8 – last 3 octets can change			
• /0 – all octets can change			

- Allows part of underlying IP to get additional next values from base IP
- Defined how many bits can change in IP
 - /0, /24, /32
- Can take 2 forms
 - /8 → 255.0.0.0
 - /16 → 255.255.0.0
 - /32 → 255.255.255.255

Understanding CIDR – Little Exercise

- $192.168.0.0/24 = \dots ?$
 - $192.168.0.0 - 192.168.0.255$ (256 IPs)
- $192.168.0.0/16 = \dots ?$
 - $192.168.0.0 - 192.168.255.255$ (65,536 IPs)
- $134.56.78.123/32 = \dots ?$
 - Just $134.56.78.123$
- $0.0.0.0/0$
 - All IPs!
- When in doubt, use this website <https://www.ipaddressguide.com/cidr>

Public vs Private IP (IPv4)

- Internet Assigned Numbers Authority (IANA) established certain blocks of IPv4 addresses for private and public addresses
 - Private IP
 - $10.0.0.0 \rightarrow 10.255.255.255$ ($10.0.0.0/8$) → in big networks
 - $172.16.0.0 \rightarrow 172.31.255.255$ ($172.16.0.0/12$) → AWS default VPC in that range
 - $192.168.0.0 \rightarrow 192.168.255.255$ ($192.168.0.0/16$) → home networks
 - Rest are public IPs

Default VPC Overview

- All new accounts have default VPC
- New EC2 instances are launched into default VPC if no subnet is specified
- Default VPC has internet connectivity and all EC2 instances have public IPv4 address
- Get public and private IPv4 DNS names

VPC in AWS – IPv4

- Can have multiple VPC in region (soft limit of 5 per region)
- Max CIDR per VPC is 5; for each CIDR:
 - Min size: /28 (16 IP)
 - Max size /16 (65536 IP)
- Because VPC is private, only private IPv4 ranges allowed
- VPC CIDR should not overlap with other networks (VPCs or corporate networks...)

VPC – Subnet (IPv4)

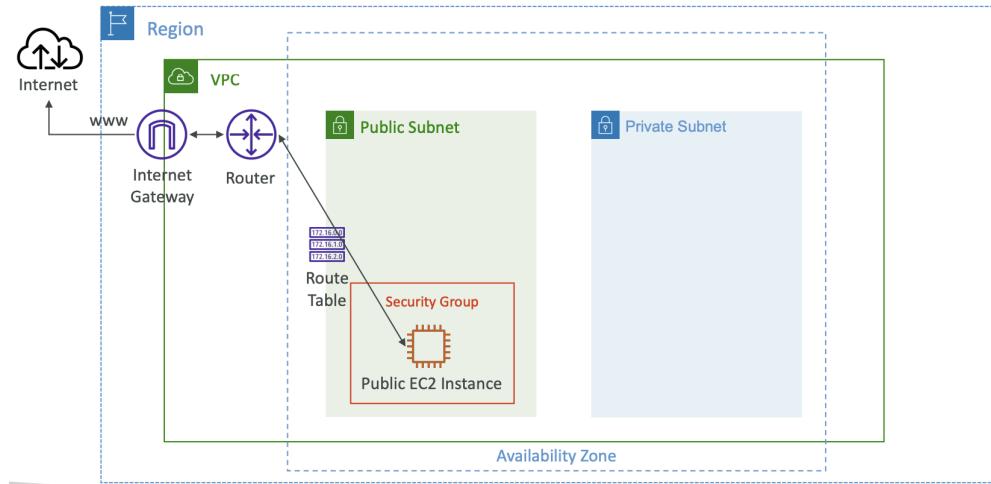
VPC – Subnet (IPv4)



- AWS reserves 5 IP addresses (first 4 & last 1) in each subnet
- These 5 IP addresses are not available for use and can't be assigned to an EC2 instance
- Example: if CIDR block 10.0.0.0/24, then reserved IP addresses are:
 - 10.0.0.0 – Network Address
 - 10.0.0.1 – reserved by AWS for the VPC router
 - 10.0.0.2 – reserved by AWS for mapping to Amazon-provided DNS
 - 10.0.0.3 – reserved by AWS for future use
 - 10.0.0.255 – Network Broadcast Address. AWS does not support broadcast in a VPC, therefore the address is reserved
- Exam Tip, if you need 29 IP addresses for EC2 instances:
 - You can't choose a subnet of size /27 (32 IP addresses, $32 - 5 = 27 < 29$)
 - You need to choose a subnet of size /26 (64 IP addresses, $64 - 5 = 59 > 29$)
- AWS reserves 5 IP addresses (first 4, last 1) in each subnet
 - Are not available for use and can't be assigned to EC2 instance
- If need 29 IP addresses for EC2 instances, can't use /27 because 2^5 IP addresses, so must use /26

Internet Gateway

Editing Route Tables

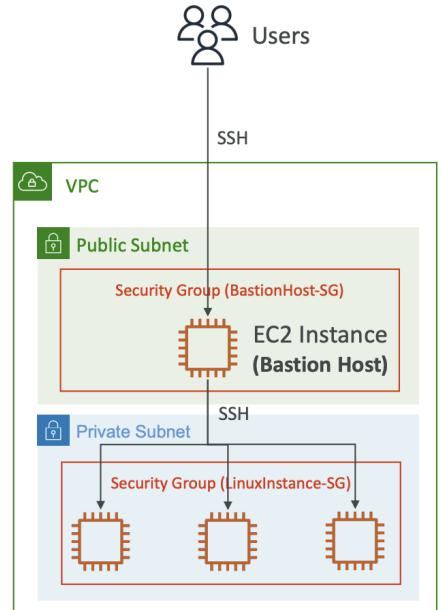


- Allows resources in VPC to connect to internet

- Do not allow internet on its own, need route table
- Scales horizontally and highly available / redundant
- Must be created separately from VPC
 - 1 VPC can only be attached to 1 IGW and vice versa

Bastion Hosts

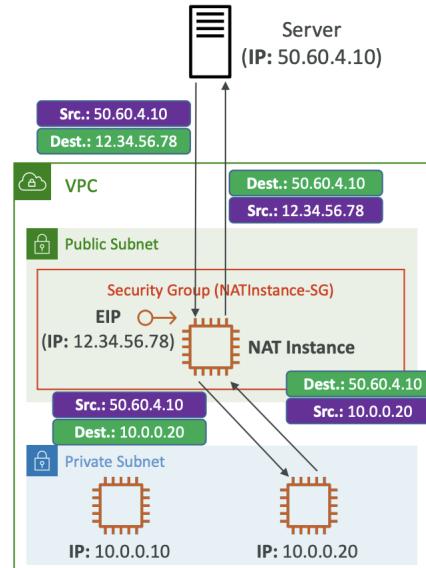
- Used to access private subnet EC2 instances
- Bastion host in the public subnet that then connects to private subnets
 - Bastion host SG must allow inbound from target on port 22 from restricted CIDR
 - SG of EC2 instances must allow SG of Bastion Host or private IP of bastion host



Nat Instance (outdated, still at exam)

NAT Instance (outdated, but still at the exam)

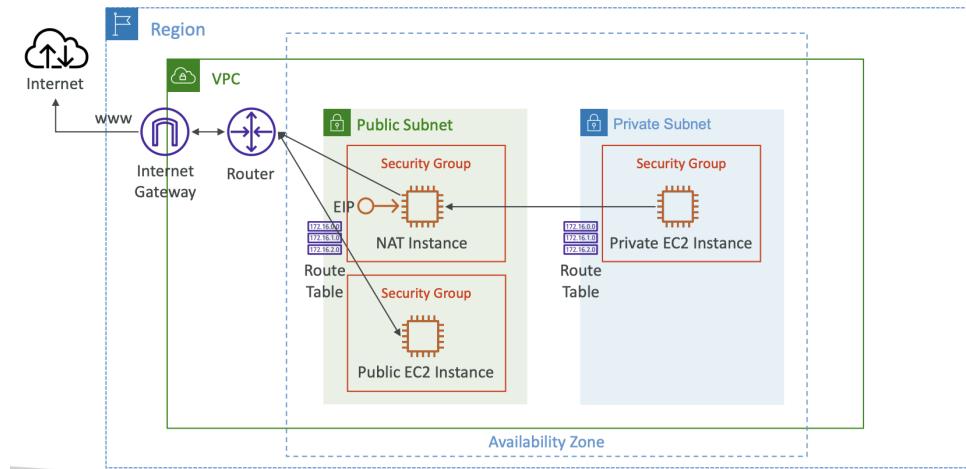
- NAT = Network Address Translation
- Allows EC2 instances in private subnets to connect to the Internet
- Must be launched in a public subnet
- Must disable EC2 setting: Source / destination Check
- Must have Elastic IP attached to it
- Route Tables must be configured to route traffic from private subnets to the NAT Instance



- Network Address Translation
- Allow EC2 instance in private subnets to connect to internet
- Must be launched in public subnet

- Must disable EC2 setting: source / destination check
 - The source / destination at each stage will get rewritten
- Must have Elastic IP attached
- Route table must be configured to route traffic from private subnets to NAT instance

NAT Instance

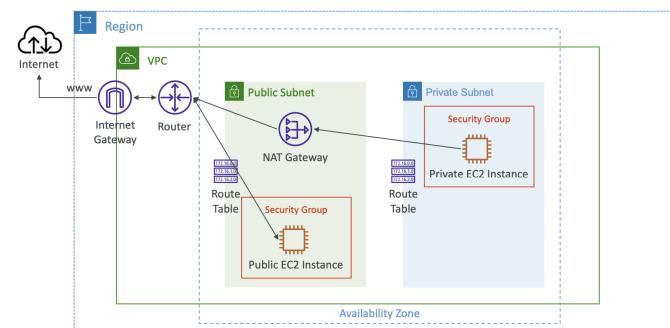


Comments

- Not highly available / resilient setup out of the box
 - Need to create ASG in multi AZ + resilient user data script
- Internet traffic bandwidth depends on EC2 instance type
- Must manage SG & rules
 - Inbound
 - Allow HTTP / S traffic coming from private subnets
 - Allow SSH from home network
 - Outbound
 - Allow HTTP / S traffic to internet

NAT Gateway

NAT Gateway

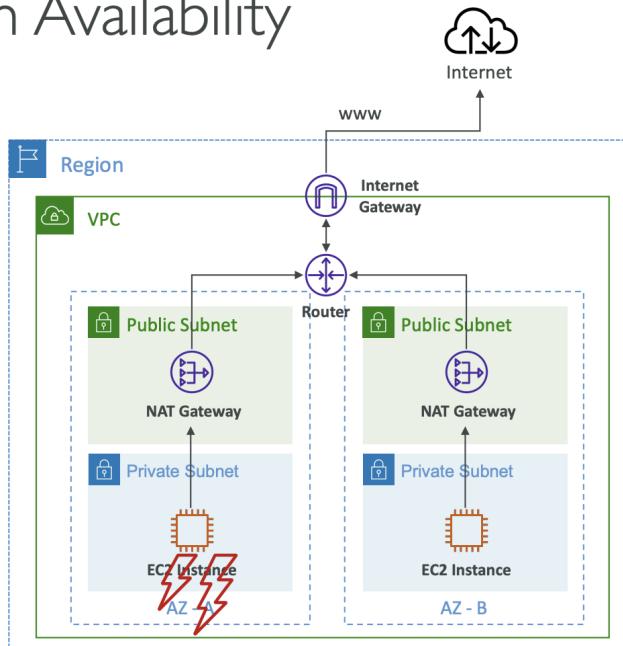


- AWS managed, higher bandwidth, high availability, no administration
 - 5 Gbps of bandwidth with auto scaling up to 100 Gbps
- Pay per hour for usage and bandwidth
- Created in specific AZ, uses Elastic IP
- Can't be used by EC2 instances in same subnet (only other subnets)
- Requires IGW (private subnet → NAT GW → IGW)
- No SG to manage

NAT with High Availability

NAT Gateway with High Availability

- NAT Gateway is resilient within a single Availability Zone
- Must create multiple NAT Gateways in multiple AZs for fault-tolerance
- There is no cross-AZ failover needed because if an AZ goes down it doesn't need NAT



- Resilient within single AZ
 - Must create multiple NAT GW in multiple AZ for fault tolerance
 - No cross AZ failover needed because if an AZ goes down it doesn't need NAT

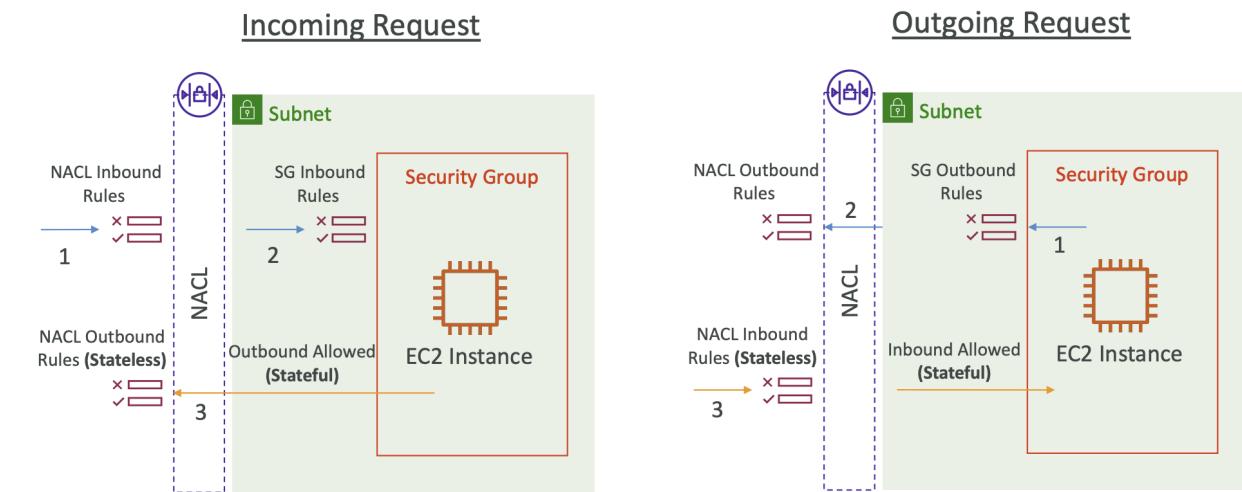
NAT Gateway vs NAT Instance

NAT Gateway vs. NAT Instance

	NAT Gateway	NAT Instance
Availability	Highly available within AZ (create in another AZ)	Use a script to manage failover between instances
Bandwidth	Up to 100 Gbps	Depends on EC2 instance type
Maintenance	Managed by AWS	Managed by you (e.g., software, OS patches, ...)
Cost	Per hour & amount of data transferred	Per hour, EC2 instance type and size, + network \$
Public IPv4	✓	✓
Private IPv4	✓	✓
Security Groups	✗	✓
Use as Bastion Host?	✗	✓

Security Group & NACLs

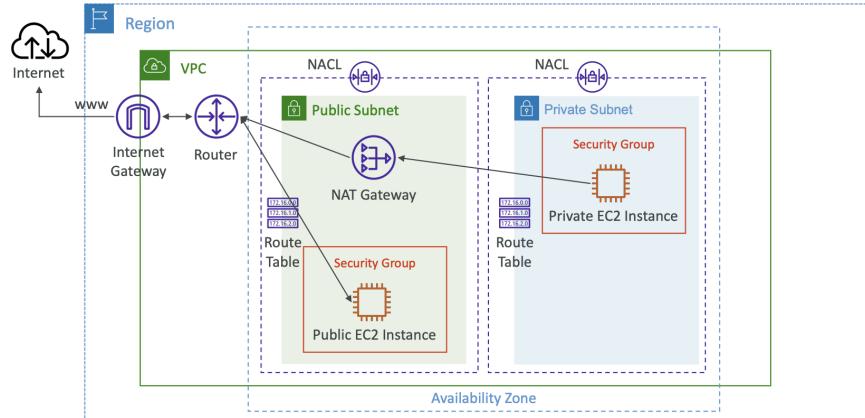
Security Groups & NACLs



- SG is stateful meaning whatever is accepted can go out. NACL is stateless where in and out must be accepted
- NACL is like firewall that control traffic to / from subnet
 - Great at blocking specific IP at the subnet level
- 1 NACL per subnet, new subnets are assigned default NACL
- Define NACL rules:
 - Rules have a number, higher precedence = lower number

- First rule matched will be the decision, where last rule is * and denies if no rule match
- Add rules by increment of 100
- Newly created NACL deny everything

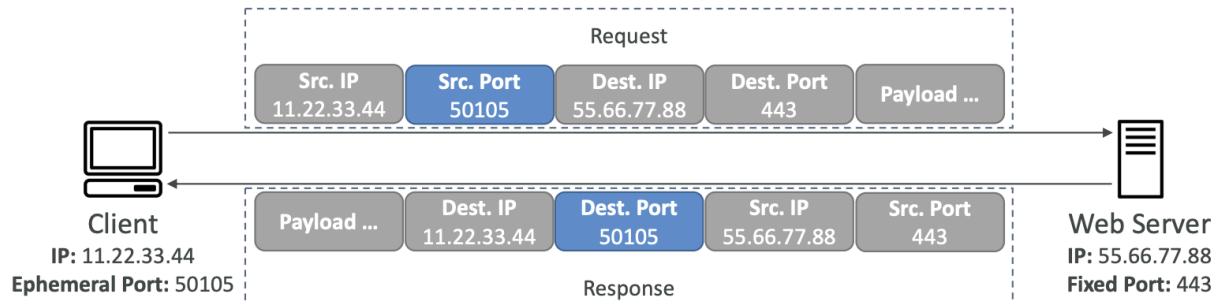
NACLs



Default NACL

- Accepts everything inbound / outbound with the subnets it's associated with
 - Do not modify default, create custom NACL

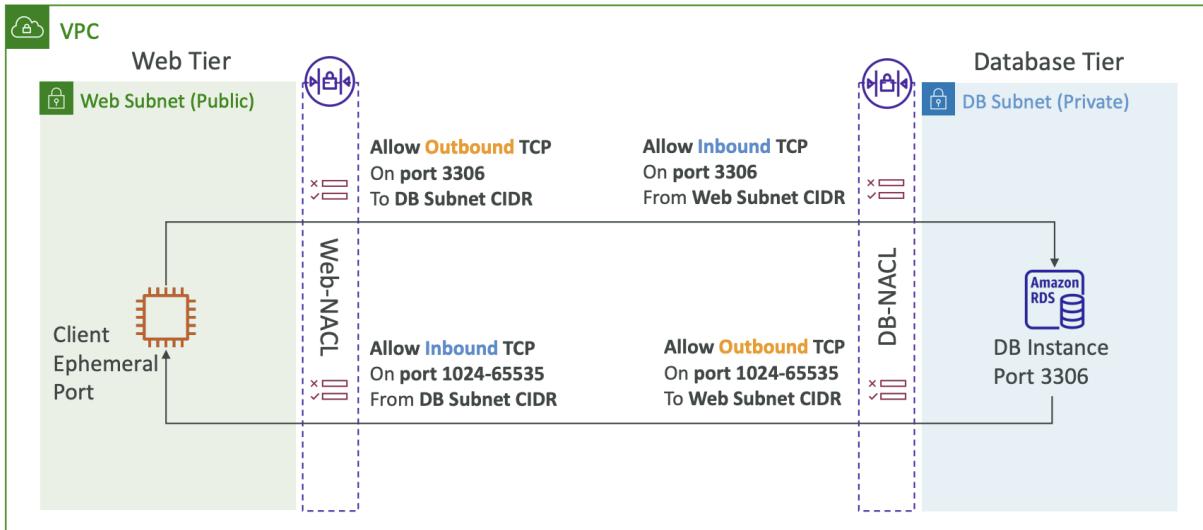
Ephemeral Ports



- For any 2 endpoints to establish a connection, they must use ports
- Clients connect to a defined port and expect a response on an ephemeral port
 - Different OS use different port ranges

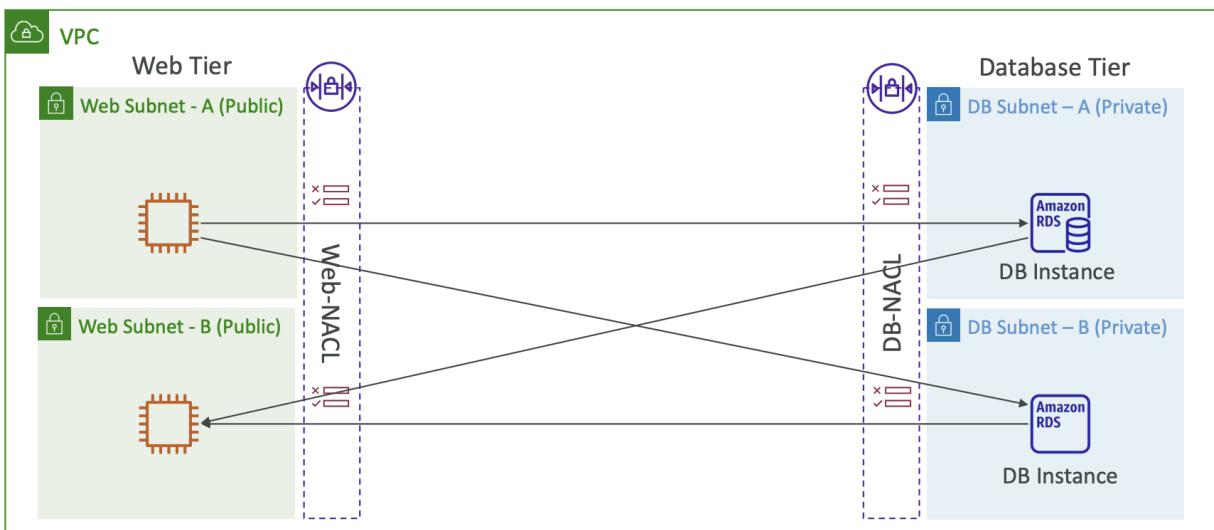
NACL with Ephemeral Ports

NACL with Ephemeral Ports



Create NACL rules for each target subnets CIDR

Create NACL rules for each target subnets CIDR



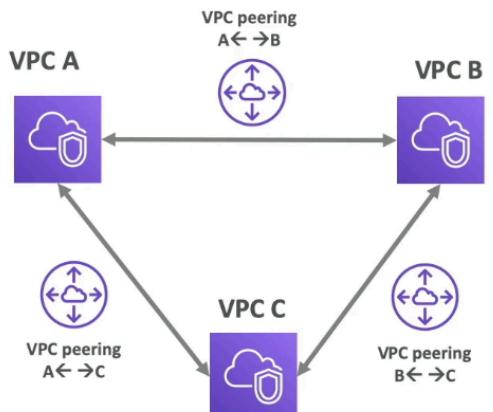
Security Group vs NACLs

Security Group vs. NACLs

Security Group	NACL
Operates at the instance level	Operates at the subnet level
Supports allow rules only	Supports allow rules and deny rules
Stateful: return traffic is automatically allowed, regardless of any rules	Stateless: return traffic must be explicitly allowed by rules (think of ephemeral ports)
All rules are evaluated before deciding whether to allow traffic	Rules are evaluated in order (lowest to highest) when deciding whether to allow traffic, first match wins
Applies to an EC2 instance when specified by someone	Automatically applies to all EC2 instances in the subnet that it's associated with

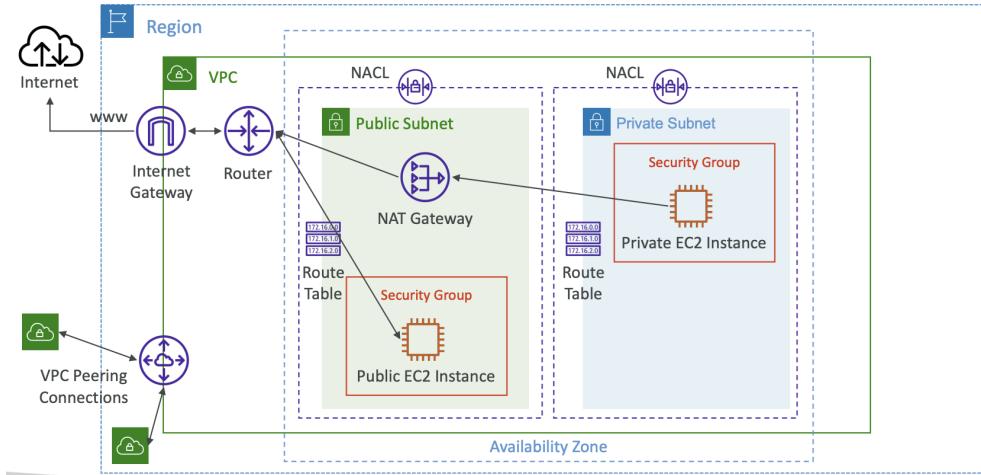
VPC Peering

- Privately connect 2 VPCs using AWS network to make them behave as if they were in the same network
- Must not have overlapping CIDRs
- Not transitive (must be established for each VPC that need to communicate)
- Must update route tables in each VPC's subnets to ensure EC2 instances can communicate with each other



VPC Peering – Good to know

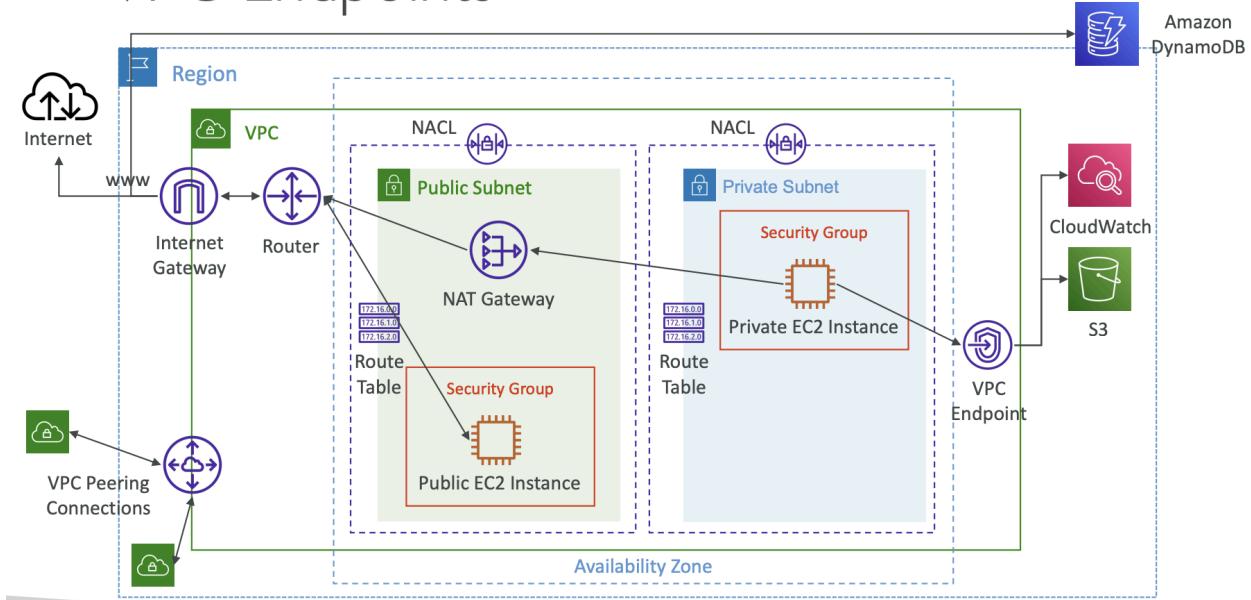
VPC Peering



- Can create VPC peering connection between VPC in different accounts / regions
- Can reference a SG in a peered VPC
 - Works cross account in same region

VPC Endpoints (AWS PrivateLink)

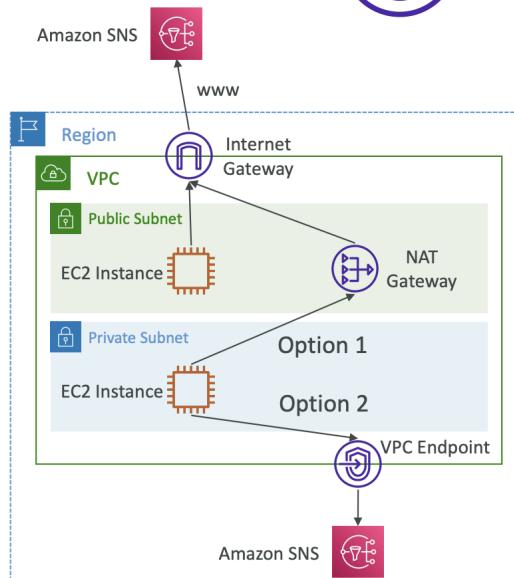
VPC Endpoints



VPC Endpoints (AWS PrivateLink)



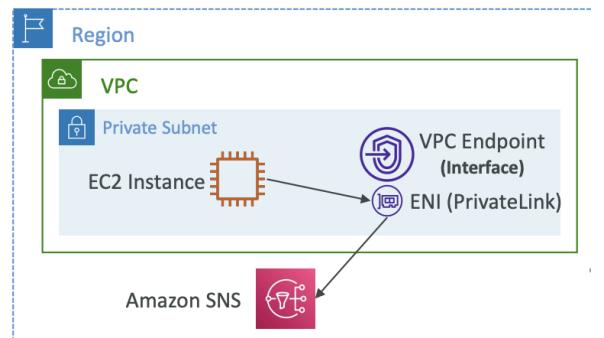
- Every AWS service is publicly exposed (public URL)
- VPC Endpoints (powered by AWS PrivateLink) allows you to connect to AWS services using a private network instead of using the public Internet
- They're redundant and scale horizontally
- They remove the need of IGW, NATGW, ... to access AWS Services
- In case of issues:
 - Check DNS Setting Resolution in your VPC
 - Check Route Tables



- Every AWS service is publicly exposed (public URL)
- VPC endpoints allows connection to AWS services using private internet
 - Redundant and scale horizontally
 - Removes the need for IGW, NATGW...
 - In case of issues:
 - Check DNS setting resolution in VPC
 - Check route table

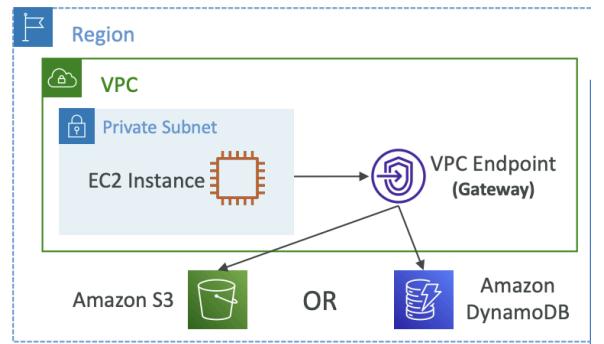
Interface Endpoint (powered by PrivateLink)

- Provisions ENI (private IP) as entry point (must attach to SG)
- Supports most AWS services
- \$ per hour + \$ per GB data processed



Gateway Endpoints

- Provisions a gateway and must be used as a target in a route table (no SG used)
- Supports S3 and DynamoDB
- Free



Gateway or Interface Endpoint for S3?

- Gateway preferred at exam, cost is free
- Interface Endpoint is preferred access is required from on premise, a different VPC or region

VPC Flow Logs

- Capture info about IP traffic going into interfaces:
 - VPC flow logs
 - Subnet flow logs
 - Elastic Network Interface flow logs
- Helps monitor / troubleshoot connectivity issues
- Captures network info from AWS managed interfaces and can be sent to other AWS resources (CloudWatch, S3)

VPC Flow Log Syntax

VPC Flow Logs Syntax

version	interface-id	dstaddr	dstport	packets	start	action
2	123456789010	eni-1235b8ca123456789	172.31.16.139	172.31.16.21	20641	ACCEPT OK
2	123456789010	eni-1235b8ca123456789	172.31.9.69	172.31.9.12	49761	REJECT OK
account-id	srcaddr	srcport	protocol	bytes	end	log-status

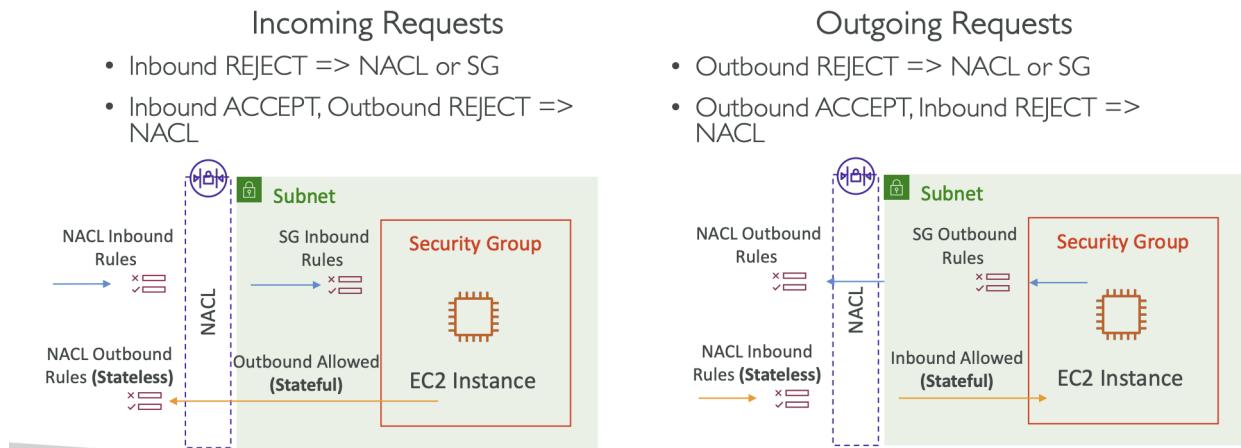
- srcaddr & dstaddr – help identify problematic IP
 - srcport & dstport – help identify problematic ports
 - Action – success or failure of the request due to Security Group / NACL
 - Can be used for analytics on usage patterns, or malicious behavior
 - Query VPC flow logs using Athena on S3 or CloudWatch Logs Insights
 - Flow Logs examples: <https://docs.aws.amazon.com/vpc/latest/userguide/flow-logs-records-examples.html>
- Query using Athena on S3 or CloudWatch Logs Insights

Troubleshoot SG & NACL Issues

- Look at action field

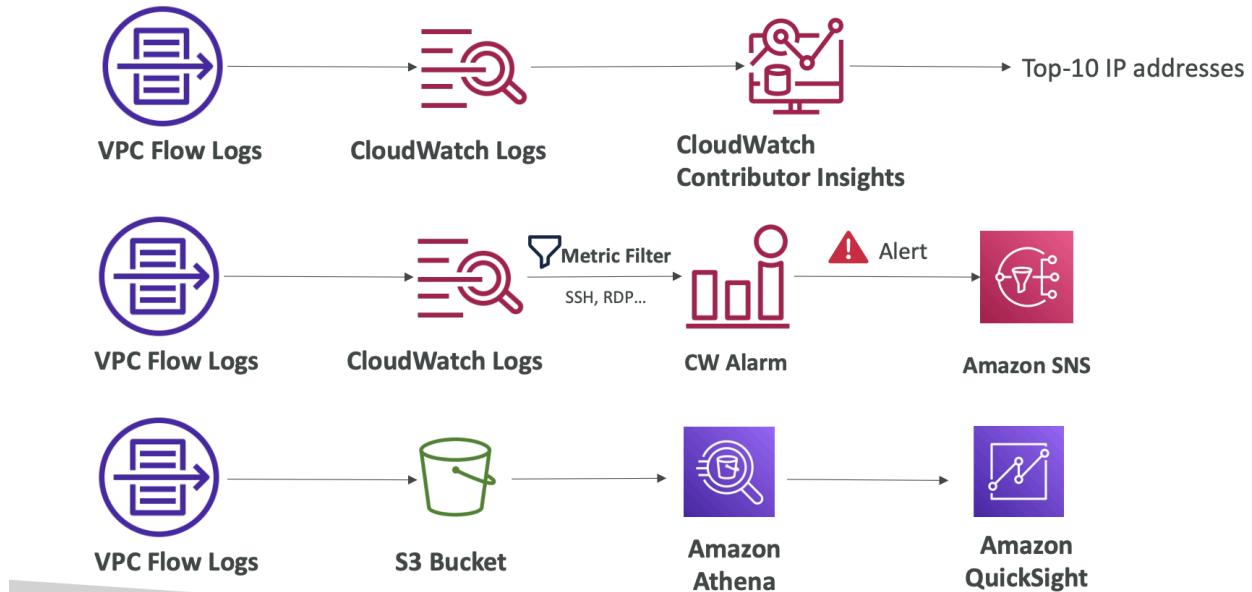
VPC Flow Logs – Troubleshoot SG & NACL issues

Look at the “ACTION” field



Architectures

VPC Flow Logs – Architectures

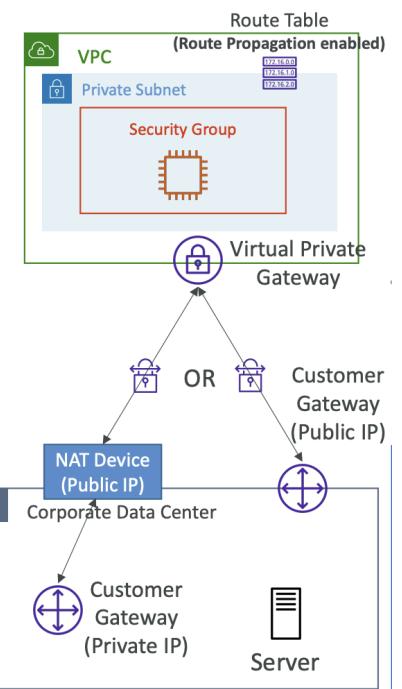


AWS Site to Site VPN

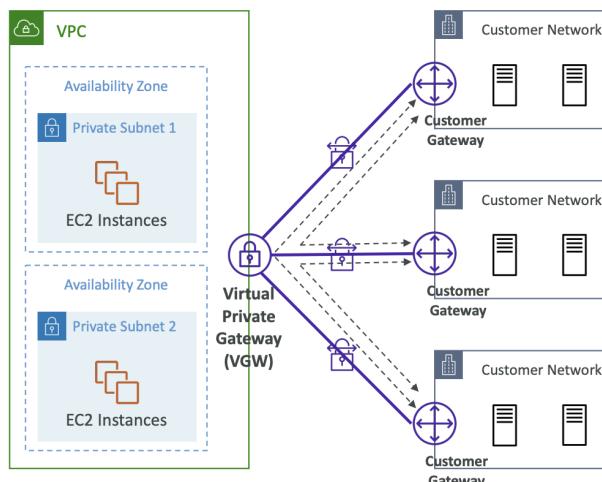
- Virtual Private Gateway (VPW)
 - VPN concentrator on AWS side of VPN connection
 - VGW created and attached to VPC from which you want to create the site to site VPN connection
 - Possibility to customize the ASN (autonomous system number)
- Customer Gateway (CGW)
 - Software or physical device on customer side of VPN connection

Site to Site VPN Connections

- Customer Gateway Device (on premise)
 - What IP to use?
 - Public internet routable IP for customer gateway device
 - If it's behind a NAT device that has NAT traversal (NAT-T), use the public IP of the NAT device
 - Private IP connection
 - **Enable Route Propagation for VPG in route table that is associated with your subnets**
 - If you need to ping EC2 instances from on premise, make sure to add ICMP protocol on the inbound of SG



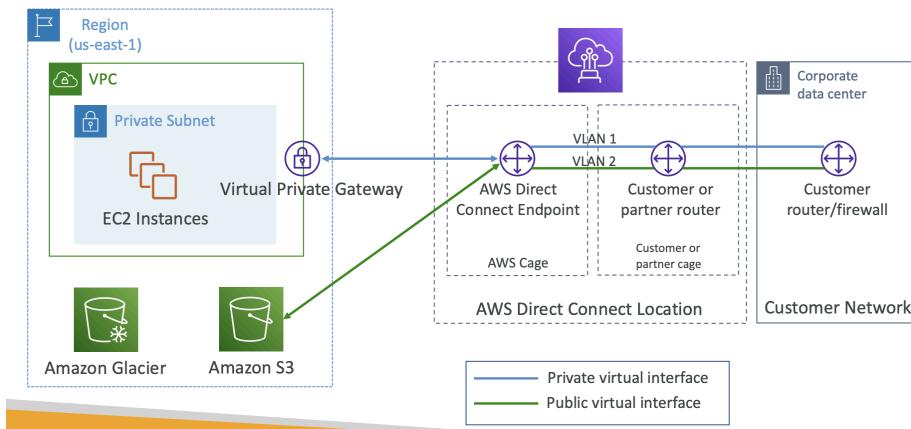
AWS VPN CloudHub



- Provides secure communication between multiple on premise sites if you have multiple VPN connections
 - VPC connection via public internet
 - Connect multiple VPN connections on same VGW, set up dynamic routing and configure route tables
- Low cost hub and spoke model for primary or secondary network connectivity between different locations (VPN only)

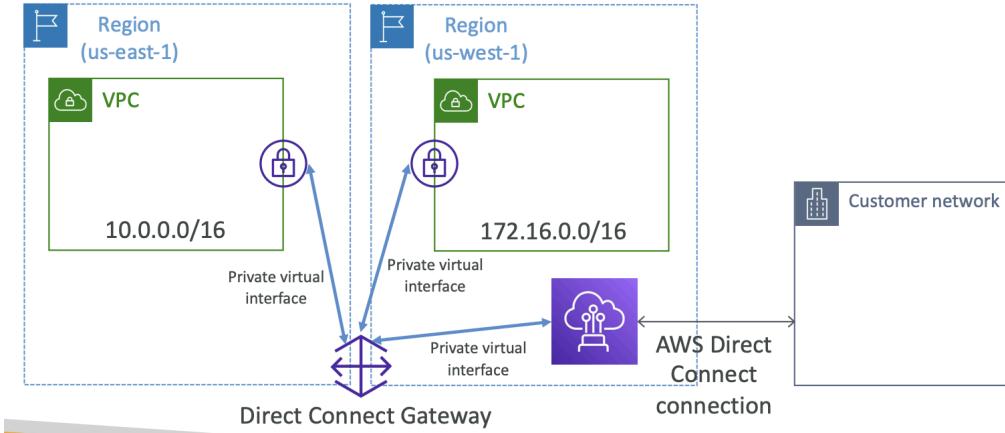
Direct Connect (DX)

Direct Connect Diagram



- Private dedicated connection from remote network to VPC; supports IPv4 and v6
 - Must be set up between on premise and AWS Direct Connect locations, 1+ month to set up
 - Access public and private resources on same connection
- Need to set up a Virtual Private Gateway on VPC
- Dedicated Connection:
 - 1, 10, 100 Gbps capacity + physical ethernet port
- Hosted Connection:
 - 50, 500 Mbps, 10 Gbps
 - Connection requests made via AWS Direct Connect Partners
 - Capacity can be added or removed on demand
- Use cases: increased bandwidth, consistent network, hybrid environments

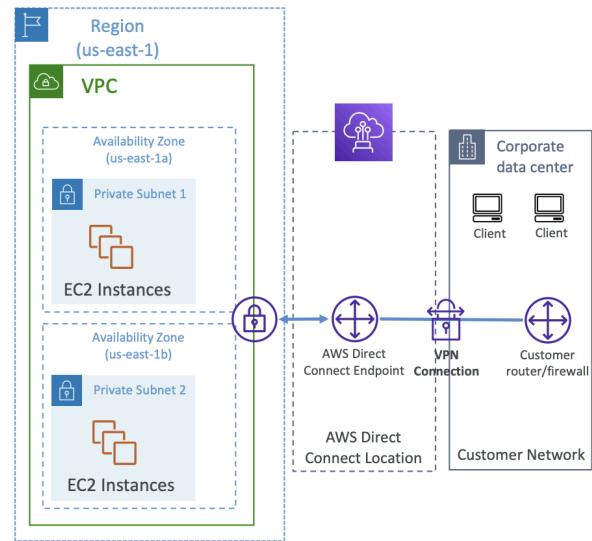
Direct Connect Gateway



- Direct Connect to 1+ VPC in many different regions (same account)

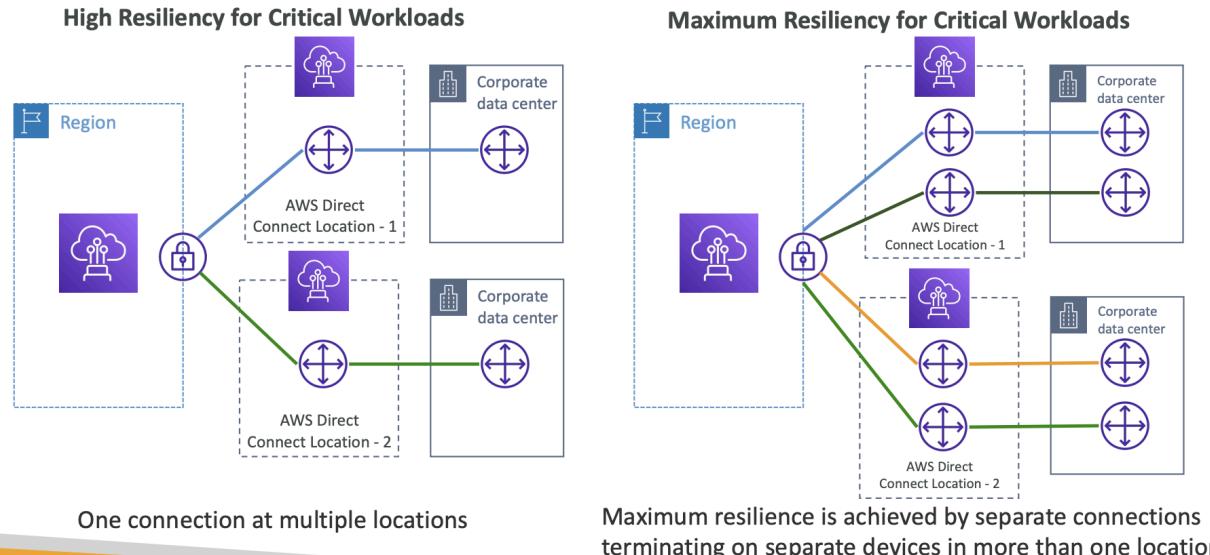
DX Encryption

- Data in transit not encrypted, but private
- AWS Direct Connect + VPN provides IPsec-encrypted private connection
 - Good for extra level of security, more complex



DX – Resiliency

Direct Connect - Resiliency



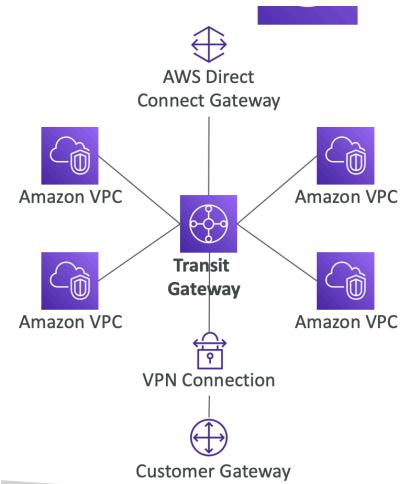
- High Resiliency for Critical Workloads
 - One connection at multiple locations
- Maximum resiliency for critical workloads
 - Separate connections terminated on separate devices in more than 1 location

Site to Site VPN Connection as Backup

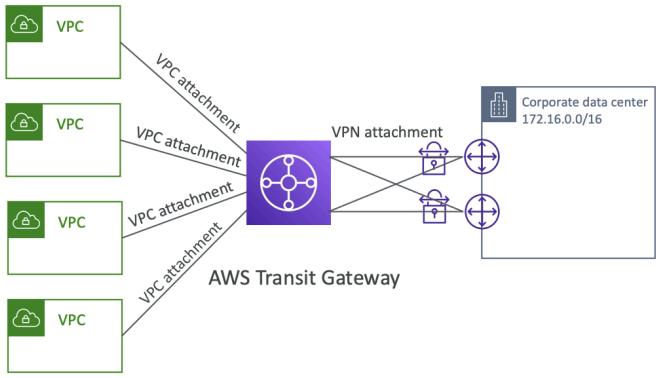
- In case Direct Connect fails, can set a backup Direct Connect connection (expensive) or site to site VPN connection

Transit Gateway

- Transitive peering between thousands of VPC and on-premise, hub and spoke (star) connection
 - Works with Direct Connect Gateway, VPN connections
 - Supports IP multicast
- Regional resource, but can work cross region
 - Share cross account using Resource Access Manager (RAM)
 - Can peer Transit Gateway across regions
- Route tables: Limit which VPC can talk with other VPC



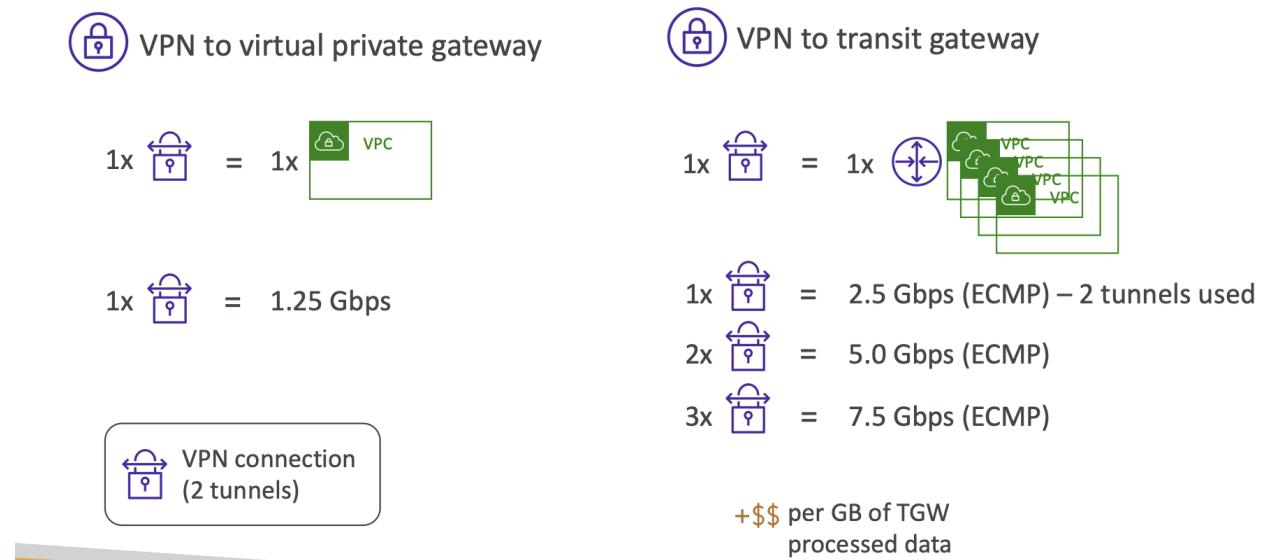
Transit Gateway: Site to Site VPN ECMP



- Equal cost multi path routing
 - Routing strategy to allow to forward a packet over multiple best path
- Use case: create multiple site to site VPN connections to increase bandwidth of connection to AWS

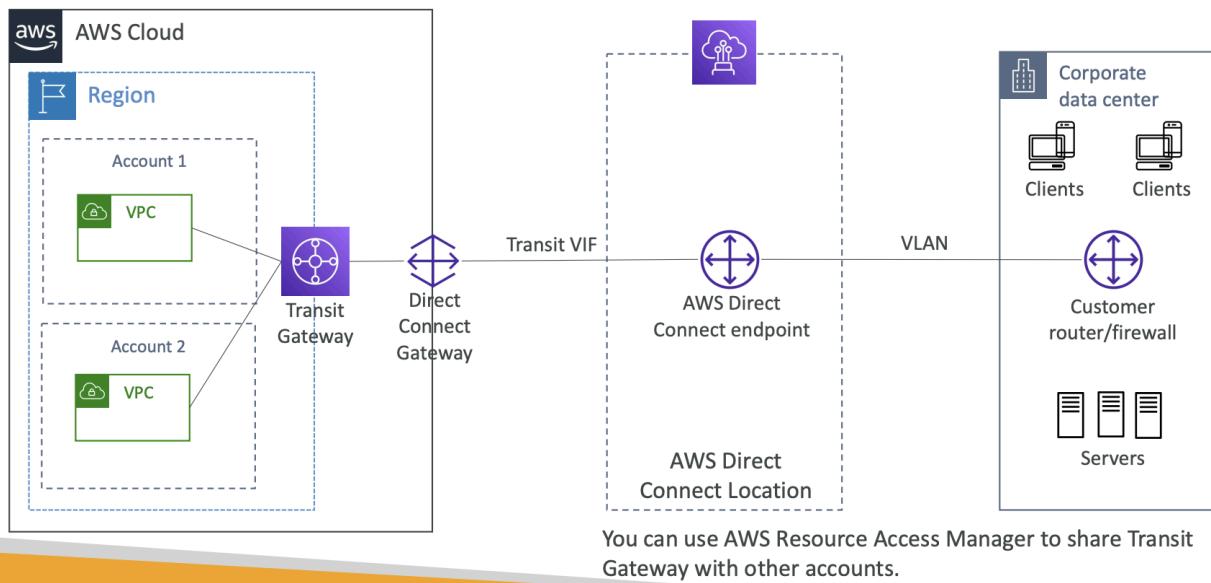
Throughput with ECMP

Transit Gateway: throughput with ECMP



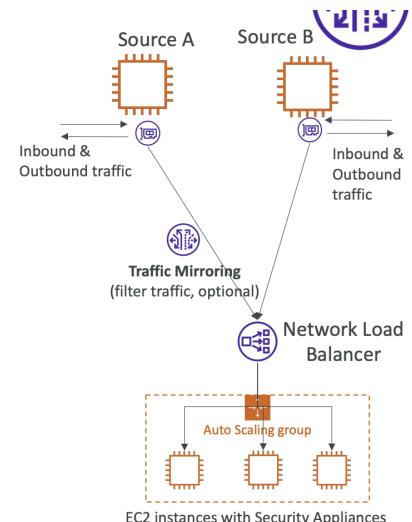
Share Direct Connect between multiple accounts

Transit Gateway – Share Direct Connect between multiple accounts



VPC Traffic Mirroring

- Capture and inspect network traffic in VPC
- Route traffic to security appliances that you manage
- Capture traffic:
 - From (source): ENI
 - Captures all packets or filter
 - To (targets): ENI or NLB
- Source and target can be in same VPC or with VPC peering
- Use case: content inspection, threat monitoring, troubleshooting

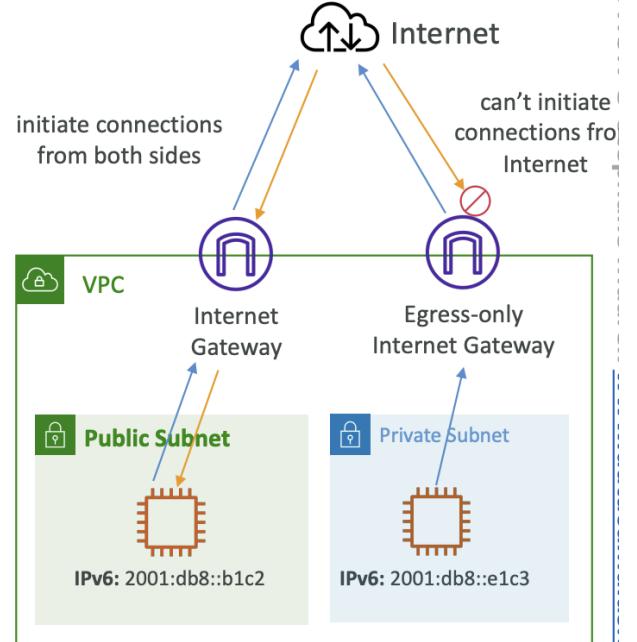


IPv6 in VPC

- Every IPv6 address in AWS is public and internet routable (no private range)
- IPv4 cannot be disabled for VPC and subnets, can enable IPv6 to operate in dual stack mode
 - EC2 instances will at least get internal IPv4 and public IPv6
 - Can communicate using either to internet via Internet Gateway

IPv6 Troubleshooting

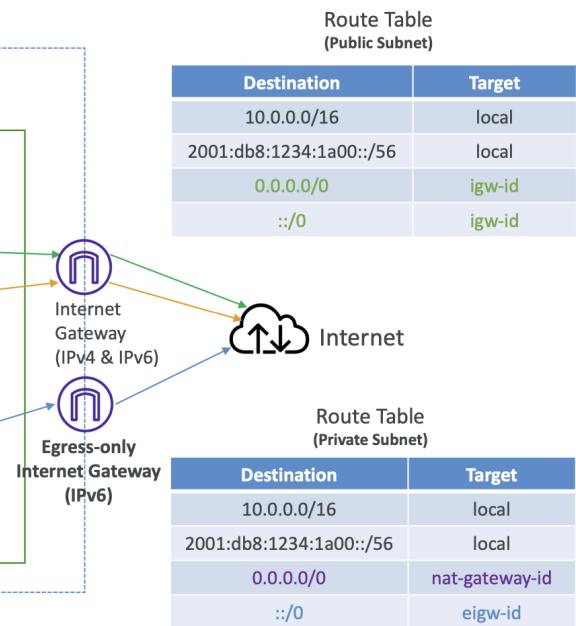
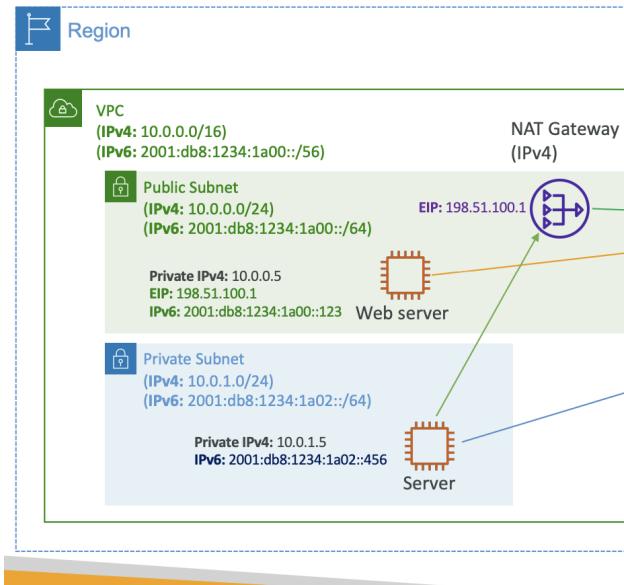
- If you cannot launch an EC2 instance in your subnet with IPv6 enabled VPC, it's because there are no available IPv4 in subnet
 - Solution: create new IPv4 CIDR in subnet



Egress Only Internet Gateway

- Used only for IPv6 traffic, similar to NAT Gateway for IPv6
- Allows instances in VPC outbound connections over IPv6 while preventing internet to initiate IPv6 connection to instances
- Must update route tables

IPv6 Routing



VPC Summary

VPC Section Summary (1/3)

- CIDR – IP Range
- VPC – Virtual Private Cloud => we define a list of IPv4 & IPv6 CIDR
- Subnets – tied to an AZ, we define a CIDR
- Internet Gateway – at the VPC level, provide IPv4 & IPv6 Internet Access
- Route Tables – must be edited to add routes from subnets to the IGW,VPC Peering Connections,VPC Endpoints, ...
- Bastion Host – public EC2 instance to SSH into, that has SSH connectivity to EC2 instances in private subnets
- NAT Instances – gives Internet access to EC2 instances in private subnets. Old, must be setup in a public subnet, disable Source / Destination check flag
- NAT Gateway – managed by AWS, provides scalable Internet access to private EC2 instances, when the target is an IPv4 address

VPC Section Summary (2/3)

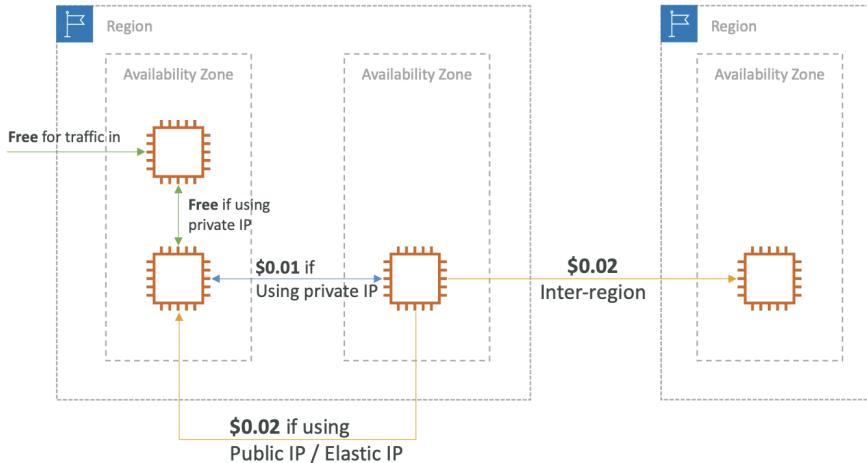
- NACL – stateless, subnet rules for inbound and outbound, don't forget Ephemeral Ports
- Security Groups – stateful, operate at the EC2 instance level
- VPC Peering – connect two VPCs with non overlapping CIDR, non-transitive
- VPC Endpoints – provide private access to AWS Services (S3, DynamoDB, CloudFormation, SSM) within a VPC
- VPC Flow Logs – can be setup at the VPC / Subnet / ENI Level, for ACCEPT and REJECT traffic, helps identifying attacks, analyze using Athena or CloudWatch Logs Insights
- Site-to-Site VPN – setup a Customer Gateway on DC, a Virtual Private Gateway on VPC, and site-to-site VPN over public Internet
- AWS VPN CloudHub – hub-and-spoke VPN model to connect your sites

VPC Section Summary (3/3)

- Direct Connect – setup a Virtual Private Gateway on VPC, and establish a direct private connection to an AWS Direct Connect Location
- Direct Connect Gateway – setup a Direct Connect to many VPCs in different AWS regions
- AWS PrivateLink / VPC Endpoint Services:
 - Connect services privately from your service VPC to customers VPC
 - Doesn't need VPC Peering, public Internet, NAT Gateway, Route Tables
 - Must be used with Network Load Balancer & ENI
- ClassicLink – connect EC2-Classic EC2 instances privately to your VPC
- Transit Gateway – transitive peering connections for VPC, VPN & DX
- Traffic Mirroring – copy network traffic from ENIs for further analysis
- Egress-only Internet Gateway – like a NAT Gateway, but for IPv6 targets

Network Costs in AWS

Networking Costs in AWS per GB - Simplified



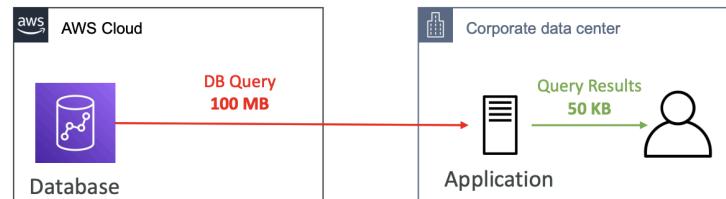
- Use Private IP instead of Public IP for good savings and better network performance
- Use same AZ for maximum savings (at the cost of high availability)

- Use private IP for savings and performance and use same AZ (cost of high availability)

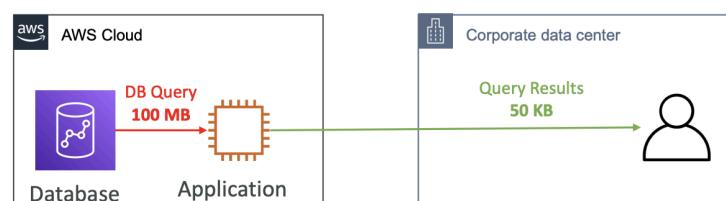
Minimize Egress traffic network cost

- Egress: outbound traffic (AWS to outside)
- Ingress: inbound (outside to AWS)
 - Keep traffic in AWS for cost saving
- Direct connection location that are co-located in same AWS region result in lower cost for egress network

Egress cost is high



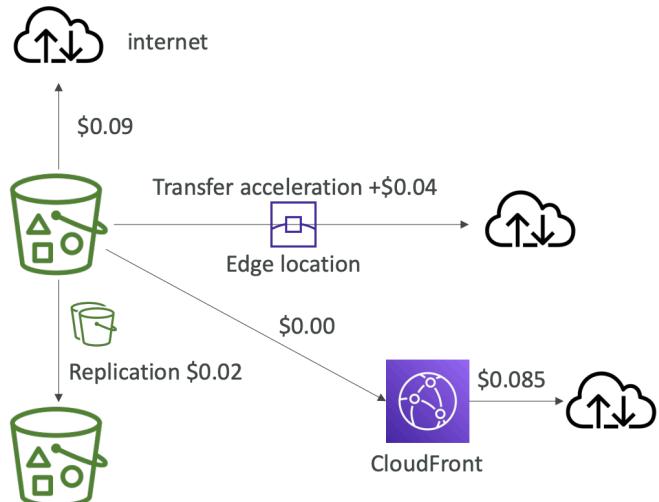
Egress cost is minimized



S3 Data Transfer Pricing

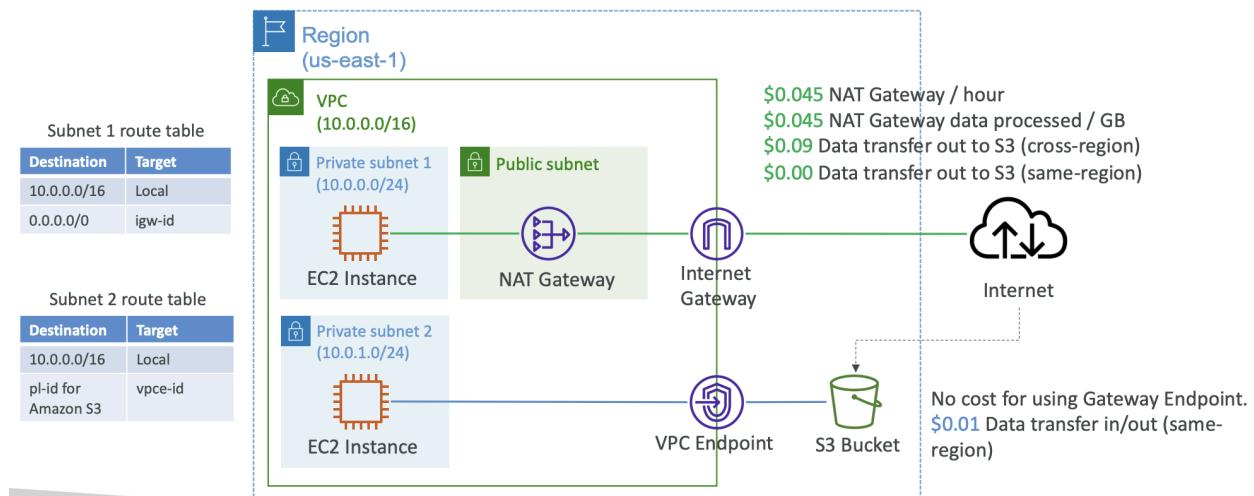
S3 Data Transfer Pricing – Analysis for USA

- S3 ingress: free
- S3 to Internet: \$0.09 per GB
- S3 Transfer Acceleration:
 - Faster transfer times (50 to 500% better)
 - Additional cost on top of Data Transfer Pricing: +\$0.04 to \$0.08 per GB
- S3 to CloudFront: \$0.00 per GB
- CloudFront to Internet: \$0.085 per GB (slightly cheaper than S3)
 - Caching capability (lower latency)
 - Reduce costs associated with S3 Requests Pricing (7x cheaper with CloudFront)
- S3 Cross Region Replication: \$0.02 per GB



- S3 ingress, to CloudFront = free
- S3 to internet, transfer acceleration, CRR = cost

Pricing: NAT Gateway vs Gateway VPC Endpoint



AWS Network Firewall

- Protect entire VPC from layer 3 to 7
- Any direction, inspecting:

- VPC to VPC traffic
- In / outbound traffic
- To / from Direct Connect & Site to Site VPN
- Internally uses Gateway LB
 - Rules centrally managed cross account by AWS Firewall Manager to apply to many VPCs

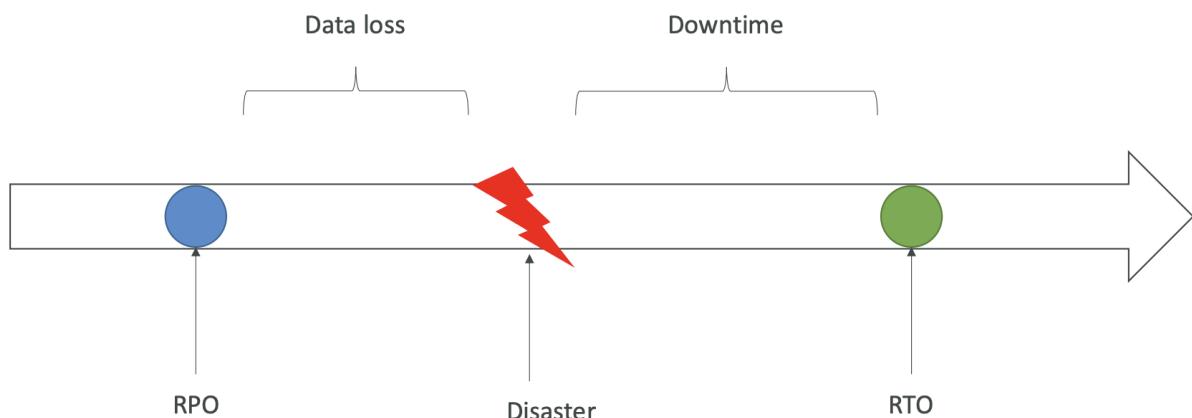
Fine Grained Controls

- Supports many rules
 - IP & port, protocol, domain level, general pattern matching
- Traffic filtering: allow, drop, alert for traffic that matches rules
- Active flow inspection to protect against network threats with intrusion prevention capabilities
- Sends logs to S3, CloudWatch Logs, Kinesis Firehose

Section 28: Disaster Recovery & Migrations

Disaster Recovery Overview

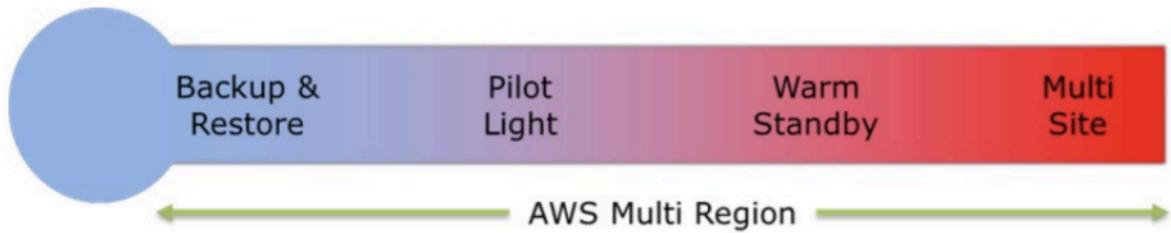
RPO and RTO



- RPO: recovery point objective
 - How much data loss can be accepted?
- RTO: recovery time objective
 - Downtime application has

Disaster Recovery Strategies

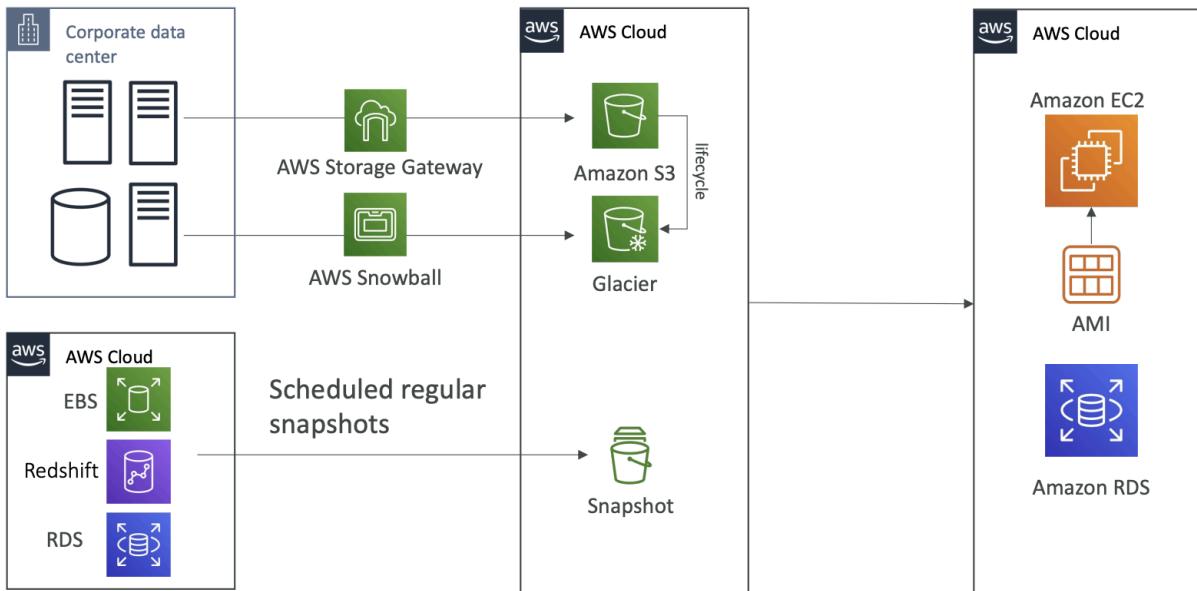
Faster RTO



- Backup and restore
- Pilot light
- Warm standby
- Hot / multi sight approach

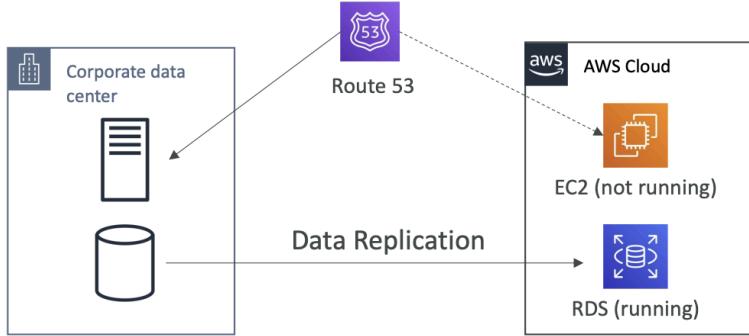
Backup and Restore (High RPO)

Backup and Restore (High RPO)



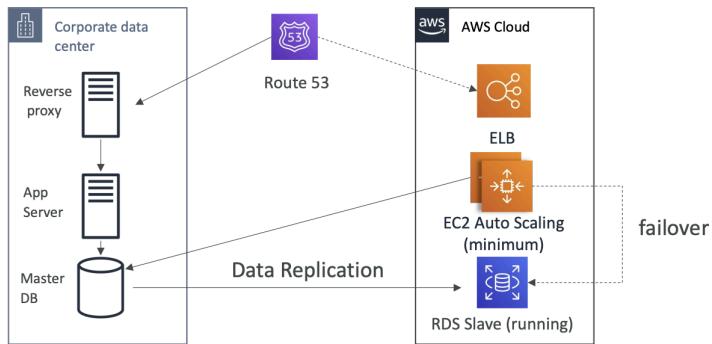
- Inexpensive, high RPO

Pilot Light



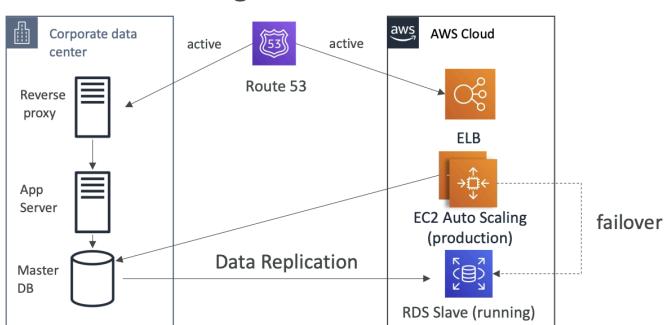
- Small version of app always running on the cloud, similar to backup and restore
- Useful for critical core and faster than backup and restore as critical core is up

Warm Standby



- Full system is running at minimum size and scaled to production load on disaster

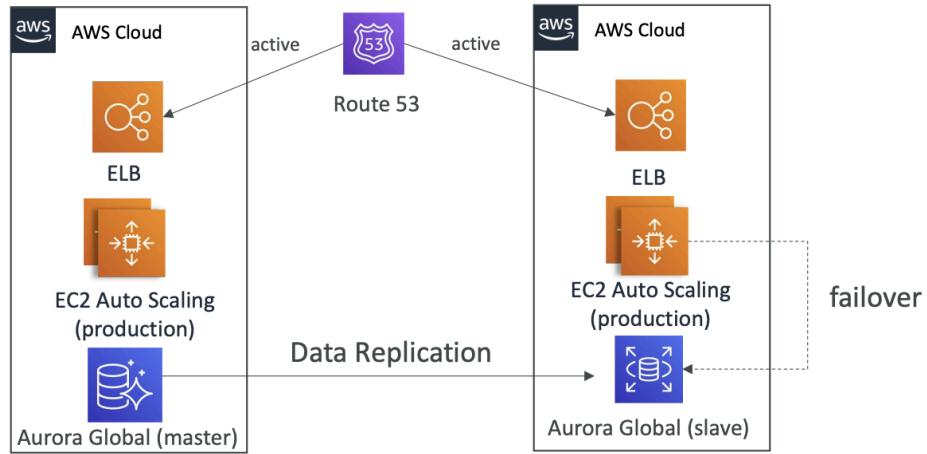
Multi Site / Hot Approach



- Very low RTO, very expensive with full production scale running AWS and on premise

All AWS Multi Region

All AWS Multi Region



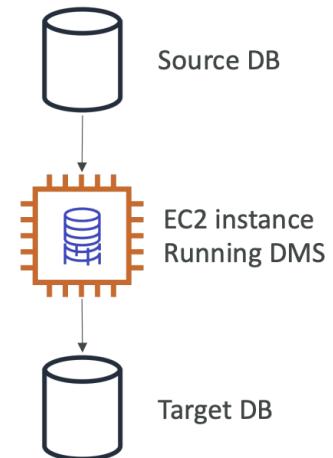
Disaster Recovery Tips

Disaster Recovery Tips

- **Backup**
 - EBS Snapshots, RDS automated backups / Snapshots, etc...
 - Regular pushes to S3 / S3 IA / Glacier, Lifecycle Policy, Cross Region Replication
 - From On-Premise: Snowball or Storage Gateway
- **High Availability**
 - Use Route53 to migrate DNS over from Region to Region
 - RDS Multi-AZ, ElastiCache Multi-AZ, EFS, S3
 - Site to Site VPN as a recovery from Direct Connect
- **Replication**
 - RDS Replication (Cross Region), AWS Aurora + Global Databases
 - Database replication from on-premises to RDS
 - Storage Gateway
- **Automation**
 - CloudFormation / Elastic Beanstalk to re-create a whole new environment
 - Recover / Reboot EC2 instances with CloudWatch if alarms fail
 - AWS Lambda functions for customized automations
- **Chaos**
 - Netflix has a "simian-army" randomly terminating EC2

Database Migration Service (DMS)

- Quickly and securely migrate DB to AWS, resilient, self healing
- Source DB remains available during migration
- Supports:
 - Homogeneous migrations (ex: Oracle to Oracle)
 - Heterogeneous migrations: (ex: Microsoft SQL Server to Aurora)
- Continuous data replication using CDC (change data capture)
- Must create EC2 instance to perform the replication tasks



DMS Sources and Targets

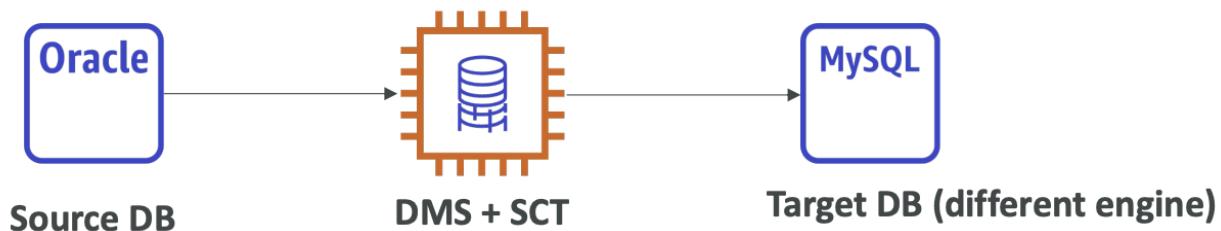
SOURCES:

- On-Premises and EC2 instances databases: *Oracle, MS SQL Server, MySQL, MariaDB, PostgreSQL, MongoDB, SAP, DB2*
- Azure: *Azure SQL Database*
- Amazon RDS: all including Aurora
- Amazon S3
- DocumentDB

TARGETS:

- On-Premises and EC2 instances databases: Oracle, MS SQL Server, MySQL, MariaDB, PostgreSQL, SAP
- Amazon RDS
- Redshift, DynamoDB, S3
- OpenSearch Service
- Kinesis Data Streams
- Apache Kafka
- DocumentDB & Amazon Neptune
- Redis & Babelfish

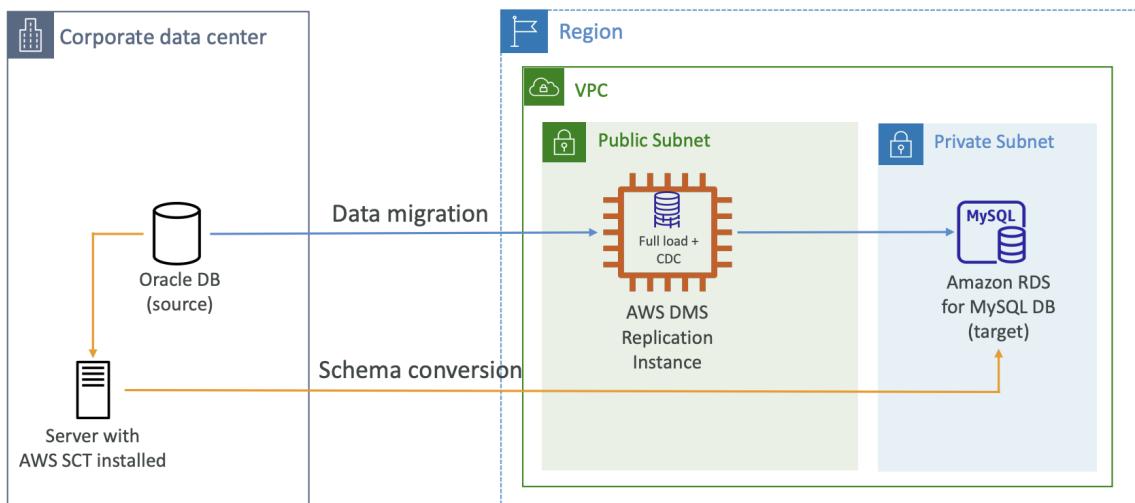
AWS Schema Conversion Tool (SCT)



- Convert DB schema from one engine to another
 - Do not need if migrating the same DB engine
 - Ex: Oracle to Aurora
- Prefer compute intensive instances to optimize data conversions

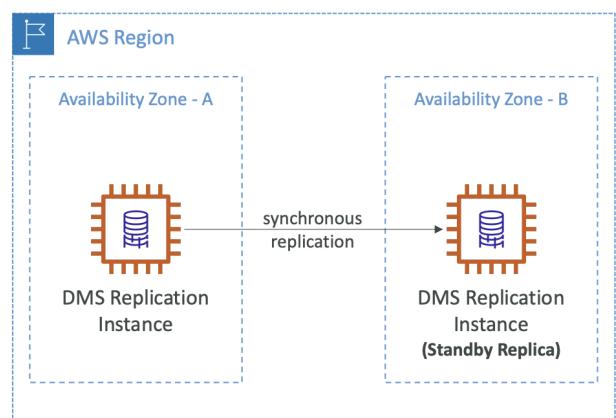
DMS Continuous Replications

DMS - Continuous Replication



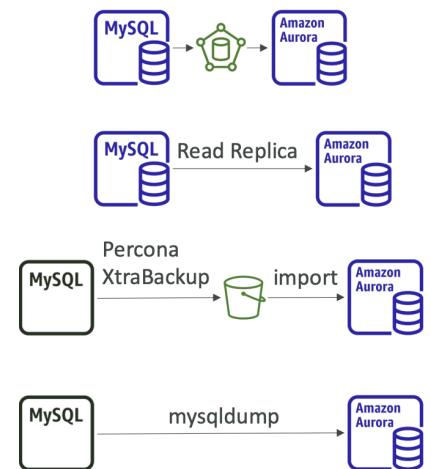
AWS DMS – Multi AZ Deployment

- When multi AZ enabled, DMS provisions and maintains a synchronously stand replica in different AZ
 - Advantages: data redundancy, minimize latency, eliminate I/O freezes



RDS & Aurora MySQL Migrations

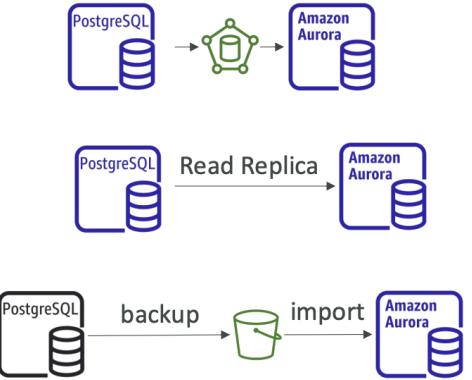
- RDS to Aurora MySQL
 - Option 1: DB Snapshots from RDS MySQL restored as MySQL Aurora DB
 - Potential downtime
 - Option 2: Create an Aurora Read Replica from your RDS MySQL, and when the replication lag is 0, promote it as its own DB cluster (can take time and cost \$)
- External MySQL to Aurora MySQL
 - Option 1:



- Use Percona XtraBackup to create file backup in S3
- Create Aurora MySQL DB from S3
- Option 2:
 - Create Aurora MySQL DB and use mysqldump utility to migrate MySQL into Aurora (slower than S3)
- Use DMS if both DB up and running

RDS & Aurora PostgreSQL Migrations

- RDS PostgreSQL to Aurora PostgreSQL \
 - Option 1: DB Snapshots from RDS PostgreSQL restored as PostgreSQL Aurora DB
 - Option 2: Create Aurora read replica from RDS PostgreSQL and when replication lag is 0, promote as own DB cluster
- External PostgreSQL to Aurora PostgreSQL
 - Create backup into S3 and import via aws_s3 Aurora extension
- Use DMS if both DB up and running

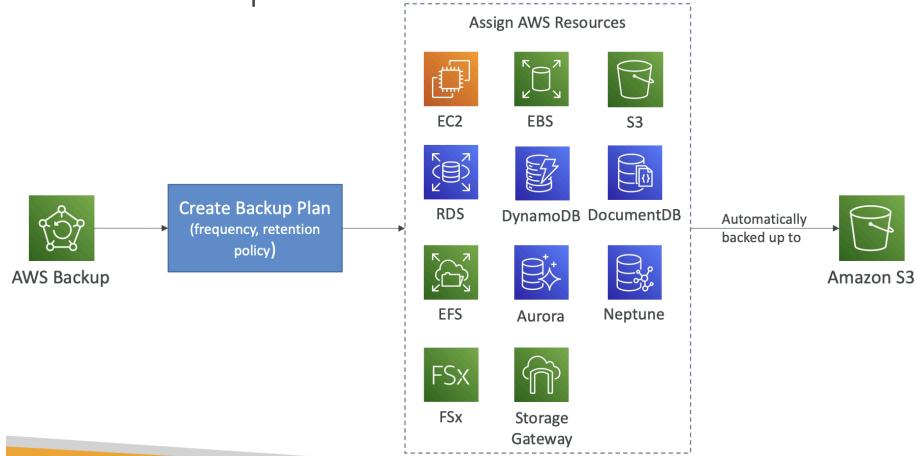


On Premise Strategy with AWS

- Ability to download Amazon Linux 2 AMI as a VM (.iso format)
 - VMWare, KVM, VirtualBox (Oracle VM), Microsoft Hyper-V
- VM Import / Export
 - Migrate existing applications into EC2
 - Create a DR repository strategy for your on-premises VMs
 - Can export back the VMs from EC2 to on-premises
- AWS Application Discovery Service
 - Gather information about your on-premises servers to plan a migration
 - Server utilization and dependency mappings
 - Track with AWS Migration Hub
- AWS Database Migration Service (DMS)
 - replicate On-premise => AWS , AWS => AWS, AWS => On-premise
 - Works with various database technologies (Oracle, MySQL, DynamoDB, etc..)
- AWS Server Migration Service (SMS)
 - Incremental replication of on-premises live servers to AWS

AWS Backup

AWS Backup



- Fully managed service that centrally manage and automate backups across AWS services
 - No need to create custom scripts and manual processes
- Supports many AWS services, cross region backups, cross account backups
- Supports point in time recovery with on demand and scheduled backups
 - Tag based backup policies
- Can create backup policy known as backup plans
 - Backup frequency
 - Backup window
 - Transition to cold storage
 - Retention period

AWS Backup Vault Lock

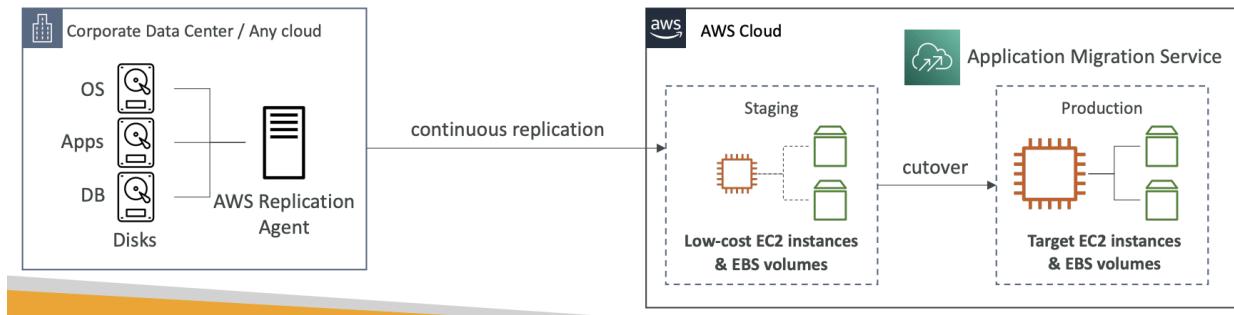
- Enforce WORM (write once read many) state for all backups that you store in AWS Backup Vault
 - Additional layer of defense to protect backups against accidental or malicious delete or updates that shorten or later retention periods
 - Root user cannot delete backup when enabled

AWS Application Discovery Service

- Plan migration projects by gathering info about on premise centers
 - Server utilization data and dependency mapping important for migrations
- Agentless Discovery (AWS Agentless Discovery Connector)
 - VM inventory, configuration, and performance history such as CPU, memory...
- Agent Based Discovery (AWS Application Discovery Agent)

- System configuration, performance, details of network connections between systems
- Resulting data can be viewed in AWS Migration Hub

AWS Application Migration Service (MGN)

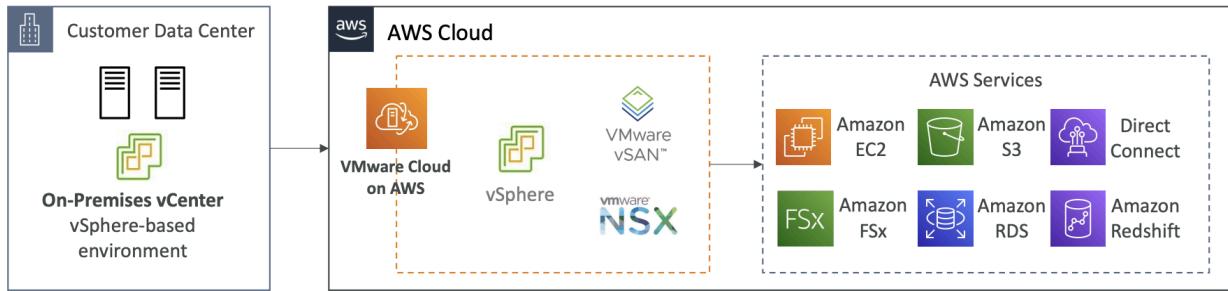


- Lift and shift (rehost) solution which simplify migrating apps to AWS
 - Converts physical, virtual, and cloud based servers to run natively on AWS
 - Supports many platforms, OS, and DB with minimal downtime, reduced cost

Transferring Large amount of Data into AWS

- Example: transfer 200TB of data in the cloud. We have a 100 Mbps internet connection.
- Over the internet / Site-to-Site VPN:
 - Immediate to setup
 - Will take $200(\text{TB}) * 1000(\text{GB}) * 1000(\text{MB}) * 8(\text{Mb}) / 100 \text{ Mbps} = 16,000,000 \text{s} = 185\text{d}$
- Over direct connect 1 Gbps:
 - Long for the one-time setup (over a month)
 - Will take $200(\text{TB}) * 1000(\text{GB}) * 8(\text{Gb}) / 1 \text{ Gbps} = 1,600,000 \text{s} = 18.5\text{d}$
- Over Snowball:
 - Will take 2 to 3 snowballs in parallel
 - Takes about 1 week for the end-to-end transfer
 - Can be combined with DMS
- For on-going replication / transfers: Site-to-Site VPN or DX with DMS or DataSync

VMware Cloud on AWS

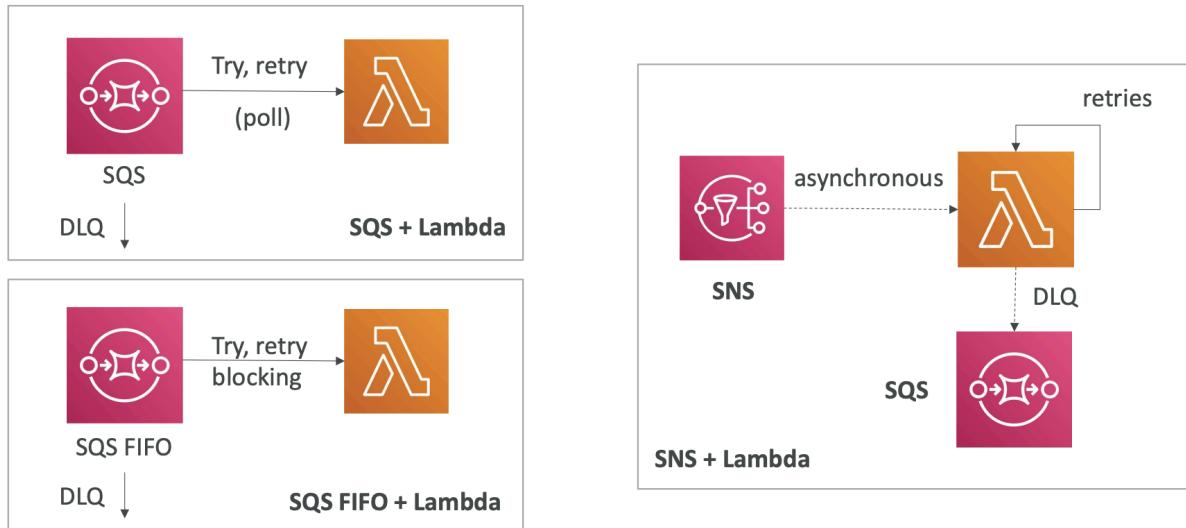


- Migrates VMware vSphere based applications and workloads to AWS
 - Can run production workloads across multiple cloud environments with disaster recovery and AWS services

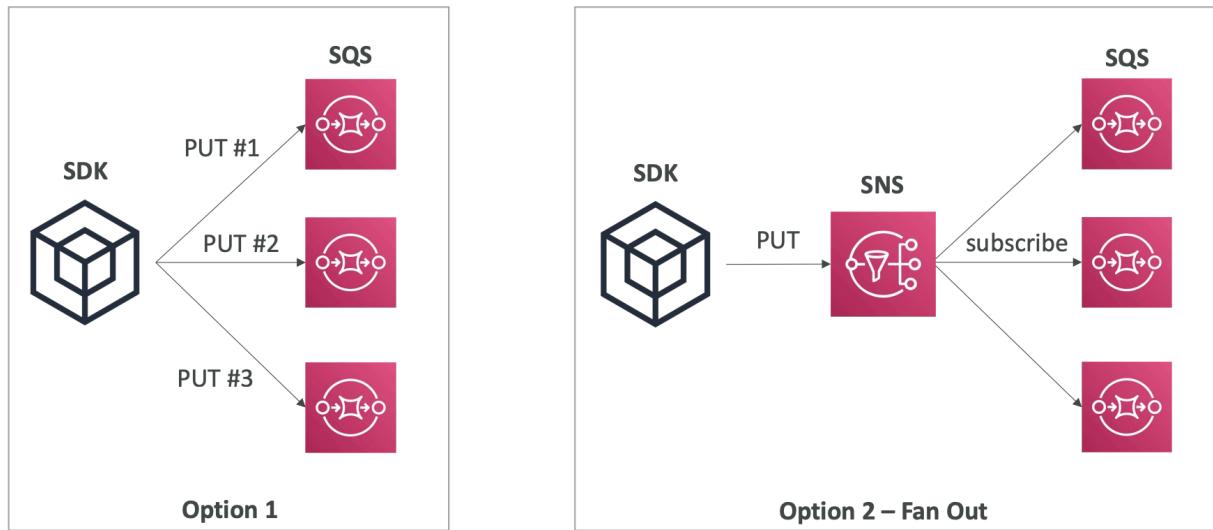
Section 29: More Solution Architectures

Event based Processing

Lambda, SNS & SQS

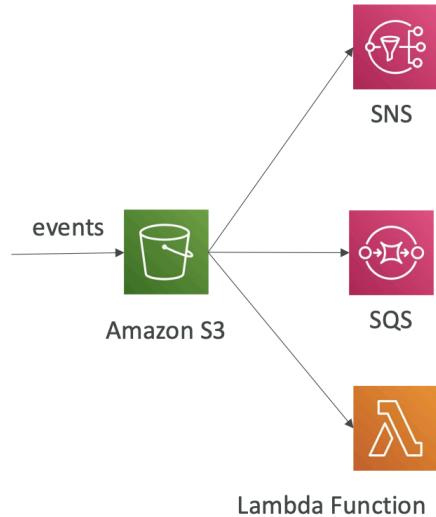


Fan Out Pattern: deliver to multiple SQS

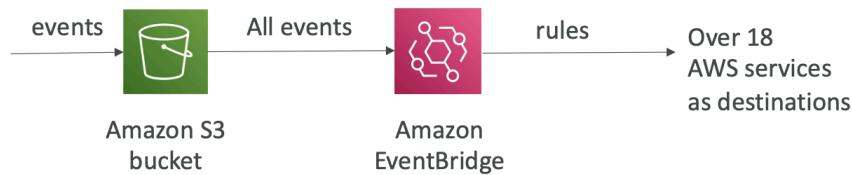


S3 Event Notifications

- S3:ObjectCreated, S3:ObjectRemoved, S3:ObjectRestore, S3:Replication...
 - Object name filtering possible (*.jpg)
 - Use case: generate thumbnails of images uploaded to S3
 - Can create as many “S3 events” as desired
-
- S3 event notifications typically deliver events in seconds but can sometimes take a minute or longer

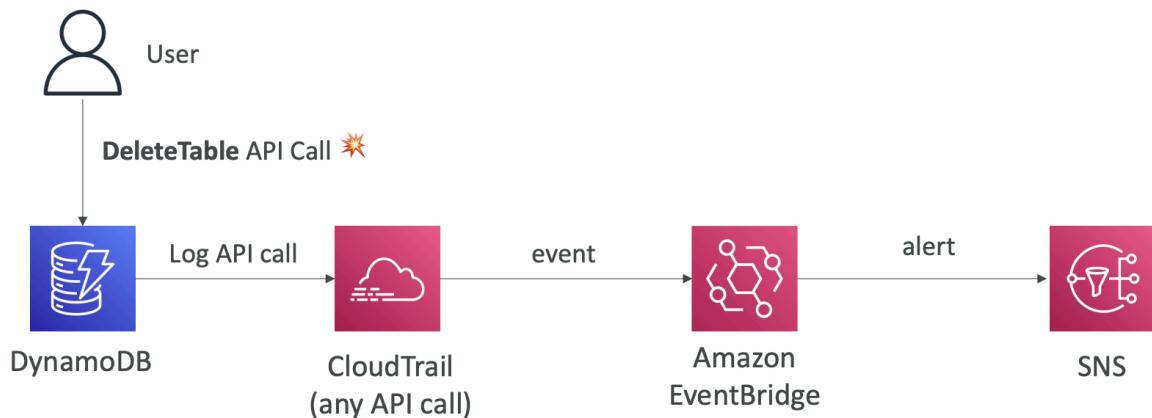


S3 Event Notifications with Amazon EventBridge



- Advanced filtering options with JSON rules (metadata, object size, name...)
- Multiple Destinations – ex Step Functions, Kinesis Streams / Firehose...
- EventBridge Capabilities – Archive, Replay Events, Reliable delivery

Amazon EventBridge – Intercept API Calls

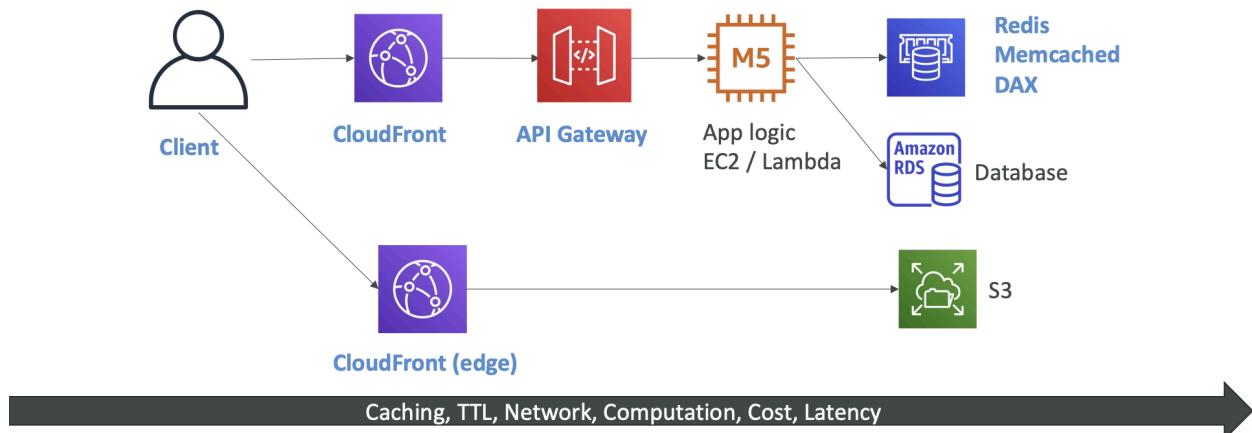


API Gateway – AWS Service Integration Kinesis Data Streams example



Caching Strategies

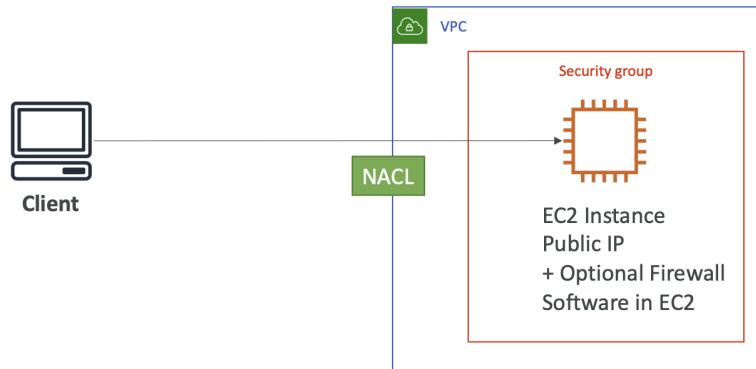
Caching Strategies



- CloudFront is cached at the edge with the potential of stale data, but has fast return to user
- API GW is regional, thus more lag between client and GW
- App logic cache via ElastiCache or DAX allows high frequent reads to be returned faster

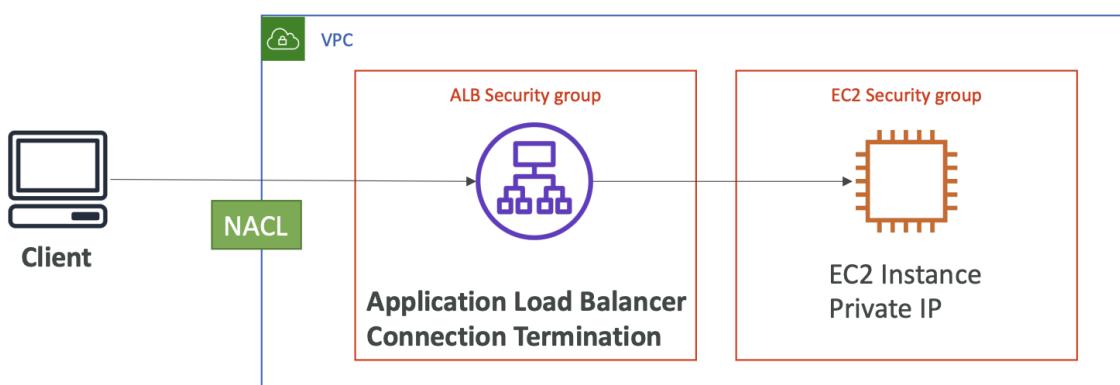
Blocking IP in AWS

Blocking an IP address

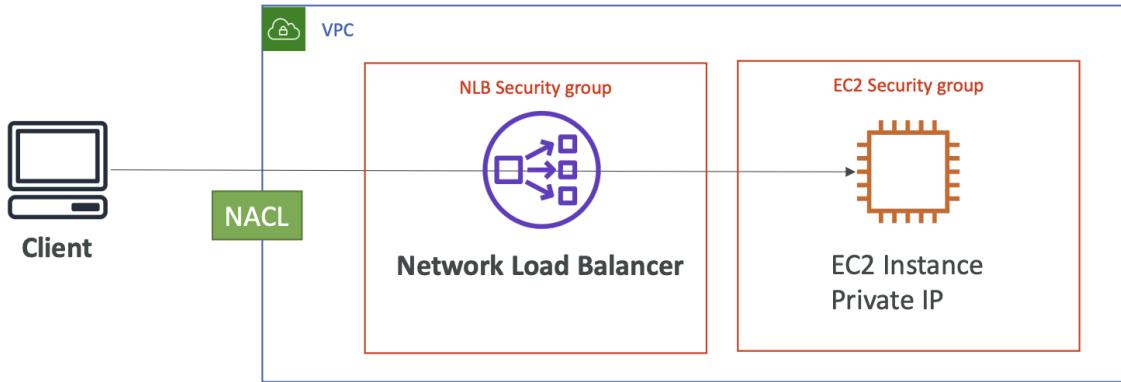


- NACL can deny, but SG can only allow. SG can allow a subset of IP, but is not as effective. Since requests reach the instance, there will be a CPI cost to process the request

Blocking an IP address – with an ALB

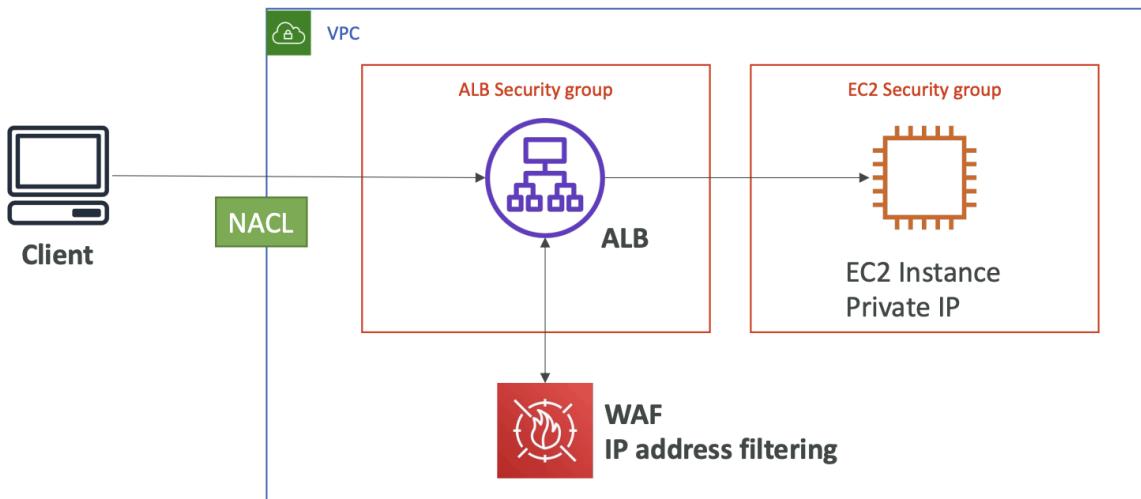


Blocking an IP address – with an NLB



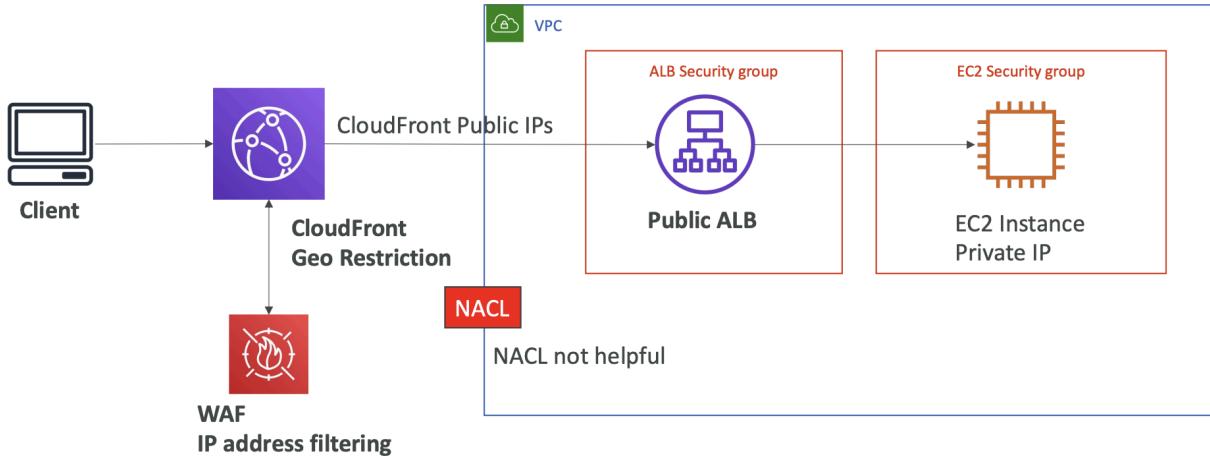
- LB will terminate the client connection and create a new connection from ALB to instance. Instance SG must allow ALB SG as source

Blocking an IP address – ALB + WAF



- WAF can have address filtering that is a service installed on ALB

Blocking an IP address – ALB, CloudFront WAF



- Using CloudFront, the ALB will only see CloudFront IP so NACL on ALB is not helpful. WAF should be on CloudFront instead to see client IPs

High Performance Computing on AWS

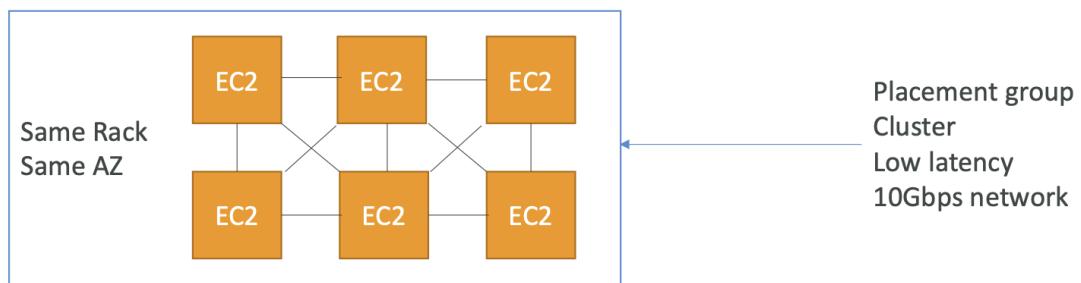
- Can create high number of resources and speed up time to results by adding more resources; pay for only what you use

Data Management & Transfer

- # Data Management & Transfer
- AWS Direct Connect:
 - Move GB/s of data to the cloud, over a private secure network
 - Snowball & Snowmobile
 - Move PB of data to the cloud
 - AWS DataSync
 - Move large amount of data between on-premises and S3, EFS, FSx for Windows

Compute and Networking

- # Compute and Networking
- EC2 Instances:
 - CPU optimized, GPU optimized
 - Spot Instances / Spot Fleets for cost savings + Auto Scaling
 - EC2 Placement Groups: Cluster for good network performance



Compute and Networking

- EC2 Enhanced Networking (SR-IOV)
 - Higher bandwidth, higher PPS (packet per second), lower latency
 - Option 1: Elastic Network Adapter (ENA) up to 100 Gbps
 - Option 2: Intel 82599 VF up to 10 Gbps – LEGACY
- Elastic Fabric Adapter (EFA)
 - Improved ENA for HPC, only works for Linux
 - Great for inter-node communications, tightly coupled workloads
 - Leverages Message Passing Interface (MPI) standard
 - Bypasses the underlying Linux OS to provide low-latency, reliable transport

Storage

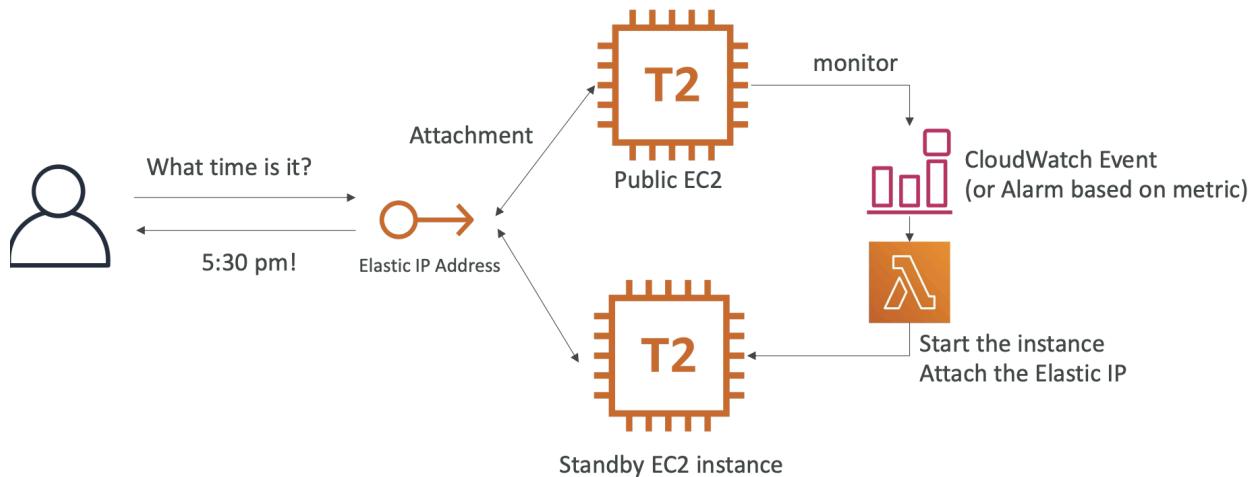
- Instance-attached storage:
 - EBS: scale up to 256,000 IOPS with io2 Block Express
 - Instance Store: scale to millions of IOPS, linked to EC2 instance, low latency
- Network storage:
 - Amazon S3: large blob, not a file system
 - Amazon EFS: scale IOPS based on total size, or use provisioned IOPS
 - Amazon FSx for Lustre:
 - HPC optimized distributed file system, millions of IOPS
 - Backed by S3

Automation and Orchestration

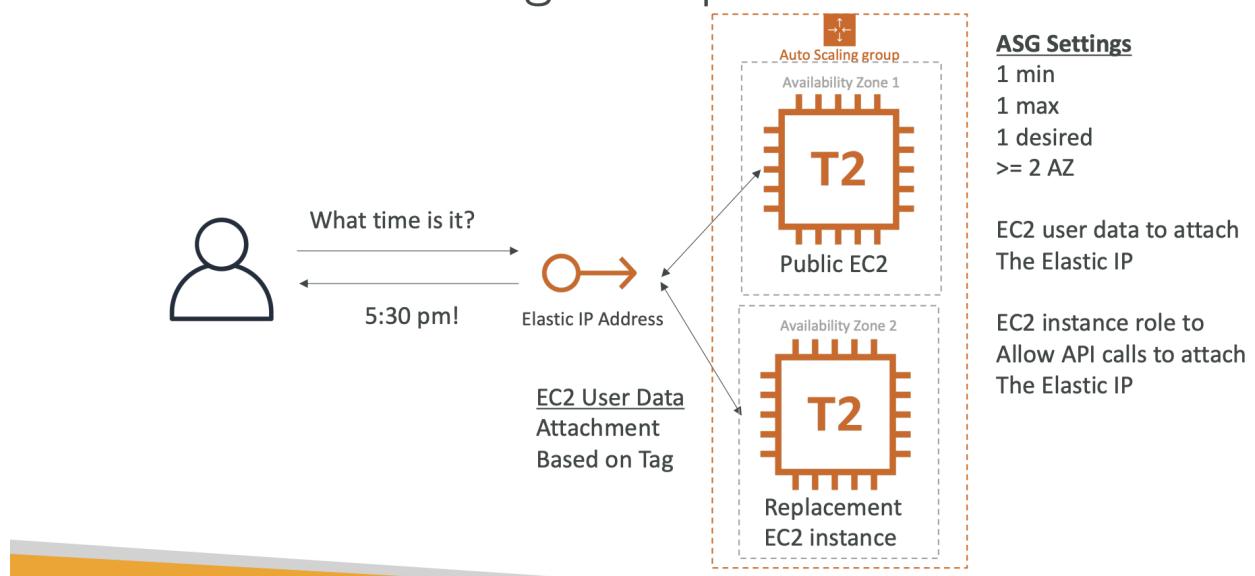
- AWS Batch
 - AWS Batch supports multi-node parallel jobs, which enables you to run single jobs that span multiple EC2 instances.
 - Easily schedule jobs and launch EC2 instances accordingly
- AWS ParallelCluster
 - Open-source cluster management tool to deploy HPC on AWS
 - Configure with text files
 - Automate creation of VPC, Subnet, cluster type and instance types
 - Ability to enable EFA on the cluster (improves network performance)

EC2 Instance High Availability

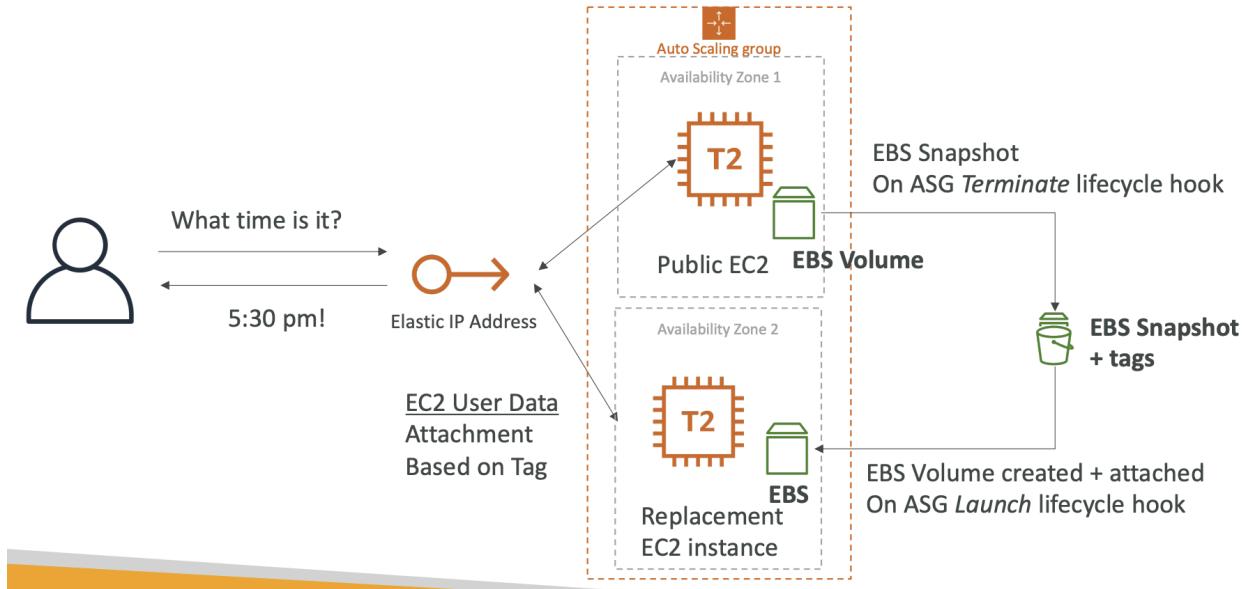
Creating a highly available EC2 instance



Creating a highly available EC2 instance
With an Auto Scaling Group



Creating a highly available EC2 instance With ASG + EBS



Section 30: Other Services

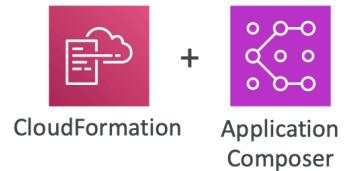
CloudFormation Intro

- Declarative way to outline AWS infrastructure via code, created in the right order, declaratively
- IaaC
 - No resources manually created, version control with code
- Cost: all resources within the stack is tagged with an identifier to see cost
 - Can estimate costs
 - Can create and destroy
- Productivity:
 - Destroy and recreate infrastructure on the fly
- Separation of concern: create many stacks for many apps and many layers
- Don't reinvent the wheel with existing templates and documentation

CF + Application Composer

CloudFormation + Application Composer

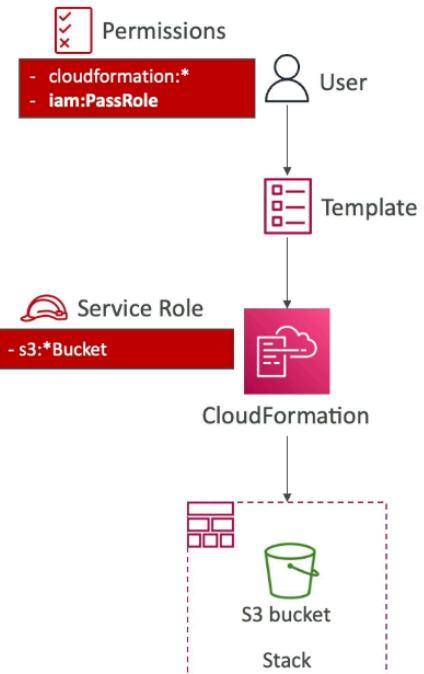
- Example: WordPress CloudFormation Stack
- We can see all the resources
- We can see the relations between the components



CF – Service Roll

CloudFormation – Service Role

- IAM role that allows CloudFormation to create/update/delete stack resources on your behalf
- Give ability to users to create/update/delete the stack resources even if they don't have permissions to work with the resources in the stack
- Use cases:
 - You want to achieve the least privilege principle
 - But you don't want to give the user all the required permissions to create the stack resources
- User must have `iam:PassRole` permissions



- IAM role that allows CF to create / update / delete stack resources on behalf
 - Gives users ability to update resources even without permissions to work with the resources in the stack

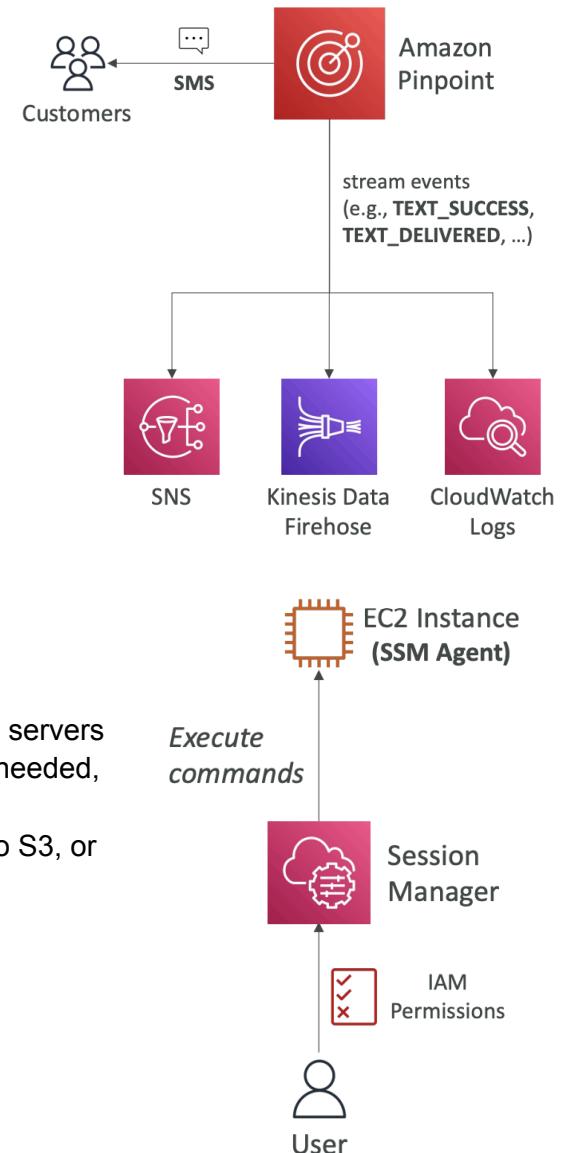
- User must have iam:PassRole Permissions

AWS Simple Email Service (SES)

- Send emails and ability receive email
 - Integrates with S3, SNS, Lambda and IAM for allowing to send emails
 - Allows in/outbound emails
 - Send emails using app via Console or API or SMTP
- Reputation dashboard, performance insights, anti spam feedback
 - Statistics such as email deliveries, bounces, email open...
- Supports DomainKeys Identified Mail (DKIM) and Sender Policy Framework (SPF)
- Flexible IP deployment: shared, dedicated, and customer owned IP
- Use case: transactional, marketing, bulk email communications

Amazon Pinpoint

- Scalable 2 way (outbound/inbound) marketing communications service that supports email, SMS, push, voice, and in app messaging
 - Can segment and personalize messages
 - Can receive replies
 - Scales to billions of messages
- Use cases: run campaigns by sending marketing, bulk, transactional SMS messages
- vs SNS or SES?
 - SNS & SES you manage each message's audience, content and delivery schedule
 - Pinpoint can create message templates, delivery schedules, highly-targeted segments and full campaigns

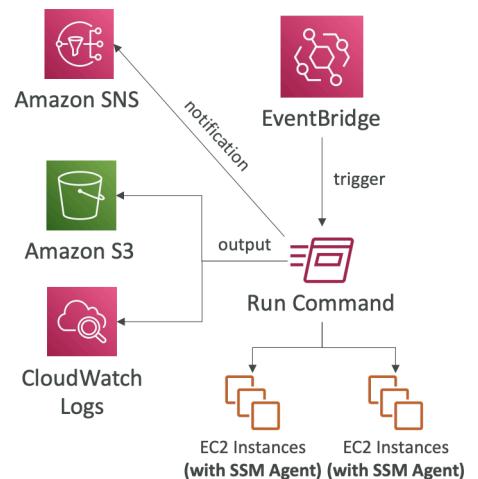


Systems Manager – SSM Session Manager

- Allows to start a secure shell to EC2 and on premise servers without SSH, bastion host, or SSH keys (no port 22 needed, better security)
 - Supports all OS, can send session log data to S3, or CloudWatch Logs

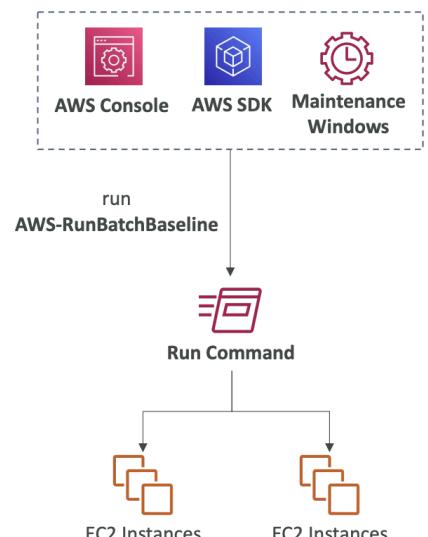
SSM – Run Command

- Execute document (script) or just run a command
 - Run command across multiple instances (using resource groups)
 - No need for SSH
 - Command output sent to S3 or CloudWatch Logs
 - Send notifications to SNS about command status
 - Integrated with IAM & CloudTrail
 - Can be invoked via EventBridge



SSM – Patch Manager

- Automated patching managed instances
 - OS updates, app updates, security updates
 - Supports EC2 instances, on premise; all OS
- Patch on demand or schedule via Maintenance Windows
- Scan instances and generate patch compliance report (missing patches)



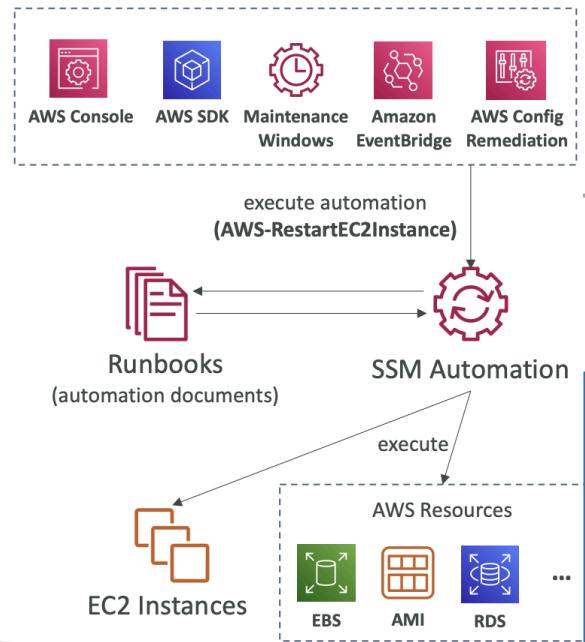
SSM – Maintenance Windows



- Schedule for when to perform actions on instances
 - Ex: OS patching, update drivers, install software...
- Contains:
 - Schedule, duration, set of registered instances, set of registered tasks

SSM – Automation

- Simplified common maintenance and deployment tasks of EC2 instances and other AWS resources
 - Ex: restart instances, create AMI, EBS snapshot
- Automation Runbook: SSM Documents for predefined actions on EC2 instances or AWS resources



Cost Explorer

- Visualize, understand and manage AWS costs and usage over time
- Create custom reports that analyze cost and usage data
 - Analyze data at high level: total costs and usage across all accounts
 - Hourly, monthly, resource level granularity
- Choose optimal savings plan
- Forecast usage up to 12 months based on previous usage

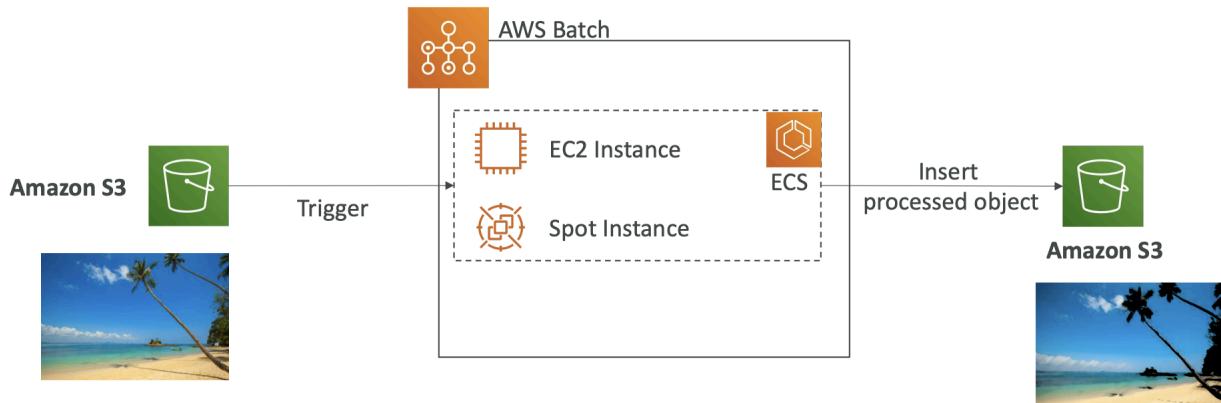
AWS Cost Anomaly Detection

- Continuously monitor cost and usage using ML to detect unusual spends
 - Learns unique historical patterns to detect one time cost spike and / or continuous cost increases (no threshold needed)
- Monitor AWS services, member accounts, cost allocation tags, cost categories
- Sends anomaly detection report with root cause analysis
 - Notified with individual alerts or daily / weekly summary (SNS)

AWS Batch

- Fully managed batch processing at any scale; can efficiently run 100,000s of computing batch jobs on AWS
- Batch is a job with a start and end (not continuous)
 - Batch will dynamically launch EC2 instances or Spot Instances
 - AWS Batch provisions right amount of compute / memory
 - Just submit or schedule batch jobs and AWS Batch does the rest
- Batch jobs defined as docker images and run on ECS
- Helpful for cost optimizations and focusing less on infrastructure

AWS Batch – Simplified Example



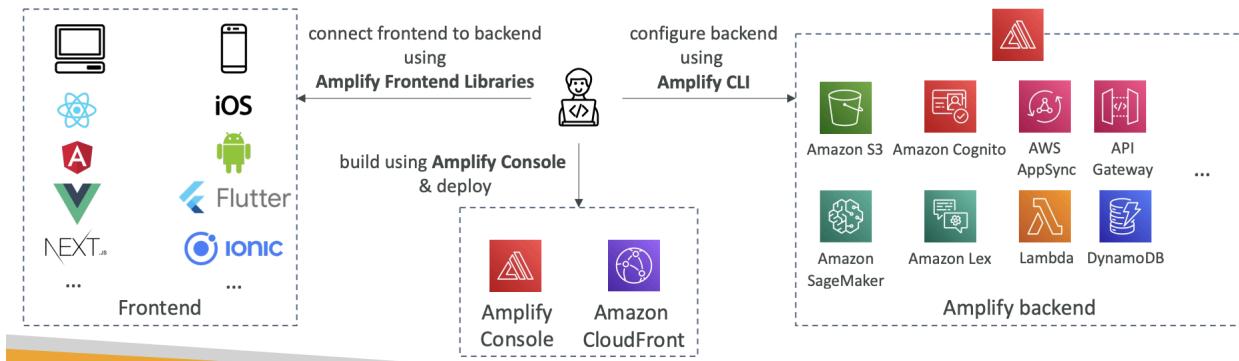
Batch vs Lambda

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

Amazon AppFlow

- Fully managed integration service that enables you to securely transfer data between SaaS apps and AWS
 - Source: Salesforce, SAP, Slack, ServiceNow...
 - Destination: S3, Redshift, or non-AWS like Snowflake and Salesforce
 - Frequency: scheduled, response to events, or on demand
 - Data transformation capabilities like filtering and validation
 - Encrypted over public internet or privately over AWS PrivateLink
 - No time writing integrations and leverage APIs

AWS Amplify



- Create mobile and web apps → Elastic beanstalk for mobile and web apps powered by CloudFormation
- Relies on DynamoDB, AppSync, Cognito, S3 as backend with any frontend library

Section 31: Whitepapers and Architectures

Well Architected Framework

Well Architected Framework General Guiding Principles

- <https://aws.amazon.com/architecture/well-architected>
- Stop guessing your capacity needs
- Test systems at production scale
- Automate to make architectural experimentation easier
- Allow for evolutionary architectures
 - Design based on changing requirements
- Drive architectures using data
- Improve through game days
 - Simulate applications for flash sale days

Well Architected Framework

6 Pillars

- 1) Operational Excellence
 - 2) Security
 - 3) Reliability
 - 4) Performance Efficiency
 - 5) Cost Optimization
 - 6) Sustainability
-
- They are not something to balance, or trade-offs, they're a synergy

AWS Well Architected Tool

AWS Well-Architected Tool



- Free tool to review your architectures against the 6 pillars Well-Architected Framework and adopt architectural best practices
- How does it work?
 - Select your workload and answer questions
 - Review your answers against the 6 pillars
 - Obtain advice: get videos and documentations, generate a report, see the results in a dashboard
- Let's have a look: <https://console.aws.amazon.com/wellarchitected>

Name	Overall status	High risks	Medium risks	Improvement status	Last updated
Internal Employee Portal	Answered	13	2	None	Nov 24, 2018 3:40 PM UTC-8
Mobile app - Android	Answered	9	1	None	Nov 24, 2018 3:43 PM UTC-8
Mobile app - iOS	Answered	0	1	None	Nov 24, 2018 3:49 PM UTC-8
Retail Website- EU	Unanswered	0	0	None	Nov 24, 2018 3:52 PM UTC-8
Retail Website- North America	Unanswered	0	0	None	Nov 24, 2018 3:19 PM UTC-8

AWS Trusted Advisor

Trusted Advisor

- No need to install anything – high level AWS account assessment
- Analyze your AWS accounts and provides recommendation on 6 categories:
 - Cost optimization
 - Performance
 - Security
 - Fault tolerance
 - Service limits
 - Operational Excellence
- Business & Enterprise Support plan
 - Full Set of Checks
 - Programmatic Access using AWS Support API



Checks

▶ Amazon EBS Public Snapshots

Checks the permission settings for your Amazon Elastic Block Store snapshots. 0 EBS snapshots are marked as public.

▶ Amazon RDS Public Snapshots

Checks the permission settings for your Amazon Relational Database Service snapshots. 0 RDS snapshots are marked as public.

▶ IAM Use

This check is intended to discourage the use of root access keys. At least one IAM user has been created for this account.