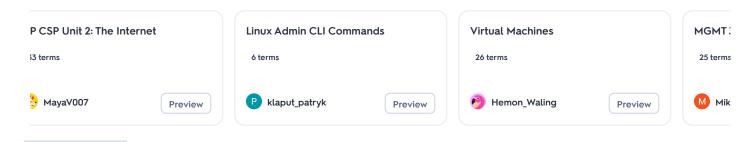
AWS Certified Solutions Architect

Students also viewed



Terms in this set (554)

	Each availability zone (AZ) is one or more
AWS Availability Zones	discrete data centers with redundant power,
	networking, and connectivity
	They're separate from each other, so that
	they're isolated from disasters
	Root account created by default, shouldn't be used or shared
AAA 11 9 Curawaa	· Users are people within your organization, and can be grouped
IAM: Users & Groups	Groups only contain users, not other groups
	· Users don't have to belong to a group, and user can belong to multiple groups
	-Consists of
	 Version: policy language version, always include "2012 -10 - 17"
	· Id: an identifier for the policy (optional)
	Statement: one or more individual statements (required)
IAM Policies Structure	-Statements consists of
	Sid: an identifier for the statement (optional)
	Effect: whether the statement allows or denies access (Allow, Deny)
	Principal: account/user/role to which this policy applied to
	• Action: list of actions this policy allows or denies •Resource: list of resources to which
	the actions applied to
	Condition: conditions for when this policy is in effect (optional)
	- IAM Credentials Report (account-level) • a report that lists all your account's users
	and the status of their various credentials
AM Security Tools	- IAM Access Advisor (user-level) · Access advisor shows the service permissions
	granted to a user and when those services were last accessed. • You can use this
	information to revise your policies.

	Don't use the root account except for AWS account setup
	• One physical user = One AWS user
	· Assign users to groups and assign permissions to groups
	· Create a strong password policy
IAM Guidelines and Best Practices	· Use and enforce the use of Multi Factor Authentication (MFA)
	· Create and use Roles for giving permissions to AWS services
	Use Access Keys for Programmatic Access (CLI / SDK)
	Audit permissions of your account with the IAM Credentials Report
	Never share IAM users & Access Keys
	It is possible to bootstrap our instances using an EC2 User data script.
	· bootstrapping means launching commands when a machine starts · That script is only
	run once at the instance first start
EC2 User Data	• EC2 user data is used to automate boot tasks such as:
	• Installing updates • Installing software • Downloading common files from the internet •
	Anything you can think of
	• The EC2 User Data Script runs with the root user
	AWS has the following naming convention:
	m5.2xlarge
EC2 Instance Types - Overview	• m: instance class
Loz matarice types overview	• 5: generation (AWS improves them over time)
	• 2xlarge: size within the instance class
	-
EC2 - Instance Type: General purpose	Great for a diversity of workloads such as web servers or code repositories
	Great for compute-intensive tasks that require high performance
	processors:
	Batch processing workloads
EC2 Instance Types - Compute Optimized	Media transcoding
LCZ instance Types Compute Optimized	High performance web servers
	· High performance computing (HPC)
	Scientific modeling & machine learning
	Dedicated gaming servers
	Fast performance for workloads that process large data sets in memory
	• Use cases:
EC2 Instance Types - Memory Optimized	High performance, relational/non-relational databases
LC2 Instance Types - Memory Optimized	Distributed web scale cache stores
	 In-memory databases optimized for BI (business intelligence)
	 Applications performing real-time processing of big unstructured data
	Great for storage-intensive tasks that require high, sequential read and write
	access to large data sets on local storage
	· Use cases:
	High frequency online transaction processing (OLTP) systems
EC2 Instance Types - Storage Optimized	• Relational & NoSQL databases
	Cache for in-memory databases (for example, Redis)
	Data warehousing applications
	Distributed file systems
	-Security Groups are fundamental of network security in AWS
	Secondy Groups are formamental of fletwork seconds in Avra
	-They control how traffic is allowed into or out of our EC2 Instances
Introduction to Security Groups	
	-Security groups only contain allow rules
	-Security groups rules can reference by IP or by security group
	, , , , , , , , , , , , , , , , , , ,

	Con he attached to multiple instances
Security Groups - Good to know	Can be attached to multiple instances Lacked down to a region (VRC combination)
	Locked down to a region / VPC combination
	Does live "outside" the EC2 - if traffic is blocked the EC2 instance won't see it
	It's good to maintain one separate security group for SSH access
	If your application is not accessible (time out), then it's a security group issue
	· If your application gives a "connection refused" error, then it's an application
	error or it's not launched
	All inbound traffic is blocked by default
	All outbound traffic is authorised by default
	· 22 = SSH (Secure Shell) - log into a Linux instance
	· 21 = FTP (File Transfer Protocol) - upload files into a file share
Classic Ports to know	· 22 = SFTP (Secure File Transfer Protocol) - upload files using SSH
Classic Folks to know	• 80 = HTTP - access unsecured websites
	• 443 = HTTPS - access secured websites
	• 3389 = RDP (Remote Desktop Protocol) - log into a Windows instance
	On-Demand Instances - short workload, predictable pricing, pay by second
	· Reserved (1 & 3 years)
	· Reserved Instances - long workloads
	Convertible Reserved Instances - long workloads with flexible instances
EC2 Instances Purchasing Options	· Savings Plans (1 & 3 years) -commitment to an amount of usage, long workload
	Spot Instances - short workloads, cheap, can lose instances (less reliable)
	Dedicated Hosts - book an entire physical server, control instance placement
	Dedicated Instances - no other customers will share your hardware
	· Capacity Reservations - reserve capacity in a specific AZ for any duration
	• Up to 72% discount compared to On-demand
	You reserve a specific instance attributes (Instance Type, Region, Tenancy, OS)
	Reservation Period - 1 year (+discount) or 3 years (+++discount)
	Payment Options - No Upfront (+), Partial Upfront (++), All Upfront (+++)
	Reserved Instance's Scope - Regional or Zonal (reserve capacity in an AZ)
EC2 Reserved Instances	Recommended for steady-state usage applications (think database)
	You can buy and sell in the Reserved Instance Marketplace
	Convertible Reserved Instance
	Can change the EC2 instance type, instance family, OS, scope and tenancy
	• Up to 66% discount
	• Get a discount based on long-term usage (up to 72% - same as RIs)
	• Commit to a certain type of usage (\$10/hour for 1 or 3 years)
	Usage beyond EC2 Savings Plans is billed at the On-Demand price
EC2 Savings Plans	Locked to a specific instance family & AWS region (e.g., M5 in us-east-1)
	Flexible across:
	Instance Size (e.g., m5.xlarge, m5.2xlarge)
	· OS (e.g., Linux, Windows)
	Tenancy (Host, Dedicated, Default)
	Tenancy (11034, Dedicated, Default)

	Can get a discount of up to 90% compared to On-demand
	Instances that you can "lose" at any point of time if your max price is less than the current spot price
	The MOST cost-efficient instances in AWS
EC2 Spot Instances	Useful for workloads that are resilient to failure
	Batch jobs
	Data analysis
	· Image processing
	Any distributed workloads
	Workloads with a flexible start and end time
	Not suitable for critical jobs or databases
	A physical server with EC2 instance capacity fully dedicated to your use
	· Allows you address compliance requirements and use your existing serverbound
	software licenses (per-socket, per-core, pe—VM software licenses)
	Purchasing Options:
	On-demand - pay per second for active Dedicated Host
	Reserved - 1 or 3 years (No Upfront, Partial Upfront, All Upfront)
	• The most expensive option
EC2 Dedicated Hosts	Useful for software that have complicated licensing model (BYOL - Bring Your
	Own License)
	Or for companies that have strong regulatory or compliance needs
	Instances run on hardware that's
	dedicated to you
	May share hardware with other
	instances in same account
	No control over instance placement (can move hardware after Stop / Start)
	(can move naroware arter stop / start)
	- Reserve On-Demand instances capacity in a specific AZ for any duration
	- You always have access to EC2 capacity when you need it
	- No time commitment (create/cancel anytime), no billing discounts
EC2 Capacity Reservations	- Combine with Regional Reserved Instances and Savings Plans to benefit from billing discounts
	- You're charged at On-Demand rate whether you run instances or not
	- Suitable for short-term, uninterrupted workloads that needs to be in a specific AZ
	Solitable for short term, offinterropted workloads that freeds to be in a specific AZ
	On demand: coming and staying in resort
	whenever we like, we pay the full price
	· Reserved: like planning ahead and if we plan to
	stay for a long time, we may get a good discount.
	Savings Plans: pay a certain amount per hour for
	certain period and stay in any room type (e.g.,
Which purchasing option is right for me?	King, Suite, Sea View,)
	Spot instances: the hotel allows people to bid for the ampty rooms and the highest hidder keeps the
	the empty rooms and the highest bidder keeps the
	rooms. You can get kicked out at any time
	Dedicated Hosts: We book an entire building of the resort
	· Capacity Reservations: you book a room for a
	period with full price even you don't stay in it
	paration of price of on you don't diag in it

EC2 Spot Instance Requests	Can get a discount of up to 90% compared to On-Demand Define max spot price and get the instance while current spot price < max The hourly spot price varies based on offer and capacity If the current spot price > your max price you can choose to stop or terminate your instance with a 2 minutes grace period. Other strategy: Spot Block- "block" spot instance during a specified time frame (1 to 6 hours) without interruptions • In rare situations, the instance may be reclaimed Used for batch jobs, data analysis, or workloads that are resilient to failures. Not great for critical jobs or databases
Spot Fleets	Spot Fleets = set of Spot Instances + (optional) On-Demand Instances The Spot Fleet will try to meet the target capacity with price constraints Define possible launch pools: instance type (m5.large), OS, Availability Zone Can have multiple launch pools, so that the fleet can choose Spot Fleet stops launching instances when reaching capacity or max cost Strategies to allocate Spot Instances: lowestPrice: from the pool with the lowest price (cost optimization, short workload) diversified: distributed across all pools (great for availability, long workloads) capacityOptimized: pool with the optimal capacity for the number of instances priceCapacityOptimized (recommended): pools with highest capacity available, then select the pool with the lowest price (best choice for most workloads) Spot Fleets allow us to automatically request Spot Instances with the lowest price
Elastic IPs	 When you stop and then start an EC2 instance, it can change its public IP. If you need to have a fixed public IP for your instance, you need an Elastic IP An Elastic IP is a public IPv4 IP you own as long as you don't delete it You can attach it to one instance at a time With an Elastic IP address, you can mask the failure of an instance or software by rapidly remapping the address to another instance in your account. You can only have 5 Elastic IP in your account (you can ask AWS to increase that).
EC2 Hibernate	On start, the following happens: First start: the OS boots & the EC2 User Data script is run Following starts: the OS boots up Then your application starts, caches get warmed up, and that can take time! 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
EC2 Hibernate - Good to know	 Supported Instance Families - C3, C4, C5, I3, M3, M4, R3, R4, T2, T3, Instance RAM Size - must be less than 150 GB. Instance Size - not supported for bare metal instances. AMI - Amazon Linux 2, Linux AMI, Ubuntu, RHEL, CentOS & Windows Root Volume - must be EBS, encrypted, not instance store, and large Available for On-Demand, Reserved and Spot Instances An instance can NOT be hibernated more than 60 days

	 It's a network drive (i.e. not a physical drive) It uses the network to communicate the instance, which means there might be a bit of latency It can be detached from an EC2 instance and attached to another one quickly
EBS Volume	 It's locked to an Availability Zone (AZ) An EBS Volume in us-east-la cannot be attached to us-east-lb To move a volume across, you first need to snapshot it Have a provisioned capacity (size in GBs, and IOPS) You get billed for all the provisioned capacity
	· You can increase the capacity of the drive over time
EBS - Delete on Termination attribute	 Controls the EBS behaviour when an EC2 instance terminates By default, the root EBS volume is deleted (attribute enabled) By default, any other attached EBS volume is not deleted (attribute disabled) This can be controlled by the AWS console / AWS CLI Use case: preserve root volume when instance is terminated
EBS Snapshots	 Make a backup (snapshot) of your EBS volume at a point in time Not necessary to detach volume to do snapshot, but recommended Can copy snapshots across AZ or Region
EBS Snapshots Features	EBS Snapshot Archive Move a Snapshot to an "archive tier" that is 75% cheaper Takes within 24 to 72 hours for restoring the archive Recycle Bin for EBS Snapshots Setup rules to retain deleted snapshots so you can recover them after an accidental deletion
	 Specify retention (from 1 day to 1 year) Fast Snapshot Restore (FSR) Force full initialization of snapshot to have no latency on the first use (\$\$\$)
AMI Overview	 AMI = Amazon Machine Image AMI are a customization of an EC2 instance You add your own software, configuration, operating system, monitoring Faster boot / configuration time because all your software is pre-packaged AMI are built for a specific region (and can be copied across regions) You can launch EC2 instances from: A Public AMI: AWS provided Your own AMI: you make and maintain them yourself An AWS Marketplace AMI: an AMI someone else made (and potentially sells)
AMI Process (from an EC2 instance)	Start an EC2 instance and customize it Stop the instance (for data integrity) Build an AMI - this will also create EBS snapshots Launch instances from other AMIs
EC2 Instance Store	Better I/O performance EC2 Instance Store lose their storage if they're stopped (ephemeral) Good for buffer / cache / scratch data / temporary content Risk of data loss if hardware fails Backups and Replication are your responsibility
EBS Volume Types	Only gp2/gp3 and io1/io2 can be used as boot volumes
EBS Volume Types Use cases General Purpose SSD	Cost effective storage, low-latency • System boot volumes, Virtual desktops, Development and test environments • 1 GiB - 16 TiB

	Critical business applications with sustained IOPS performance
EBS Volume Types Use cases	Or applications that need more than 16,000 IOPS
Provisioned IOPS (PIOPS) SSD	Great for databases workloads (sensitive to storage perf and consistency)
	Supports EBS Multi-attach
	· Attach the same EBS volume to multiple EC2
	instances in the same AZ
	Each instance has full read & write permissions
	to the high-performance volume
	• Use case:
EBS Multi-Attach - io1/io2 family	Achieve higher application availability in clustered
	Linux applications (ex: Teradata)
	Applications must manage concurrent write operations
	• Up to 16 EC2 Instances at a time
	· Must use a file system that's cluster-aware (not
	XFS, EXT4, etc)
	When you create an encrypted EBS volume, you get the following:
	Data at rest is encrypted inside the volume
	· All the data in flight moving between the instance and the volume is encrypted
	• All snapshots are encrypted
	• All volumes created from the snapshot
EBS Encryption	Encryption and decryption are handled transparently (you have nothing to
	do)
	• Encryption has a minimal impact on latency
	• EBS Encryption leverages keys from KMS (AES-256)
	Copying an unencrypted snapshot allows encryption Spanshots of encrypted yellumos are encrypted.
	Snapshots of encrypted volumes are encrypted
	· Create an EBS snapshot of the volume
Encryption: encrypt an unencrypted EBS	Encrypt the EBS snapshot (using copy)
volume	· Create new ebs volume from the snapshot (the volume will also be
	encrypted)
	Now you can attach the encrypted volume to the original instance
	· Managed NFS (network file system) that can be mounted on many EC2
	• EFS works with EC2 instances in multi-AZ
Amazon EFS - Elastic File System	· Highly available, scalable, expensive (3x gp2), pay per use
	Use cases: content management, web serving, data sharing, Wordpress
	Use cases: content management, web serving, data sharing, Wordpress
	Performance Mode (set at EFS creation time)
	· General Purpose (default) - latency-sensitive use cases (web server, CMS, etc)
	Max I/O - higher latency, throughput, highly parallel (big data, media processing)
	· Throughput Mode
	Bursting - 1TB = 50MiB/s + burst of up to 100MiB/s
EFS - Performance & Storage Classes	• Provisioned - set your throughput regardless of storage size, ex: 1 GiB/s for 1 TB
	storage
	· Elastic - automatically scales throughput up or down based on your workloads
	· Up to 3GiB/s for reads and 1GiB/s for writes
	· Used for unpredictable workloads
	Storage Tiers (lifecycle management feature
	- move file after N days)
	Standard: for frequently accessed files
EFS - Storage Classes	• Infrequent access (EFS-IA): cost to retrieve files,
	lower price to store. Enable EFS-IA with a Lifecycle
	Policy
	· · · · · · · · · · · · · · · · · · ·

EBS vs EFS - Elastic Block Storage	EBS volumes	
	one instance (except multi-attach io1/io2)	
	· are locked at the Availability Zone (AZ) level	
	· gp2: IO increases if the disk size increases	
	· iol: can increase IO independently	
	• To migrate an EBS volume across AZ	
	· Take a snapshot	
	• Restore the snapshot to another AZ	
	• EBS backups use IO and you shouldn't run	
	them while your application is handling a lot	
	of traffic	
	Root EBS Volumes of instances get	
	terminated by default if the EC2 instance	
	gets terminated. (you can disable that)	
	Mounting 100s of instances across AZ	
	• EFS share website files (WordPress)	
	Only for Linux Instances (POSIX)	
EBS vs EFS - Elastic File System	• EFS has a higher price point than EBS	
	Can leverage EFS-IA for cost savings	
	Remember: EFS vs Instance Store	
	Scalability means that an application / system can handle greater loads	
	by adapting.	
Scalability & High Availability	• There are two kinds of scalability:	
3	· Vertical Scalability	
	Horizontal Scalability (= elasticity)	
	Scalability is linked but different to High Availability	
	· High Availability usually goes hand in	
	hand with horizontal scaling	
Lligh Availability	· High availability means running your	
High Availability	application / system in at least 2 data	
	centers (== Availability Zones)	
	• The goal of high availability	
	Load Balances are servers that forward traffic to multiple	
	servers (e.g., EC2 instances) downstream	
	Spread load across multiple downstream instances	
	• Expose a single point of access (DNS) to your application	
	Seamlessly handle failures of downstream instances	
What is load balancing?	Do regular health checks to your instances	
	Provide SSL termination (HTTPS) for your websites	
	• Enforce stickiness with cookies	
	· High availability across zones	
	Separate public traffic from private traffic	
	An Elastic Load Balancer is a managed load balancer	
	AWS guarantees that it will be working	
	AWS takes care of upgrades, maintenance, high availability	
	AWS provides only a few configuration knobs	
	It costs less to setup your own load balancer but it will be a lot more effort	
Why use an Elastic Load Balancer?	on your end	
	It is integrated with many AWS offerings / services	
	• EC2, EC2 Auto Scaling Groups, Amazon ECS	
	AWS Certificate Manager (ACM), CloudWatch	
	Route 53, AWS WAF, AWS Global Accelerator	
	noste 30, And har, And Global Accelerator	

	· Health Checks are crucial for Load Balancers
Health Checks	· They enable the load balancer to know if instances it forwards traffic to
	are available to reply to requests
	• The health check is done on a port and a route (/health is common)
	· If the response is not 200 (OK), then the instance is unhealthy
	Application load balancers is Layer 7 (HTTP)
	Load balancing to multiple HTTP applications across machines
	(target groups)
	Load balancing to multiple applications on the same machine
	(ex: containers)
	Routing tables to different target groups:
	Routing based on path in URL (example.com/users & example.com/posts)
	Routing based on hostname in URL (one.example.com & other.example.com)
	Routing based on Query String, Headers
	(example.com/users?id=123ℴ=false)
Application Load Balancer (v2)	ALB are a great fit for micro services & container-based application
	(example: Docker & Amazon ECS)
	Has a port mapping feature to redirect to a dynamic port in ECS
	In comparison, we'd need multiple Classic Load Balancer per application
	EC2 instances (can be managed by an Auto Scaling Group) - HTTP
	• ECS tasks (managed by ECS itself) - HTTP
	Lambda functions - HTTP request is translated into a JSON event
	• IP Addresses - must be private IPs
	• ALB can route to multiple target groups
	Health checks are at the target group level
	Fixed hostname (XXX.region.elb.amazonaws.com)
Application Load Balancer (v2)	• The application servers don't see the IP of the client directly
Good to Know	• The true IP of the client is inserted in the header X-Forwarded-For
	We can also get Port (X-Forwarded-Port) and proto (X-Forwarded-Proto)
	Network load balancers (Layer 4) allow to:
	Forward TCP & UDP traffic to your instances
	Handle millions of request per seconds
Network Load Balancer (v2)	 Less latency ~100 ms (vs 400 ms for ALB)
	NLB has one static IP per AZ, and supports assigning Elastic IP
	(helpful for whitelisting specific IP)
	Not included in the AWS free tier
	Deploy, scale, and manage a fleet of 3rd party
	network virtual appliances in AWS
	• Example: Firewalls, Intrusion Detection and
	Prevention Systems, Deep Packet Inspection
	Systems, payload manipulation,
Gateway Load Balancer	· Operates at Layer 3 (Network Layer) - IP
	Packets
	Combines the following functions:
	· Transparent Network Gateway - single
	entry/exit for all traffic
	· Load Balancer - distributes traffic to your virtual
	appliances
	· Uses the GENEVE protocol on port 6081

Sticky Sessions (Session Affinity)	It is possible to implement stickiness so that the same client is always redirected to the same instance behind a load balancer This works for Classic Load Balancer, Application Load Balancer, and Network Load Balancer For both CLB & ALB, the "cookie" used for stickiness has an expiration date you control Application-based Cookies Custom cookie Generated by the target Can include any custom attributes required by the application Cookie name must be specified individually for each target group Don't use AWSALB, AWSALBAPP, or AWSALBTG (reserved for use by the ELB) Application cookie Generated by the load balancer Cookie name is AWSALBAPP Duration-based Cookies Cookie generated by the load balancer
	Cookie name is AWSALB for ALB, AWSELB for CLB
Cross-Zone Load Balancing	With Cross Zone Load Balancing: each load balancer instance distributes evenly across all registered instances in all AZ Application Load Balancer • Enabled by default (can be disabled at the Target Group level) • No charges for inter AZ data • Network Load Balancer & Gateway Load Balancer • Disabled by default • You pay charges (\$) for inter AZ data if enabled • Classic Load Balancer • Disabled by default • No charges for inter AZ data if enabled
	-An SSL Certificate allows traffic between your clients and your load balancer to be encrypted in transit (in-flight encryption) -SSL refers to Secure Sockets Layer, used to encrypt connections
SSL/TLS - Basics	-TLS refers to Transport Layer Security, which is newer version -Nowadays, TLS certificates are mainly used, but people still refer as SSL
	-Public SSL certificates are issued by Certificate Authorities (CA)
	-Comodo, Symantec, GoDaddy, GlobalSign, Digicert, Letsencrypt, etc
	-SSL certificates have an expiration date (you set) and must be renewed
Load Balancer - SSL Certificates	The load balancer uses an X.509 certificate (SSL/TLS server certificate) · You can manage certificates using ACM (AWS Certificate Manager) · You can create upload your own certificates alternatively

SSL - Server Name Indication (SNI)	SNI solves the problem of loading multiple SSL certificates onto one web server (to serve multiple websites) It's a "newer" protocol, and requires the client to indicate the hostname of the target server in the initial SSL handshake The server will then find the correct certificate, or return the default one Note: Only works for ALB & NLB (newer generation), CloudFront
Elastic Load Balancers - SSL Certificates	Classic Load Balancer (vI) Support only one SSL certificate Must use multiple CLB for multiple hostname with multiple SSL certificates Application Load Balancer (v2) Supports multiple listeners with multiple SSL certificates Uses Server Name Indication (SNI) to make it work Network Load Balancer (v2) Supports multiple listeners with multiple SSL certificates Uses Server Name Indication (SNI) to make it work
Connection Draining	 Time to complete "in-flight requests" while the instance is de-registering or unhealthy Stops sending new requests to the EC2 instance which is de-registering Between 1 to 3600 seconds (default: 300 seconds) Can be disabled (set value to 0)
Auto Scaling Group Attributes	A Launch Template (older "Launch Configurations" are deprecated) • AMI + Instance Type • EC2 User Data • EBS Volumes • Security Groups • SSH Key Pair • IAM Roles for your EC2 Instances • Network + Subnets Information • Load Balancer Information • Min Size / Max Size / Initial Capacity • Scaling Policies
Auto Scaling Groups - Dynamic Scaling Policies	Target Tracking Scaling Most simple and easy to set-up Example: I want the average ASG CPU to stay at around 40% Simple / Step Scaling When a CloudWatch alarm is triggered (example CPU > 70%), then add 2 units When a CloudWatch alarm is triggered (example CPU < 30%), then remove 1 Scheduled Actions Anticipate a scaling based on known usage patterns Example: increase the min capacity to 10 at 5 pm on Fridays Predictive scaling: continuously forecast load and schedule scaling ahead
Good metrics to scale on	 CPUUtilization: Average CPU utilization across your instances RequestCountPerTarget: to make sure the number of requests per EC2 instances is stable Average Network In / Out (if you're application is network bound) Any custom metric (that you push using CloudWatch)

	After a scaling activity happens, you are in
	the cooldown period (default 300
	seconds)
	During the cooldown period, the ASG will
Auto Scaling Groups - Scaling Cooldowns	not launch or terminate additional
l l l l l l l l l l l l l l l l l l l	instances (to allow for metrics to stabilize)
	· Advice: Use a ready-to-use AMI to reduce
	configuration time in order to be serving
	request fasters and reduce the cooldown
	period
	RDS is a managed service:
	· Automated provisioning, OS patching
	· Continuous backups and restore to specific timestamp (Point in Time Restore)!
	Monitoring dashboards
Advantage over using RDS versus deploying	Read replicas for improved read performance
DB on EC2	Multi AZ setup for DR (Disaster Recovery)
	Maintenance windows for upgrades
	Scaling capability (vertical and horizontal)
	· Storage backed by EBS (gp2 or io1)
	Helps you increase storage on your RDS DB instance
	dynamically
	· When RDS detects you are running out of free database
	storage, it scales automatically
	Avoid manually scaling your database storage
DDC Chamana Alaba Caalina	· You have to set Maximum Storage Threshold (maximum
RDS - Storage Auto Scaling	limit for DB storage)
	Automatically modify storage if:
	• Free storage is less than 10% of allocated storage
	· Low-storage lasts at least 5 minutes
	6 hours have passed since last modification
	Useful for applications with unpredictable workloads
	-SYNC replication
	-One DNS name - automatic app failover to standby
	-Increase availability
	-Failover in case of loss of AZ, loss of network, instance or storage failure
RDS Multi AZ (Disaster Recovery)	-No manual intervention in apps
	-Not used for scaling
	-Multi-AZ replication is free
	-Note: The Read Replicas be setup as Multi AZ for Disaster Recovery (DR)
	2. 3. 2. 3. 2. 3. 2. 3. 2. 3. 2. 3. 2. 3. 3. 3. 3. 3. 3. 3. 3. 3. 3. 3. 3. 3.

- Access the underlying EC2 instance using SSH or SSM Session Manager - De-activate Automation Mode to perform your customization, better to take a DB snapshot before RDS vs RDS-custom: - RDS entire database and the OS to be managed by AWS - RDS Custom: full admin access to the underlying OS and the database Aurora is a proprietary technology from AWS (not open sourced) - Postgres and MySQL are both supported as Aurora D8 (that means your drivers will work as if Aurora was a Postgres or MySQL database) - Aurora is 'AWS cloud optimized' and claims 5x performance improvement over MySQL on RDS, over 5x the performance of Postgres on RDS - Aurora storage automatically grows in increments of IDGB, up to I28 TB Aurora can have up to 15 replicas and the replication process is faster than MySQL (sub ID in septica lag) - Fallover in Aurora is instantaneous. It's HA (High Availability) native. 6 copies of your data across 3 A.Z: - 4 copies out of 6 needed for writes - 3 copies out of 1 need for reads - Self healing with peer-10-peer replication - Storage is striped across IDDs of volumes - One Aurora instance takes writes (master) - Automated fallower for master in less than 30 seconds - Master + up to IS Aurora Read Replicas serve reads - Support for Cross Region Replication - Automated fallover - Backup and Recovery - Bac		
Custom: access to the underlying database in AWS Custom: access to the underlying database and OS so you can: - Configure settings - Install praches - Enable native features - Access the underlying EC2 instance using SSH or SSM Session Manager - De-activate Automation Mode to perform your customization, better to take a DB snapshot before RDS vs RDS-custom - RDS entire database and the OS to be managed by AWS - RDS Custom full admin access to the underlying OS and the database Aurora is a proprietary technology from AWS (not open sourced) - Postgres and HySCII are both supported as Autora DB (that means your drivers will work as if Aurora was a Postgres or MySCII database) - Aurora is AWS could opinized and claims Sx performance improvement over MySCII on RDS, over 3x the performance of Postgres on RDS - Aurora as have up to 15 replicas and the replication process is faster than MySCII (dub 10 ms replica log) - Fallover in Aurora is intantanarous. Its HA (High Availability) native. d copies of your data across 3 AZ: - 4 copies out of an eeded for writes - 3 copies out of an eeded for writes - 3 copies out of an eeded for writes - 3 copies out of an eeded for writes - 3 copies out of an eeded for eads - Self healing with peer-to-peer replication - Storage is striped across 100s of volumes - Aurora and a Pusturance takes writes (master) - Automated fallover for master in less than 30 seconds - Master- up to 15 Aurora Read Replicas - Service and a second to the self- of the striped across 100s of volumes - Automated fallover for master in less than 30 seconds - Master- up to 15 Aurora Read Replicas - Service and the replication - Aurora - Custom Endpoints - Automated Aurora - Push-button scaling - Automated Balchover - Backup and Recovery - Isolation and security - Industry compliance - Push-button scaling - Automated Balchover - Backup and Recovery - Industry compliance - Backtrack: restore data at any point of time without using backups - Define a subset of Aurora Instances as a Custom Endpoint		- Managed Oracle and Microsoft SQL Server Database with OS and database
Custom: access to the underlying database and OS so you can - Configure settings - Install patches - Enable native features - Access the underlying EC2 instance using SSH or SSM Session Manager - De-activate Automation Mode to perform your customization, better to take a DB snapshot before RDS vs RDS-custom: - RDS, entire database and the OS to be managed by AWS - RDS custom full admin access to the underlying OS and the database Autors is a proprietary technology from AWS (not open sourced) - Postgress and MySCJ are both supported as Auvora DB that means your drivers will work as if Autors ava as Postgres or MySCJ that means your drivers will work as if Autors ava as Postgres or MySCJ and the database Amazon Aurora Amazon Aurora Amazon Aurora Arac an Save up to 15 replicas and the replication process is faster than MySCJ, or MSC by or Sx the performance of Postgres on RDS - Aurora can have up to 15 replicas and the replication process is faster than MySCJ, or MSCJ, or MSC by a stop and the replication process is faster than MySCJ, or MSCJ, or		customization
- Configure settings - Instalt patches - Enable native features - Access the underlying ECZ instance using SSH or SSM Session Manager - De-activate Automation Mode to perform your customization, better to take a DB snapshot before RDS vs RDS-custom: - R		- RDS: Automates setup, operation, and scaling of database in AWS
- Configure settings - Instalt patches - Instalte patches - Enable native features - Access the underlying EC2 instance using SSH or SSM Session Manager - De-activate Automation Mode to perform your customization, better to take a DB snapshot before RDS vs RDS-custom: - RDS entire database and the OS to be managed by AWS - RDS Customs full admin access to the underlying OS and the database Aurora is a proprietary technology from AWS (not open sourced) - Postgres and MySQL are both supported as Aurora DB (that means your drivers will work as if Aurora was a Postgres or MySQL database) - Aurora is if Aurora was a Postgres or MySQL database) - Aurora is if Aurora was a Postgres or MySQL database - Aurora a Stonage automatically grows in increments of 10GB, up to 128 TB Aurora atorage automatically grows in increments of 10GB, up to 128 TB Aurora ana have up to 15 replicas and the replication process is faster than MySQL (sub 10 ms repl		
## Custom Install platches		, ,
### Paper Paper		
- Access the underlying EC2 instance using SSH or SSM Session Manager - De-activate Automation Mode to perform your customization, better to take a DB snapshot before RDS vs RDS-custom: - RDS-entire database and the OS to be managed by AWS - RDS Custom: full admin access to the underlying OS and the database Aurora is a proprietary technology from AWS (not open sourced) - Postgres and MySQL are both supported as Aurora DB (that means your drivers will work as if Aurora was a Postgres or MySQL database) - Aurora is 1 Aurora was a Postgres or MySQL database) - Aurora is 1 Aurora was a Postgres or MySQL database) - Aurora to 1 Aurora was a Postgres or MySQL database) - Aurora to 1 Aurora was a Postgres or MySQL database) - Aurora to 1 Aurora was a Postgres or MySQL database) - Aurora to 1 Aurora was a Postgres or MySQL database) - Aurora to 1 Aurora was a Postgres or MySQL database) - Aurora to 1 Aurora was a Postgres or MySQL database) - Aurora to 1 Aurora was a Postgres or MySQL database) - Aurora to 1 Aurora was a Postgres or MySQL database) - Aurora to 1 Aurora was a Postgres or MySQL database) - Aurora database or MySQL are both supported as Aurora Database - Aurora database or MySQL are both supported as Aurora Database - Aurora database or MySQL are both supported as Aurora Hostgres or MySQL database) - Aurora database or MySQL are both supported as Aurora Hostgres or MySQL database) - Aurora database and the replication - Storage is striped across 100s of volumes - Aurora High Availability and Read Scaling - Automated failover for master in less than - 30 seconds - Master + up to 15 Aurora Read Replicas - Serve reads - Support for Cross Region Replication - Automated failover - Backup and Recovery - Isolation and security - Indiustry compliance - Push-button scaling - Automated Patching with Zero Downtime - Advanced Monitoring - Routine Maintenance - Backtrack: restore data at any point of time without using backups - Define a subset of Aurora Instances as a Custom Endpoint - Example: Run analytical que		
- De-activate Automation Mode to perform your customization, better to take a DB snapshot before RDS vs RDS-custom: - RDS: entire database and the OS to be managed by AWS - RDS Custom: full admin access to the underlying OS and the database Aurora is a proprietary technology from AWS (not open sourced) - Postgres and MySQL are both supported as Aurora DB (that means your drivers will work as if Aurora was a Postgress or MySQL database) - Aurora is 'AWS cloud optimized' and claims 5x performance improvement over MySQL on RDS, over 3x the performance of Postgres on RDS - Aurora storage automatically grows in increments of 10GB, up to 128 TB Aurora can have up to 15 replicas and the replication process is faster than MySQL (sub to 10x replica tag) - Fallover in Aurora is instantaneous. It's HA (High Availability) native. 6 copies of your data across 3 AZ: - 4 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 4 copies out of 6 needed for writes - 5 set healing with peer-to-peer replication - 5 torage is striped across 100s of volumes - 0 ne Aurora Instance takes writes (master) - 4 utomated fail clover for master in less than - 30 seconds - Master + up to 15 Aurora Read Replicas - serve reads - 5 support for Cross Region Replication - 4 utomated fail-over - 8 ackup and Recovery - 8 location and security - 1 industry compliance - 9 ush button scaling - Automated Patching with Zero Downtime - 4 dvanced Monitoring - Routine Walneranace - 8 acktrack restore data at any point of time without using backups - Define a subset of Aurora Instances as a Custom Endpoint - 5 cample: Run analytical queries on specific replicas	RDS Custom	
snapshot before RDS vs RDS-custom: - RDS: entire database and the OS to be managed by AWS - RDS Custom: full admin access to the underlying OS and the database Aurora is a proprietary technology from AWS (not open sourced) - Postgres and MySQL are both supported as Aurora DB (that means your drivers will work as if Aurora was a Postgres or MySQL database) - Aurora is "aWS cloud optimized" and claims 5x performance improvement over MySQL on RDS, over 5x the performance of Postgres on RDS - Aurora storage automatically grows in increments of 106B, up to 128 TB Aurora can have up to 15 replicas and the replication process is faster than MySQL (sub 10 ms replica lag) - Fallover in Aurora is instantaneous. It's HA (High Availability) native. d copies of your data across 3 AZ! - 4 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies of volumes - One Aurora instance takes writes (master) - Automated failover for master in less than 30 seconds - Master + up to 15 Aurora Read Replicas - Support for Cross Region Replication - Automated failover for master in less than 30 seconds - Master + up to 15 Aurora Read Replicas - Support for Cross Region Replication - Automated failover for Aurora Read Replicas - Support for Cross Region Replication - Automated Patching with Zero Downtime - Advanced Monitoring - Routine Maintenance - Backtrack restore data at any point of time without using backups - Define a subset of Aurora Instances as a Custom Endpoint - Example. Run analytical queries on specific replicas		- Access the underlying EC2 instance using SSH or SSM Session Manager
RDS vs RDS-custom: - RDS: entire database and the OS to be managed by AWS - RDS Custom: full admin access to the underlying OS and the database Aurora is a proprietary technology from AWS (not open sourced) - Postgres and MySQL are both supported as Aurora DB (that means your drivers will work as if Aurora was a Postgres or MySQL database) - Aurora is "AWS cloud optimized" and claims 5x performance improvement over MySQL on RDS, over 3x the performance of Postgres on RDS - Aurora can have up to 15 replicas and the replication process is faster than MySQL (sub 10 ms replica lag) - Failover in Aurora is instantaneous. It's HA (High Availability) native. 6 copies of your data across 3 x2: - 4 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for master in less than 30 seconds - Master - up to 15 Aurora Read Replicas serve reads - Automate failover for master in less than 30 seconds - Master - up to 15 Aurora Read Replicas server leads - Support for Cross Region Replication - Automate fail-over - Backup and Recovery - Industry compliance - Push-button scaling - Automated Patching with Zero Downtime - Advanced Monitoring - Routine Maintenance - Backtrack: restore data at any point of time without using backups - Define a subset of Aurora Instances as a Custom Endpoint - Example: Run analytical queries on specific replicas		- De-activate Automation Mode to perform your customization, better to take a DB
- RDS: entire database and the OS to be managed by AWS - RDS Custom: full admin access to the underlying OS and the database Aurora is a proprietary technology from AWS (not open sourced) - Postgres and MySQL are both supported as Aurora D8 (that means your drivers will work as if Aurora was a Postgres or MySQL database) - Aurora is 'AWS cloud optimized' and claims 'Sx performance improvement over MySQL on RDS, over 3x the performance of Postgres on RDS - Aurora storage automatically grows in increments of 1006, up to 126 TB, - Aurora can have up to 15 replicas and the replication process is faster than MySQL (sub 10 ms replica tag) - Failover in Aurora is instantaneous. It's HA (High Availability) native. 6 copies of your data across 3 AZ: - 4 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 4 valuability and Read Scaling - 4 valuability and Read Scaling - 5 valuability and Read Scaling - 6 valuability and Read Scaling - 6 valuability and Read Scaling - 6 valuability and Read Scaling - 7 valuability and Read Scaling - 8 valuability and Read Scaling - 9 valuability and Read Scaling - 9 valuability and Read Scaling - 1 valuability and Read Scaling - 2 valuability and Read Scaling - 3 valuability and Read Scaling - 4 valuability and Read Scaling - 5 valuability and Read Scaling - 6 valuability and Read Scaling - 6 valuability and Read Scaling - 7 valuability and Read Scaling - 8 valuability and Read Scaling - 9 val		snapshot before
- RDS Custom: full admin access to the underlying OS and the database Aurora is a proprietary technology from AWS (not open sourced) - Postgres and MySQL are both supported as Aurora DB (that means your drivers will work as if Aurora was a Postgres or MySQL database) - Aurora is 'AWS cloud optimized' and claims 5x performance improvement over MySQL on RDS, over 3x the performance of Postgres on RDS - Aurora storage automatically grows in increments of 10GB, up to 128 TB Aurora can have up to 15 replicas and the replication process is faster than MySQL (sub 10 ms replica lag) - Failover in Aurora is instantaneous. It's HA (High Availability) native. 6 copies of your data across 3 AZ: - 4 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for reads - Self healing with peer-to-peer replication - Storage is striped across 100s of volumes - One Aurora Instance takes writes (master) - Automated failover for master in less than - 30 seconds - Master - up to 15 Aurora Read Replicas - Serve reads - Support for Cross Region Replication - Automatic fail-over - Backup and Recovery - Isolation and security - Industry compliance - Push-button scaling - Automated Patching with Zero Downtime - Advanced Monitoring - Routine Maintenance - Backura and Recovery - Backup and Backup and Recovery - Backup and Recovery - Backup and Re		RDS vs RDS-custom:
- RDS Custom: full admin access to the underlying OS and the database Aurora is a proprietary technology from AWS (not open sourced) - Postgres and MySQL are both supported as Aurora DB (that means your drivers will work as if Aurora was a Postgres or MySQL database) - Aurora is "AWS cloud optimized" and claims 5x performance improvement over MySQL on RDS, over 5x the performance of Postgres on RDS - Aurora storage automatically grows in increments of 10GB, up to 128 TB, - Aurora can have up to 15 replicas and the replication process is faster than MySQL (sub 10 ms replica lag) - Failover in Aurora is instantaneous. It's HA (High Availability) native. 6 copies of your data across 3 AZ: - 4 copies out of 6 need for reads - Self healing with peer-to-peer replication - Storage is striped across 100s of volumes - One Aurora Instance takes writes (master) - Automated failover for master in less than - 30 seconds - Master - up to 15 Aurora Read Replicas serve reads - Support for Cross Region Replication - Automatic fail-over - Backup and Recovery - Isolation and security - Industry compliance - Push-button scaling - Automated Patching with Zero Downtime - Advanced Monitoring - Routine Maintenance - Backtrack: restore data at any point of time without using backups - Define a subset of Aurora Instances as a Custom Endpoint - Example: Run analytical queries on specific replicas		
Aurora is a proprietary technology from AWS (not open sourced) Postgres and MySQL are both supported as Aurora DB (that means your drivers will work as if Aurora was a Postgres or MySQL database) Aurora is "AWS cloud optimized" and claims 5x performance improvement over MySQL on RDS, over 3x the performance of Postgres on RDS Aurora storage automatically grows in increments of IGGB, up to 128 TB. Aurora can have up to 15 replicas and the replication process is faster than MySQL (sub ID ms replica lag) Failover in Aurora is instantaneous. It's HA (High Availability) native. 6 copies of your data across 3 AZ: 4 copies out of 6 needed for writes 3 copies out of 6 needed for writes 3 copies out of 6 needed for writes 5 seth healing with peer-to-peer replication 5 storage is striped across 100s of volumes Aurora High Availability and Read Scaling Aurora instance takes writes (master) Automated failover for master in less than 30 seconds Master • up to 15 Aurora Read Replicas serve reads • Support for Cross Region Replication Automatic fail-over Backup and Recovery Isolation and security Industry compliance Push-button scaling Automated Patching with Zero Downtime Advanced Monitoring Routine Maintenance Backtrack: restore data at any point of time without using backups Define a subset of Aurora Instances as a Custom Endpoint Example: Run analytical queries on specific replicas		- 1
Postgres and MySQL are both supported as Aurora DB (that means your drivers will work as if Aurora was a Postgres or MySQL database) Aurora is "AWS cloud optimized" and claims 5x performance improvement over MySQL on RDS, over 5x the performance of Postgres on RDS Aurora can have up to 15 replicas and the replication process is faster than MySQL (sub 10 ms replica lag) Failover in Aurora is instantaneous. It's HA (High Availability) native. 6 copies of your data across 3 AZ: 4 copies out of 6 needed for writes 5 copies out of 6 needed for writes 5 copies out of 6 needed for reads 5 set freating with peer-to-peer replication Storage is striped across 100s of volumes One Aurora Instance takes writes (master) Automated failover for master in less than 30 seconds Master • up to 15 Aurora Read Replicas serve reads Support for Cross Region Replication Automatic fail-over Backup and Recovery Isolation and security Industry compliance Peat-tuon scaling Automated Patching with Zero Downtime Advanced Monitoring Routine Maintenance Backtrack: restore data at any point of time without using backups Define a subset of Aurora Instances as a Custom Endpoint Example: Run analytical queries on specific replicas		NEW COSTONIA TO CARGONIA ACCESS TO THE OTHER TYPING GO AND THE GALABAGE
drivers will work as if Aurora was a Postgres or MySQL database) - Aurora is 'AWS cloud optimized' and claims 5x performance improvement over MySQL on RDS, over 3x the performance of Postgres on RDS - Aurora storage automatically grows in increments of 100.8, up to 128 TB. - Aurora can have up to 15 replicas and the replication process is faster than MySQL (sub 10 ms replica lag) - Fallover in Aurora is instantaneous. It's HA (High Availability) native. 6 copies of your data across 3 AZ: - 4 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 4 copies out of 6 needed for writes - 5 copies out of 6 needed for writes - 5 copies out of 6 needed for writes - 5 copies out of 6 needed for writes - 5 copies out of 6 needed for writes - 5 copies out of 6 needed for writes - 5 copies out of 6 needed for writes - 6 copies out of 6 needed		
Aurora is "AWS cloud optimized" and claims 5x performance improvement over MySQL on RDS, over 3x the performance of Postgres on RDS - Aurora storage automatically grows in increments of 10GB, up to 128 TB Aurora can have up to 15 replicas and the replication process is faster than MySQL (sub 10 ms replica lag) - Faillover in Aurora is instantaneous. It's HA (High Availability) native. d copies of your data across 3 AZ: - 4 copies out of 6 need for reads - 3 copies out of 6 need for reads - Self healing with peer-to-peer replication - Storage is striped across 100s of volumes - One Aurora Instance takes writes (master) - Automated failover for master in less than 30 seconds - Master - up to 15 Aurora Read Replicas serve reads - Support for Cross Region Replication - Automatic fail-over - Backup and Recovery - Isolation and security - Industry compliance - Push-button scaling - Automated Patching with Zero Downtime - Advanced Monitoring - Routine Maintenance - Backtrack: restore data at any point of time without using backups - Define a subset of Aurora Instances as a Custom Endpoint - Example: Run analytical queries on specific replicas		· Postgres and MySQL are both supported as Aurora DB (that means your
Amazon Aurora over MySQL on RDS, over 3x the performance of Postgres on RDS · Aurora storage automatically grows in increments of 10GB, up to 128 TB. · Aurora can have up to 15 replicas and the replication process is faster than MySQL (sub 10 ms replica lag) · Fallover in Aurora is instantaneous. It's HA (High Availability) native. decopies of your data across 3 AZ: · 4 copies out of 6 needed for writes · 3 copies out of 6 needed for reads · Self healing with peer-to-peer replication · Storage is striped across 100s of volumes · One Aurora Instance takes writes (master) · Automated failover for master in less than 30 seconds · Master + up to 15 Aurora Read Replicas serve reads · Support for Cross Region Replication - Automatic fail-over · Backup and Recovery · Isolation and security · Industry compliance · Push-button scaling · Automated Patching with Zero Downtime · Advanced Monitoring · Routine Maintenance · Backtrack: restore data at any point of time without using backups - Define a subset of Aurora Instances as a Custom Endpoint - Example: Run analytical queries on specific replicas		drivers will work as if Aurora was a Postgres or MySQL database)
- Aurora storage automatically grows in increments of 10GB, up to 128 TB Aurora can have up to 15 replicas and the replication process is faster than MySQL (sub 10 ms replica lag) - Failover in Aurora is instantaneous. It's HA (High Availability) native. 6 copies of your data across 3 AZ: - 4 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for reads - Self healing with peer-to-peer replication - Storage is striped across 100s of volumes - One Aurora Instance takes writes (master) - Automated failover for master in less than - 30 seconds - Master + up to 15 Aurora Read Replicas - Support for Cross Region Replication - Automatic fail-over - Backup and Recovery - Isolation and security - Industry compliance - Push-button scaling - Automated Patching with Zero Downtime - Advanced Monitoring - Routine Maintenance - Backtrack: restore data at any point of time without using backups - Define a subset of Aurora Instances as a Custom Endpoint - Example: Run analytical queries on specific replicas		 Aurora is "AWS cloud optimized" and claims 5x performance improvement
- Aurora can have up to 15 replicas and the replication process is faster than MySQL (sub 10 ms replica lag) - Failover in Aurora is instantaneous. It's HA (High Availability) native. 6 copies of your data across 3 AZ: - 4 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 needed for reads - Self healing with peer-to-peer replication - Storage is striped across 100s of volumes - One Aurora Instance takes writes (master) - Automated failover for master in less than - 30 seconds - Master + up to 15 Aurora Read Replicas - serve reads - Support for Cross Region Replication - Automatic fail-over - Backup and Recovery - Isolation and security - Industry compliance - Push-button scaling - Automated Patching with Zero Downtime - Advanced Monitoring - Routine Maintenance - Backtrack: restore data at any point of time without using backups - Define a subset of Aurora Instances as a Custom Endpoint - Example: Run analytical queries on specific repticas	Amazon Aurora	over MySQL on RDS, over 3x the performance of Postgres on RDS
MySQL (sub 10 ms replica lag) • Failover in Aurora is instantaneous. It's HA (High Availability) native. 6 copies of your data across 3 AZ: • 4 copies out of 6 needed for writes • 3 copies out of 6 need for reads • Self healing with peer-to-peer replication • Storage is striped across 100s of volumes Aurora High Availability and Read Scaling • One Aurora Instance takes writes (master) • Automated failover for master in less than 30 seconds • Master • up to 15 Aurora Read Replicas serve reads • Support for Cross Region Replication • Automatic fail-over • Backup and Recovery • Isolation and security • Industry compliance • Push-button scaling • Automated Patching with Zero Downtime • Advanced Monitoring • Routine Maintenance • Backtrack: restore data at any point of time without using backups • Define a subset of Aurora Instances as a Custom Endpoint • Example: Run analytical queries on specific replicas		· Aurora storage automatically grows in increments of 10GB, up to 128 TB.
Fallover in Aurora is instantaneous. It's HA (High Availability) native. 6 copies of your data across 3 A Z: 4 copies out of 6 needed for writes 3 copies out of 6 needed for writes 5 copies out of 6 needed for reads 5 elf healing with peer-to-peer replication 5 torage is striped across 100s of volumes One Aurora Instance takes writes (master) Automated failover for master in less than 30 seconds Master + up to 15 Aurora Read Replicas server reads Support for Cross Region Replication Automatic fail-over Backup and Recovery Isolation and security Industry compliance Push-button scaling Automated Patching with Zero Downtime Advanced Monitoring Routine Maintenance Backtrack: restore data at any point of time without using backups Define a subset of Aurora Instances as a Custom Endpoint Example: Run analytical queries on specific replicas		· Aurora can have up to 15 replicas and the replication process is faster than
6 copies of your data across 3 AZ: 4 copies out of 6 needed for writes 3 copies out of 6 needed for writes 3 copies out of 6 needed for reads Self healing with peer-to-peer replication Storage is striped across 100s of volumes One Aurora Instance takes writes (master) Automated failover for master in less than 30 seconds Master * up to 15 Aurora Read Replicas serve reads Support for Cross Region Replication Automatic fail-over Backup and Recovery Isolation and security Industry compliance Push-button scaling Automated Patching with Zero Downtime Advanced Monitoring Routine Maintenance Backtrack: restore data at any point of time without using backups Define a subset of Aurora Instances as a Custom Endpoint Example: Run analytical queries on specific replicas		MySQL (sub 10 ms replica lag)
- 4 copies out of 6 needed for writes - 3 copies out of 6 needed for writes - 3 copies out of 6 need for reads - Self healing with peer-to-peer replication - Storage is striped across 100s of volumes - One Aurora Instance takes writes (master) - Automated failover for master in less than - 30 seconds - Master + up to 15 Aurora Read Replicas - Support for Cross Region Replication - Automatic fail-over - Backup and Recovery - Isolation and security - Industry compliance - Push-button scaling - Automated Patching with Zero Downtime - Advanced Monitoring - Routine Maintenance - Backtrack: restore data at any point of time without using backups - Define a subset of Aurora Instances as a Custom Endpoint - Example: Run analytical queries on specific replicas		· Failover in Aurora is instantaneous. It's HA (High Availability) native.
- 3 copies out of 6 need for reads - Self healing with peer-to-peer replication - Storage is striped across 100s of volumes - One Aurora Instance takes writes (master) - Automated failover for master in less than - 30 seconds - Master + up to 15 Aurora Read Replicas - Support for Cross Region Replication - Automatic fail-over - Backup and Recovery - Isolation and security - Industry compliance - Push-button scaling - Automated Patching with Zero Downtime - Advanced Monitoring - Routine Maintenance - Backtrack: restore data at any point of time without using backups - Define a subset of Aurora Instances as a Custom Endpoint - Example: Run analytical queries on specific replicas		6 copies of your data across 3 AZ:
- Self healing with peer-to-peer replication - Storage is striped across 100s of volumes - One Aurora Instance takes writes (master) - Automated failover for master in less than - 30 seconds - Master + up to 15 Aurora Read Replicas - Support for Cross Region Replication - Automatic fail-over - Backup and Recovery - Isolation and security - Industry compliance - Push-button scaling - Automated Patching with Zero Downtime - Advanced Monitoring - Routine Maintenance - Backtrack: restore data at any point of time without using backups - Define a subset of Aurora Instances as a Custom Endpoint - Example: Run analytical queries on specific replicas		· 4 copies out of 6 needed for writes
Aurora High Availability and Read Scaling One Aurora Instance takes writes (master) Automated failover for master in less than 30 seconds Master + up to 15 Aurora Read Replicas serve reads Support for Cross Region Replication Automatic fail-over Backup and Recovery Isolation and security Industry compliance Push-button scaling Automated Patching with Zero Downtime Advanced Monitoring Routine Maintenance Backtrack: restore data at any point of time without using backups Define a subset of Aurora Instances as a Custom Endpoint Example: Run analytical queries on specific replicas		· 3 copies out of 6 need for reads
Aurora High Availability and Read Scaling One Aurora Instance takes writes (master) Automated failover for master in less than 30 seconds Master + up to 15 Aurora Read Replicas serve reads Support for Cross Region Replication Automatic fail-over Backup and Recovery Isolation and security Industry compliance Push-button scaling Automated Patching with Zero Downtime Advanced Monitoring Routine Maintenance Backtrack: restore data at any point of time without using backups Define a subset of Aurora Instances as a Custom Endpoint Example: Run analytical queries on specific replicas		· Self healing with peer-to-peer replication
Automated failover for master in less than 30 seconds • Master + up to 15 Aurora Read Replicas serve reads • Support for Cross Region Replication • Automatic fail-over • Backup and Recovery • Isolation and security • Industry compliance • Push-button scaling • Automated Patching with Zero Downtime • Advanced Monitoring • Routine Maintenance • Backtrack: restore data at any point of time without using backups • Define a subset of Aurora Instances as a Custom Endpoint • Example: Run analytical queries on specific replicas		· Storage is striped across 100s of volumes
30 seconds • Master + up to 15 Aurora Read Replicas serve reads • Support for Cross Region Replication • Automatic fail-over • Backup and Recovery • Isolation and security • Industry compliance • Push-button scaling • Automated Patching with Zero Downtime • Advanced Monitoring • Routine Maintenance • Backtrack: restore data at any point of time without using backups • Define a subset of Aurora Instances as a Custom Endpoint Aurora - Custom Endpoints • Example: Run analytical queries on specific replicas	Aurora High Availability and Read Scaling	· One Aurora Instance takes writes (master)
Master + up to 15 Aurora Read Replicas serve reads Support for Cross Region Replication Automatic fail-over Backup and Recovery Isolation and security Industry compliance Push-button scaling Automated Patching with Zero Downtime Advanced Monitoring Routine Maintenance Backtrack: restore data at any point of time without using backups Define a subset of Aurora Instances as a Custom Endpoint Example: Run analytical queries on specific replicas		· Automated failover for master in less than
serve reads Support for Cross Region Replication Automatic fail-over Backup and Recovery Isolation and security Industry compliance Push-button scaling Automated Patching with Zero Downtime Advanced Monitoring Routine Maintenance Backtrack: restore data at any point of time without using backups Define a subset of Aurora Instances as a Custom Endpoint Example: Run analytical queries on specific replicas		30 seconds
- Support for Cross Region Replication - Automatic fail-over - Backup and Recovery - Isolation and security - Industry compliance - Push-button scaling - Automated Patching with Zero Downtime - Advanced Monitoring - Routine Maintenance - Backtrack: restore data at any point of time without using backups - Define a subset of Aurora Instances as a Custom Endpoint - Example: Run analytical queries on specific replicas		• Master + up to 15 Aurora Read Replicas
- Automatic fail-over - Backup and Recovery - Isolation and security - Industry compliance - Push-button scaling - Automated Patching with Zero Downtime - Advanced Monitoring - Routine Maintenance - Backtrack: restore data at any point of time without using backups - Define a subset of Aurora Instances as a Custom Endpoint - Example: Run analytical queries on specific replicas		serve reads
Backup and Recovery Isolation and security Industry compliance Push-button scaling Automated Patching with Zero Downtime Advanced Monitoring Routine Maintenance Backtrack: restore data at any point of time without using backups Define a subset of Aurora Instances as a Custom Endpoint Example: Run analytical queries on specific replicas		Support for Cross Region Replication
Isolation and security Industry compliance Push-button scaling Automated Patching with Zero Downtime Advanced Monitoring Routine Maintenance Backtrack: restore data at any point of time without using backups Define a subset of Aurora Instances as a Custom Endpoint Example: Run analytical queries on specific replicas		
Industry compliance Push-button scaling Automated Patching with Zero Downtime Advanced Monitoring Routine Maintenance Backtrack: restore data at any point of time without using backups Define a subset of Aurora Instances as a Custom Endpoint Example: Run analytical queries on specific replicas		
Features of Aurora Push-button scaling Automated Patching with Zero Downtime Advanced Monitoring Routine Maintenance Backtrack: restore data at any point of time without using backups Define a subset of Aurora Instances as a Custom Endpoint Example: Run analytical queries on specific replicas		· Isolation and security
Automated Patching with Zero Downtime Advanced Monitoring Routine Maintenance Backtrack: restore data at any point of time without using backups Define a subset of Aurora Instances as a Custom Endpoint Example: Run analytical queries on specific replicas		· Industry compliance
Advanced Monitoring Routine Maintenance Backtrack: restore data at any point of time without using backups Define a subset of Aurora Instances as a Custom Endpoint Example: Run analytical queries on specific replicas	Features of Aurora	Push-button scaling
Routine Maintenance Backtrack: restore data at any point of time without using backups Define a subset of Aurora Instances as a Custom Endpoint Example: Run analytical queries on specific replicas		Automated Patching with Zero Downtime
Backtrack: restore data at any point of time without using backups Define a subset of Aurora Instances as a Custom Endpoint Example: Run analytical queries on specific replicas		· Advanced Monitoring
Define a subset of Aurora Instances as a Custom Endpoint Aurora - Custom Endpoints Example: Run analytical queries on specific replicas		· Routine Maintenance
Aurora - Custom Endpoints • Example: Run analytical queries on specific replicas		Backtrack: restore data at any point of time without using backups
		Define a subset of Aurora Instances as a Custom Endpoint
• The Reader Endpoint is generally not used after defining Custom Endpoints	Aurora - Custom Endpoints	· Example: Run analytical queries on specific replicas
		• The Reader Endpoint is generally not used after defining Custom Endpoints

Aurora Serverless Aurora Serverless Aurora Serverless Aurora Multi-Master Aurora Multi-Master Aurora Coross Region Read Replicas: - Useful for disaster recovery - Simple to put in place - Aurora Global Aurora Global Aurora Autora Autora
usage Good for infrequent, intermittent or unpredictable workloads No capacity planning needed Pay per second, can be more cost-effective In case you want continuous write availability for the writer nodes Every node does R/W - vs promoting a Read Replica as the new master Aurora Multi-Master Aurora Cross Region Read Replicas: Useful for disaster recovery Simple to put in place Aurora Global Database (recommended): 1 Primary Region (read / write) Up to 5 secondary (read-only) regions, replication lag is less than 1 second Up to 16 Read Replicas per secondary region
Aurora Serverless Good for infrequent, intermittent or unpredictable workloads No capacity planning needed Pay per second, can be more cost-effective In case you want continuous write availability for the writer nodes Every node does R/W - vs promoting a Read Replica as the new master Aurora Cross Region Read Replicas: Useful for disaster recovery Simple to put in place Aurora Global Database (recommended): I Primary Region (read / write) Up to 5 secondary (read-only) regions, replication lag is less than 1 second Up to 16 Read Replicas per secondary region
Aurora Serverless intermittent or unpredictable workloads No capacity planning needed Pay per second, can be more cost-effective In case you want continuous write availability for the writer nodes Every node does R/W - vs promoting a Read Replica as the new master Aurora Cross Region Read Replicas: Useful for disaster recovery Simple to put in place Aurora Global Database (recommended): 1 Primary Region (read / write) Up to 5 secondary (read-only) regions, replication lag is less than 1 second Up to 16 Read Replicas per secondary region
unpredictable workloads No capacity planning needed Pay per second, can be more cost-effective In case you want continuous write availability for the writer nodes Every node does R/W - vs promoting a Read Replica as the new master Aurora Cross Region Read Replicas: Useful for disaster recovery Simple to put in place Aurora Global Database (recommended): 1 Primary Region (read / write) Up to 5 secondary (read-only) regions, replication lag is less than 1 second Up to 16 Read Replicas per secondary region
No capacity planning needed Pay per second, can be more cost-effective Aurora Multi-Master In case you want continuous write availability for the writer nodes Every node does R/W - vs promoting a Read Replica as the new master Aurora Cross Region Read Replicas: Useful for disaster recovery Simple to put in place Aurora Global Database (recommended): 1 Primary Region (read / write) Up to 5 secondary (read-only) regions, replication lag is less than 1 second Up to 16 Read Replicas per secondary region
needed Pay per second, can be more cost-effective Aurora Multi-Master In case you want continuous write availability for the writer nodes Every node does R/W - vs promoting a Read Replica as the new master Aurora Cross Region Read Replicas: Useful for disaster recovery Simple to put in place Aurora Global Database (recommended): I Primary Region (read / write) Up to 5 secondary (read-only) regions, replication lag is less than 1 second Up to 16 Read Replicas per secondary region
Pay per second, can be more cost-effective In case you want continuous write availability for the writer nodes • Every node does R/W - vs promoting a Read Replica as the new master Aurora Cross Region Read Replicas: • Useful for disaster recovery • Simple to put in place • Aurora Global Database (recommended): • 1 Primary Region (read / write) • Up to 5 secondary (read-only) regions, replication lag is less than 1 second • Up to 16 Read Replicas per secondary region
Aurora Multi-Master In case you want continuous write availability for the writer nodes Every node does R/W - vs promoting a Read Replica as the new master Aurora Cross Region Read Replicas: Useful for disaster recovery Simple to put in place Aurora Global Database (recommended): 1 Primary Region (read / write) Up to 5 secondary (read-only) regions, replication lag is less than 1 second Up to 16 Read Replicas per secondary region
Aurora Multi-Master In case you want continuous write availability for the writer nodes Every node does R/W - vs promoting a Read Replica as the new master Aurora Cross Region Read Replicas: Useful for disaster recovery Simple to put in place Aurora Global Database (recommended): 1 Primary Region (read / write) Up to 5 secondary (read-only) regions, replication lag is less than 1 second Up to 16 Read Replicas per secondary region
Aurora Multi-Master - Every node does R/W - vs promoting a Read Replica as the new master Aurora Cross Region Read Replicas: - Useful for disaster recovery - Simple to put in place - Aurora Global Database (recommended): - 1 Primary Region (read / write) - Up to 5 secondary (read-only) regions, replication lag is less than 1 second - Up to 16 Read Replicas per secondary region
- Every node does R/W - vs promoting a Read Replica as the new master Aurora Cross Region Read Replicas: - Useful for disaster recovery - Simple to put in place - Aurora Global Database (recommended): - 1 Primary Region (read / write) - Up to 5 secondary (read-only) regions, replication lag is less than 1 second - Up to 16 Read Replicas per secondary region
Useful for disaster recovery Simple to put in place Aurora Global Database (recommended): I Primary Region (read / write) Up to 5 secondary (read-only) regions, replication lag is less than 1 second Up to 16 Read Replicas per secondary region
- Simple to put in place - Aurora Global Database (recommended): - 1 Primary Region (read / write) - Up to 5 secondary (read-only) regions, replication lag is less than 1 second - Up to 16 Read Replicas per secondary region
- Aurora Global Database (recommended): - 1 Primary Region (read / write) - Up to 5 secondary (read-only) regions, replication lag is less than 1 second - Up to 16 Read Replicas per secondary region
Global Aurora - 1 Primary Region (read / write) - Up to 5 secondary (read-only) regions, replication lag is less than 1 second - Up to 16 Read Replicas per secondary region
Global Aurora • Up to 5 secondary (read-only) regions, replication lag is less than 1 second • Up to 16 Read Replicas per secondary region
than 1 second Up to 16 Read Replicas per secondary region
than 1 second • Up to 16 Read Replicas per secondary region
· Helps for decreasing latency
· Promoting another region (for disaster recovery) has an RTO
of < 1 minute
Typical cross-region replication takes less than 1 second
Enables you to add ML-based predictions to
your applications via SQL
Simple, optimized, and secure integration
between Aurora and AWS ML services
- Supported services Aurora Machine Learning
• Amazon SageMaker (use with any ML model)
Amazon Comprehend (for sentiment analysis)
· You don't need to have ML experience
 Use cases: fraud detection, ads targeting,
sentiment analysis, product recommendations
Automated backups:
Daily full backup of the database (during the backup window)
· Transaction logs are backed-up by RDS every 5 minutes
RDS Backups -=> ability to restore to any point in time (from oldest backup to 5 minutes ago)
· 1 to 35 days of retention, set 0 to disable automated backups
- Manual DB Snapshots
Manually triggered by the user
· Retention of backup for as long as you want
Automated backups
·1 to 35 days (cannot be disabled)
- point-in-time recovery in that timeframe
- point-in-time recovery in that timeframe - Manual DB Snapshots
Aurora Backups • point-in-time recovery in that timeframe

	Restoring a RDS / Aurora backup or a snapshot creates a new database
	· Restoring MySQL RDS database from S3
	· Create a backup of your on-premises database
	· Store it on Amazon S3 (object storage)
RDS & Aurora Restore options	Restore the backup file onto a new RDS instance running MySQL
	Restoring MySQL Aurora cluster from S3
	· Create a backup of your on-premises database using Percona XtraBackup
	Store the backup file on Amazon S3
	Restore the backup file onto a new Aurora cluster running MySQL
	Create a new Aurora DB Cluster from an
	existing one
	· Faster than snapshot & restore
	· Uses copy-on-write protocol
	· Initially, the new DB cluster uses the same data
Aurora Database Cloning	volume as the original DB cluster (fast and efficient
	- no copying is needed)
	When updates are made to the new DB cluster
	data, then additional storage is allocated and data is
	copied to be separated
	At-rest encryption:
	- Database master & replicas encryption using AWS KMS - must be defined as launch
	time
	- If the master is not encrypted, the read replicas cannot be encrypted
	- To encrypt an un-encrypted database, go through a DB snapshot & restore as
	encrypted
	In-flight encryption: TLS-ready by default, use the AWS TLS root certificates client-side.
RDS & Aurora Security	
	IAM Authentication: IAM roles to connect to your database (instead of username/pw).
	Security Groups: Control Network access to your RDS / Aurora DB.
	No SSH available except on RDS Custom.
	Audit Logs can be enabled and sent to CloudWatch Logs for longer retention.
	-Fully managed database proxy for RDS
	-Fully managed database proxy for RDS -Allows apps to pool and share DB connections established with the database
	-Improving database efficiency by reducing the stress on database resources (e.g.,
	CPU, RAM) and minimize open connections (and timeouts)
Amoron DDC Brown	-Serverless, autoscaling, highly available (multi-AZ)
Amazon RDS Proxy	-Reduced RDS & Aurora failover time by up 66%
	-Supports RDS (MySQL, PostgreSQL, MariaDB) and Aurora (MySQL, PostgreSQL)
	-No code changes required for most apps
	-Enforce IAM Authentication for DB, and securely store credentials in AWS Secrets
	Manager
	-RDS Proxy is never publicly accessible (must be accessed from VPC)

Amazon ElastiCache Overview	The same way RDS is to get managed Relational Databases ElastiCache is to get managed Redis or Memcached Caches are in-memory databases with really high performance, low latency Helps reduce load off of databases for read intensive workloads Helps make your application stateless AWS takes care of OS maintenance / patching, optimizations, setup, configuration, monitoring, failure recovery and backups Using ElastiCache involves heavy application code changes Applications queries ElastiCache, if not available, get from RDS and store in
ElastiCache Solution Architecture - DB Cache	ElastiCache. - Helps relieve load in RDS - Cache must have an invalidation strategy to make sure only the most current data is used in there.
ElastiCache Solution Architecture - User Session Store	User logs into any of the application The application writes the session data into ElastiCache The user hits another instance of our application The instance retrieves the data and the user is already logged in
ElastiCache - Redis	REDIS Multi AZ with Auto-Failover Read Replicas to scale reads and have high availability Data Durability using AOF persistence Backup and restore features Supports Sets and Sorted Sets
ElastiCache - Memcached	 Multi-node for partitioning of data (sharding) No high availability (replication) Non persistent No backup and restore Multi-threaded architecture
ElastiCache - Cache Security	ElastiCache supports IAM Authentication for Redis IAM policies on ElastiCache are only used for AWS API-level security Redis AUTH You can set a "password/token" when you create a Redis cluster This is an extra level of security for your cache (on top of security groups) Support SSL in flight encryption

	Lazy Loading: all the read data is
	cached, data can become stale in
	cache
	Write Through: Adds or update
Patterns for ElastiCache	data in the cache when written
according for Etasticache	to a DB (no stale data)
	Session Store: store temporary
	session data in a cache (using
	TTL features)
	Domain Registrar: Amazon Route 53, GoDaddy,
	DNS Records: A, AAAA, CNAME, NS, Zone File: contains DNS records
DNS Terminologies	
	 Name Server: resolves DNS queries (Authoritative or Non-Authoritative) Top Level Domain (TLD): .com, .us, .in, .gov, .org,
	Second Level Domain (SLD): amazon.com, google.com,
	A highly available, scalable, fully
	managed and Authoritative DNS
	Authoritative = the customer (you)
	can update the DNS records
Amazon Route 53	• Route 53 is also a Domain Registrar
	Ability to check the health of your
	resources
	• The only AWS service which
	provides 100% availability SLA
	· A - maps a hostname to IPv4
	• AAAA - maps a hostname to IPv6
	CNAME - maps a hostname to another hostname
	The target is a domain name which must have an A or AAAA record
Route 53 - Record Types	Can't create a CNAME record for the top node of a DNS namespace (Zone)
,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	Apex)
	Example: you can't create for example.com, but you can create for
	www.example.com
	· NS - Name Servers for the Hosted Zone
	Control how traffic is routed for a domain
	A container for records that define how to route traffic to a domain and
	its subdomains
	 Public Hosted Zones - contains records that specify how to route
	traffic on the Internet (public domain names)
Route 53 - Hosted Zones	application1.mypublicdomain.com
	· Private Hosted Zones - contain records that specify how you route
	traffic within one or more VPCs (private domain names)
	application1.company.internal
	· You pay \$0.50 per month per hosted zone
	High TTL - e.g., 24 hr
	· Less traffic on Route 53
	Possibly outdated records
Double 57 December TTI (Time - To Live)	• Low TTL - e.g., 60 sec.
Route 53 - Records TTL (Time To Live)	- More traffic on Route 53 (\$\$)
	• Records are outdated for less
	· Records are obtuated for tess
	time

CNAME vs Alias	CNAME: Points a hostname to any other hostname. (app.mydomain.com => blabla.anything.com) ONLY FOR NON ROOT DOMAIN (aka. something.mydomain.com) Alias: Points a hostname to an AWS Resource (app.mydomain.com => blabla.amazonaws.com) Works for ROOT DOMAIN and NON ROOT DOMAIN (aka mydomain.com) Free of charge
	Native health check
Route 53 - Alias Records	Maps a hostname to an AWS resource An extension to DNS functionality Automatically recognizes changes in the resource's IP addresses Unlike CNAME, it can be used for the top node of a DNS namespace (Zone Apex), e.g.: example.com Alias Record is always of type A/AAAA for AWS resources (IPv4 / IPv6) You can't set the TTL
Route 53 - Alias Records Targets	Elastic Load Balancers CloudFront Distributions API Gateway Elastic Beanstalk environments S3 Websites VPC Interface Endpoints Global Accelerator accelerator Route 53 record in the same hosted zone
Routing Policies - Weighted	Control the % of the requests that go to each specific resource Assign each record a relative weight: traffic (%) = !"#\$%& '() * +,"-##-)"-(). /01 (' 22 &%" 3"#\$%&+ '() 22)"-().+ Weights don't need to sum up to 100 DNS records must have the same name and type Can be associated with Health Checks Use cases: load balancing between regions, testing new application versions
Routing Policies - Latency-based	Redirect to the resource that has the least latency close to us Super helpful when latency for users is a priority Latency is based on traffic between users and AWS Regions Germany users may be directed to the US (if that's the lowest latency) Can be associated with Health Checks (has a failover capability)

	HTTP Health Checks are only for public
	resources
	Health Check => Automated DNS Failover:
	Health checks that monitor an endpoint
	(application, server, other AWS resource)
Davida 57 Haaliba Chaalia	2. Health checks that monitor other health
Route 53 - Health Checks	checks (Calculated Health Checks)
	3. Health checks that monitor CloudWatch
	Alarms (full control !!) - e.g., throttles of
	DynamoDB, alarms on RDS, custom metrics,
	(helpful for private resources)
	Health Checks are integrated with CW
	metrics
	About 15 global health checkers will check the
	endpoint health
	· Healthy/Unhealthy Threshold - 3 (default)
	· Interval - 30 sec (can set to 10 sec - higher cost)
	Supported protocol: HTTP, HTTPS and TCP
Harith Charles Maritan as Furdinaint	• If > 18% of health checkers report the endpoint is
Health Checks - Monitor an Endpoint	healthy, Route 53 considers it Healthy. Otherwise, it's
	Unhealthy
	· Ability to choose which locations you want Route 53 to
	USE
	· Health Checks pass only when the endpoint
	responds with the 2xx and 3xx status codes
	Combine the results of multiple Health
	Checks into a single Health Check
	· You can use OR, AND, or NOT
Route 53 - Calculated Health Checks	• Can monitor up to 256 Child Health Checks
	Specify how many of the health checks need
	to pass to make the parent pass
	Different from Latency-based!
	This routing is based on user location
	Specify location by Continent, Country
	or by US State (if there's overlapping,
	most precise location selected)
Routing Policies - Geolocation	Should create a "Default" record (in
	case there's no match on location)
	Use cases: website localization, restrict
	content distribution, load balancing,
	Can be associated with Health Checks
	Route traffic to your resources based on the geographic location of users and
	resources Ability to chiff more traffic to resources based on the defined bias
	Ability to shift more traffic to resources based on the defined bias To change the size of the goographic ragion, specify bias values.
	• To change the size of the geographic region, specify bias values:
Routing Policies - Geoproximity	• To expand (1 to 99) - more traffic to the resource
	• To shrink (-1 to -99) - less traffic to the resource
	Resources can be:
	AWS resources (specify AWS region) Non AWS resources (specify Latitude and Longitude)
	 Non-AWS resources (specify Latitude and Longitude) You must use Route 53 Traffic Flow to use this feature
	• 100 most use noute 33 maint i tow to use this reature

	Routing is based on clients' IP addresses
	· You provide a list of CIDRs for your clients
	and the corresponding endpoints/locations
	(user-IP-to-endpoint mappings)
Routing Policies - IP-based Routing	· Use cases: Optimize performance, reduce
	network costs
	• Example: route end users from a particular
	ISP to a specific endpoint
	Use when routing traffic to multiple resources
	· Route 53 return multiple values/resources
Routing Policies - Multi-Value	· Can be associated with Health Checks (return only values for healthy resources)
	 Up to 8 healthy records are returned for each Multi-Value query
	Multi-Value is not a substitute for having an ELB
	You buy or register your domain name with a Domain Registrar typically by
	paying annual charges (e.g., GoDaddy, Amazon Registrar Inc.,)
	• The Domain Registrar usually provides you with a DNS service to manage
	your DNS records
Domain Registar vs. DNS Service	But you can use another DNS service to manage your DNS records
	• Example: purchase the domain from GoDaddy and use Route 53 to manage
	your DNS records
	Amazon
	If you buy your domain on a 3rd party registrar, you can still use
	Route 53 as the DNS Service provider
	Create a Hosted Zone in Route 53
3rd Party Registrar with Amazon Route 53	2. Update NS Records on 3rd party website to use Route 53 Name
Sid Faity Registral With Amazon Roote 33	Servers
	Domain Registrar != DNS Service
	But every Domain Registrar usually comes with some DNS features
	, , ,
	ELB sticky sessions
	• Web clients for storing cookies and making our web app stateless
	• ElastiCache
	• For storing sessions (alternative: DynamoDB)
3-tier architectures for web applications	• For caching data from RDS
	• Multi AZ
	·RDS
	· For storing user data
	· Read replicas for scaling reads
	· Multi AZ for disaster recovery
	EC2 Instances:
	· Use a Golden AMI: Install your applications, OS dependencies etc beforehand
	and launch your EC2 instance from the Golden AMI
	Bootstrap using User Data: For dynamic configuration, use User Data scripts
Instantiating Applications quickly	Hybrid: mix Golden AMI and User Data (Elastic Beanstalk)
	• RDS Databases:
	• Restore from a snapshot: the database will have schemas and data ready!
	· EBS Volumes:
	• Restore from a snapshot: the disk will already be formatted and have data!

Elastic Beanstalk - Overview	Elastic Beanstalk is a developer centric view of deploying an application on AWS It uses all the component's we've seen before: EC2, ASG, ELB, RDS, Managed service Automatically handles capacity provisioning, load balancing, scaling, application
	health monitoring, instance configuration, Just the application code is the responsibility of the developer We still have full control over the configuration Beanstalk is free but you pay for the underlying instances
	Application: collection of Elastic Beanstalk components (environments,
	versions, configurations,)
	Application Version: an iteration of your application code
Elastic Beanstalk - Components	• Environment
	Collection of AWS resources running an application version (only one application
	version at a time)
	• Tiers: Web Server Environment Tier & Worker Environment Tier
	You can create multiple environments (dev, test, prod,)
	· Go
	· Java SE
	· Java with Tomcat
	· .NET Core on Linux
	· .NET on Windows Server
	· Node.js
Elastic Beanstalk - Supported Platforms	·PHP
	• Python
	· Ruby
	Packer Builder
	· Single Container Docker
	Multi-container Docker
	Preconfigured Docker
	Amazon S3 allows people to store objects (files) in "buckets" (directories)
	Buckets must have a globally unique name (across all regions all accounts)
	Buckets are defined at the region level
	S3 looks like a global service but buckets are created in a region
	Naming convention
Amazon S3 - Buckets	No uppercase, No underscore
<u></u>	· 3-63 characters long
	• Not an IP
	Must start with lowercase letter or number
	Must NOT start with the prefix xn
	Floor To Four that the prent Air

	Objects (files) have a Key
	• The key is the FULL path:
	· s3://my-bucket/my_file.txt
	· s3://my-bucket/my_folder1/another_folder/my_file.txt
	• The key is composed of prefix + object name
	· s3://my-bucket/my_folder1/another_folder/my_file.txt
67 01: 1	· There's no concept of "directories" within buckets
Amazon S3 - Objects	(although the UI will trick you to think otherwise)
	Object values are the content of the body:
	• Max. Object Size is 5TB (5000GB)
	If uploading more than 5GB, must use "multi-part upload"
	Metadata (list of text key / value pairs - system or user metadata)
	• Tags (Unicode key / value pair - up to 10) - useful for security / lifecycle
	Version ID (if versioning is enabled)
	User-Based
	• IAM Policies - which API calls should be allowed for a specific user from IAM
	• Resource-Based
	Bucket Policies - bucket wide rules from the S3 console - allows cross account
	Object Access Control List (ACL) - finer grain (can be disabled)
Amazon S3 - Security	
	Bucket Access Control List (ACL) - less common (can be disabled) Note: an IANA principal can access an S.7 abject if
	Note: an IAM principal can access an S3 object if The area IAM principals and ALLOW it OR the presented policy ALLOW it.
	The user IAM permissions ALLOW it OR the resource policy ALLOWS it
	• AND there's no explicit DENY
	Encryption: encrypt objects in Amazon S3 using encryption keys
	JSON based policies
	Resources: buckets and objects
S3 Bucket Policies	• Effect: Allow / Deny
33 BUCKET FOLICIES	· Actions: Set of API to Allow or Deny
	· Principal: The account or user to apply the
	policy to
	If you get a 403 Forbidden error, make sure the bucket
Amazon S3 - Static Website Hosting	policy allows public reads!
	"• You can version your files in Amazon S3
	· It is enabled at the bucket level
	• Same key overwrite will change the "version": 1, 2, 3
	• It is best practice to version your buckets
	Protect against unintended deletes (ability to restore a version)
Amazon S3 - Versioning	Easy roll back to previous version
	Notes:
	Any file that is not versioned prior to enabling versioning will
	have version "null"
	Suspending versioning does not delete the previous versions"
	Must enable Versioning in source and destination buckets
	Cross-Region Replication (CRR)
	- Same-Region Replication (SRR)
	Buckets can be in different AWS accounts
	· Copying is asynchronous
Amazon S3 - Replication (CRR & SRR)	 Must give proper IAM permissions to \$3
	• Use cases:
	· CRR - compliance, lower latency access, replication across
	accounts
	· SRR - log aggregation, live replication between production and test
	accounts
	accounts

	After you enable Replication, only new objects are replicated Optionally, you can replicate existing objects using S3 Batch Replication
	 Replicates existing objects and objects that failed replication For DELETE operations
Amazon S3 - Replication (Notes)	Can replicate delete markers from source to target (optional setting)
	Deletions with a version ID are not replicated (to avoid malicious deletes)
	• There is no "chaining" of replication
	• If bucket 1 has replication into bucket 2, which has replication into bucket 3
	•Then objects created in bucket 1 are not replicated to bucket 3
	Durability:
	· High durability (99.99999999, 11 9's) of objects across multiple AZ
	· If you store 10,000,000 objects with Amazon S3, you can on average expect to
	incur a loss of a single object once every 10,000 years
S3 Durability and Availability	· Same for all storage classes
	· Availability:
	· Measures how readily available a service is
	· Varies depending on storage class
	• Example: S3 standard has 99.99% availability = not available 53 minutes a year
	99.99% Availability
	· Used for frequently accessed data
S3 Standard - General Purpose	· Low latency and high throughput
33 Standard General Forpose	Sustain 2 concurrent facility failures
	 Use Cases: Big Data analytics, mobile & gaming applications, content
	distribution
	For data that is less frequently accessed, but requires rapid access when needed
	· Lower cost than S3 Standard
	Amazon S3 Standard-Infrequent Access (S3 Standard-IA)
	• 99.9% Availability
S3 Storage Classes - Infrequent Access	Use cases: Disaster Recovery, backups
	· Amazon S3 One Zone-Infrequent Access (S3 One Zone-IA)
	· High durability (99.99999999) in a single AZ; data lost when AZ is destroyed
	• 99.5% Availability
	Use Cases: Storing secondary backup copies of on-premises data, or data you can
	recreate
	Low-cost object storage meant for archiving / backup
	Pricing: price for storage + object retrieval cost
	Amazon S3 Glacier Instant Retrieval Millionand retrieval great for data accessed and a greater.
	 Millisecond retrieval, great for data accessed once a quarter Minimum storage duration of 90 days
Amazon S3 Glacier Storage Classes	Amazon S3 Glacier Flexible Retrieval (formerly Amazon S3 Glacier):
Amazon 33 Glacier Storage Classes	• Expedited (1 to 5 minutes), Standard (3 to 5 hours), Bulk (5 to 12 hours) - free
	Minimum storage duration of 90 days
	Amazon S3 Glacier Deep Archive - for long term storage:
	Standard (12 hours), Bulk (48 hours)
	Minimum storage duration of 180 days
	Small monthly monitoring and auto-tiering fee
	Moves objects automatically between Access Tiers based on usage
	• There are no retrieval charges in S3 Intelligent-Tiering
	Frequent Access tier (automatic): default tier
S3 Intelligent-Tiering	Infrequent Access tier (automatic): objects not accessed for 30 days
	Archive Instant Access tier (automatic): objects not accessed for 90 days
	Archive Access tier (optional): configurable from 90 days to 700+ days
	Deep Archive Access tier (optional): config. from 180 days to 700+ days
	,

ou can transition objects between
orage classes
For infrequently accessed object,
ove them to Standard IA
For archive objects that you don't
eed fast access to, move them to
acier or Glacier Deep Archive
Moving objects can be automated
ing a Lifecycle Rules
ansition Actions - configure objects to transition to another storage class
Move objects to Standard IA class 60 days after creation
Move to Glacier for archiving after 6 months
Expiration actions - configure objects to expire (delete) after some time
Access log files can be set to delete after a 365 days
Can be used to delete old versions of files (if versioning is enabled)
Can be used to delete incomplete Multi-Part uploads
Rules can be created for a certain prefix (example: s3://mybucket/mp3/*)
Rules can be created for certain objects Tags (example: Department: Finance)
our application on EC2 creates images thumbnails after profile
notos are uploaded to Amazon S3. These thumbnails can be easily
created, and only need to be kept for 60 days. The source images
ould be able to be immediately retrieved for these 60 days, and
terwards, the user can wait up to 6 hours. How would you design
is?
53 source images can be on Standard, with a lifecycle configuration to
ansition them to Glacier after 60 days
33 thumbnails can be on One-Zone IA, with a lifecycle configuration to
pire them (delete them) after 60 days
rule in your company states that you should be able to recover your
eleted S3 objects immediately for 30 days, although this may happen
rely. After this time, and for up to 365 days, deleted objects should
e recoverable within 48 hours.
e recoverable within 48 hours. Enable S3 Versioning in order to have object versions, so that "deleted
e recoverable within 48 hours. Enable S3 Versioning in order to have object versions, so that "deleted Djects" are in fact hidden by a "delete marker" and can be recovered
e recoverable within 48 hours. Enable S3 Versioning in order to have object versions, so that "deleted objects" are in fact hidden by a "delete marker" and can be recovered fransition the "noncurrent versions" of the object to Standard IA
e recoverable within 48 hours. Enable S3 Versioning in order to have object versions, so that "deleted Djects" are in fact hidden by a "delete marker" and can be recovered
e recoverable within 48 hours. Enable S3 Versioning in order to have object versions, so that "deleted objects" are in fact hidden by a "delete marker" and can be recovered fransition the "noncurrent versions" of the object to Standard IA
e recoverable within 48 hours. Enable S3 Versioning in order to have object versions, so that "deleted objects" are in fact hidden by a "delete marker" and can be recovered fransition the "noncurrent versions" of the object to Standard IA fransition afterwards the "noncurrent versions" to Glacier Deep Archive
e recoverable within 48 hours. Enable S3 Versioning in order to have object versions, so that "deleted objects" are in fact hidden by a "delete marker" and can be recovered fransition the "noncurrent versions" of the object to Standard IA fransition afterwards the "noncurrent versions" to Glacier Deep Archive the pour decide when to transition objects to
e recoverable within 48 hours. Enable S3 Versioning in order to have object versions, so that "deleted objects" are in fact hidden by a "delete marker" and can be recovered fransition the "noncurrent versions" of the object to Standard IA fransition afterwards the "noncurrent versions" to Glacier Deep Archive elp you decide when to transition objects to e right storage class
e recoverable within 48 hours. Enable S3 Versioning in order to have object versions, so that "deleted objects" are in fact hidden by a "delete marker" and can be recovered fransition the "noncurrent versions" of the object to Standard IA fransition afterwards the "noncurrent versions" to Glacier Deep Archive elp you decide when to transition objects to e right storage class Recommendations for Standard and Standard
e recoverable within 48 hours. Enable S3 Versioning in order to have object versions, so that "deleted objects" are in fact hidden by a "delete marker" and can be recovered fransition the "noncurrent versions" of the object to Standard IA fransition afterwards the "noncurrent versions" to Glacier Deep Archive elp you decide when to transition objects to e right storage class Recommendations for Standard and Standard
e recoverable within 48 hours. Enable S3 Versioning in order to have object versions, so that "deleted objects" are in fact hidden by a "delete marker" and can be recovered fransition the "noncurrent versions" of the object to Standard IA fransition afterwards the "noncurrent versions" to Glacier Deep Archive elp you decide when to transition objects to e right storage class Recommendations for Standard and Standard Does NOT work for One-Zone IA or Glacier
e recoverable within 48 hours. Enable S3 Versioning in order to have object versions, so that "deleted objects" are in fact hidden by a "delete marker" and can be recovered fransition the "noncurrent versions" of the object to Standard IA fransition afterwards the "noncurrent versions" to Glacier Deep Archive elp you decide when to transition objects to e right storage class Recommendations for Standard and Standard Does NOT work for One-Zone IA or Glacier Report is updated daily

S3 - Requester Pays	In general, bucket owners pay for all Amazon S3 storage and data transfer costs associated with their bucket · With Requester Pays buckets, the requester instead of the bucket owner pays the cost of the request and the data download from the bucket · Helpful when you want to share large datasets with other accounts · The requester must be authenticated in AWS (cannot be anonymous)
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
S3 - Baseline Performance	Amazon S3 automatically scales to high request rates, latency 100-200 ms • Your application can achieve at least 3,500 PUT/COPY/POST/DELETE or 5,500 GET/HEAD requests per second per prefix in a bucket. • There are no limits to the number of prefixes in a bucket. • Example (object path => prefix): • bucket/folder1/sub1/file => /folder1/sub1/ • bucket/folder1/sub2/file => /folder1/sub2/ • bucket/1/file => /1/ • bucket/2/file => /2/ • If you spread reads across all four prefixes evenly, you can achieve 22,000 requests per second for GET and HEAD
S3 Performance	Multi-Part upload: • recommended for files > 100MB, must use for files > 5GB • Can help parallelize uploads (speed up transfers) S3 Transfer Acceleration • Increase transfer speed by transferring file to an AWS edge location which will forward the data to the S3 bucket in the target region • Compatible with multi-part upload
Public vs Private IP addresses	Public IP: Public IP means the machine can be identified on the internet (WWW) Must be unique across the whole web (not two machines can have the same public IP). Can be geo-located easily Private IP: Private IP means the machine can only be identified on a private network only The IP must be unique across the private network BUT two different private networks (two companies) can have the same IPs. Machines connect to WWW using a NAT + internet gateway (a proxy) Only a specified range of IPs can be used as private IP

RDS Read Replicas for read scalability	 - Up to 5 read replicas - Within AZ, Cross AZ or Cross Region - Replication is async so reads are eventually consistent - Replicas can be promoted to their own DB - Applications must update the connection string to leverage read replicas In AWS there's a network cost when data goes from one AZ to another - For RDS Read Replicas within the same region, you don't pay that fee
Routing Policies - Simple	Typically, route traffic to a single resource Can specify multiple values in the same record If multiple values are returned, a random one is chosen by the client When Alias enabled, specify only one AWS resource Can't be associated with Health Checks
Health Checks - Private Hosted Zones	Health Checks - Private Hosted Zones Route 53 health checkers are outside the VPC They can't access private endpoints (private VPC or on-premises resource) You can create a CloudWatch Metric and associate a CloudWatch Alarm, then create a Health Check that checks the alarm itself
S3 Performance - S3 Byte-Range Fetches	Parallelize GETs by requesting specific byte ranges Better resilience in case of failures
S3 Select & Glacier Select	 Retrieve less data using SQL by performing server side filtering Can filter by rows & columns (simple SQL statements) Less network transfer, less CPU cost client-side
S3 Batch Operations	Perform bulk operations on existing S3 objects with a single request, example: • Modify object metadata & properties • Copy objects between S3 buckets • Encrypt un-encrypted objects • Modify ACLs, tags • Restore objects from S3 Glacier • Invoke Lambda function to perform custom action on each object
Amazon S3 - Object Encryption	Server-Side Encryption (SSE) Server-Side Encryption with Amazon S3-Managed Keys (SSE-S3) - Enabled by Default Encrypts S3 objects using keys handled, managed, and owned by AWS Server-Side Encryption with KMS Keys stored in AWS KMS (SSE-KMS) Leverage AWS Key Management Service (AWS KMS) to manage encryption keys Server-Side Encryption with Customer-Provided Keys (SSE-C) When you want to manage your own encryption keys Client-Side Encryption

Amazon S3 Encryption - SSE-S3	Encryption using keys handled, managed, and owned by AWS Object is encrypted server-side Encryption type is AES-256 Must set header "x-amz-server-side-encryption": "AES256" Enabled by default for new buckets & new objects
Amazon S3 Encryption - SSE-KMS	Encryption using keys handled and managed by AWS KMS (Key Management Service) · KMS advantages: user control + audit key usage using CloudTrail · Object is encrypted server side · Must set header "x-amz-server-side-encryption": "aws:kms"
SSE-KMS Limitation	If you use SSE-KMS, you may be impacted by the KMS limits · When you upload, it calls the GenerateDataKey KMS API · When you download, it calls the Decrypt KMS API · Count towards the KMS quota per second (5500, 10000, 30000 req/s based on region) · You can request a quota increase using the Service Quotas Console
Amazon S3 Encryption - SSE-C	Server-Side Encryption using keys fully managed by the customer outside of AWS · Amazon S3 does NOT store the encryption key you provide · HTTPS must be used · Encryption key must provided in HTTP headers, for every HTTP request made
Amazon S3 Encryption - Client-Side Encryption	Use client libraries such as Amazon S3 Client-Side Encryption Library Clients must encrypt data themselves before sending to Amazon S3 Clients must decrypt data themselves when retrieving from Amazon S3 Customer fully manages the keys and encryption cycle
Amazon S3 - Encryption in transit (SSL/TLS)	Encryption in flight is also called SSL/TLS Amazon S3 exposes two endpoints: HTTP Endpoint - non encrypted HTTPS Endpoint - encryption in flight HTTPS is recommended HTTPS is mandatory for SSE-C Most clients would use the HTTPS endpoint by default
Amazon S3 - Default Encryption vs. Bucket Policies	SSE-S3 encryption is automatically applied to new objects stored in S3 bucket Optionally, you can "force encryption" using a bucket policy and refuse any API call to PUT an S3 object without encryption headers (SSE-KMS or SSE-C) Note: Bucket Policies are evaluated before "Default Encryption"
What is CORS?	Cross-Origin Resource Sharing (CORS) Origin = scheme (protocol) + host (domain) + port example: https://www.example.com (implied port is 443 for HTTPS, 80 for HTTP) Web Browser based mechanism to allow requests to other origins while visiting the main origin Same origin: http://example.com/app1 & http://example.com/app2 Different origins: http://www.example.com & http://other.example.com The requests won't be fulfilled unless the other origin allows for the requests, using CORS Headers (example: Access-Control-Allow-Origin)
Amazon S3 - CORS	 If a client makes a cross-origin request on our S3 bucket, we need to enable the correct CORS headers It's a popular exam question You can allow for a specific origin or for * (all origins)

Amazon S3 - MFA Delete	MFA will be required to: Permanently delete an object version Suspend Versioning on the bucket MFA won't be required to: Enable Versioning List deleted versions To use MFA Delete, Versioning must be enabled on the bucket Only the bucket owner (root account) can enable/disable MFA Delete
S3 Access Logs	For audit purpose, you may want to log all access to S3 buckets Any request made to S3, from any account, authorized or denied, will be logged into another S3 bucket That data can be analyzed using data analysis tools The target logging bucket must be in the same AWS region
S3 Access Logs: Warning	Do not set your logging bucket to be the monitored bucket It will create a logging loop, and your bucket will grow exponentially
Amazon S3 - Pre-Signed URLs	Generate pre-signed URLs using the S3 Console, AWS CLI or SDK URL Expiration S3 Console - 1 min up to 720 mins (12 hours) AWS CLI - configure expiration withexpires-in parameter in seconds (default 3600 secs, max. 604800 secs ~ 168 hours) Users given a pre-signed URL inherit the permissions of the user that generated the URL for GET / PUT Examples: Allow only logged-in users to download a premium video from your S3 bucket Allow an ever-changing list of users to download files by generating URLs dynamically Allow temporarily a user to upload a file to a precise location in your S3 bucket
S3 Glacier Vault Lock	 Adopt a WORM (Write Once Read Many) model Create a Vault Lock Policy Lock the 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 a WORM (Write Once Read Many) model Block an object version deletion for a specified amount of time Retention mode - Compliance: Object versions can't be overwritten or deleted by any user, including the root user Objects retention modes can't be changed, and retention periods can't be shortened Retention mode - Governance: Most users can't overwrite or delete an object version or alter its lock settings Some users have special permissions to change the retention or delete the object Retention Period: protect the object for a fixed period, it can be extended Legal Hold: protect the object indefinitely, independent from retention period can be freely placed and removed using the s3:PutObjectLegalHold IAM permission
S3 - Access Points	Access Points simplify security management for S3 Buckets Each Access Point has: its own DNS name (Internet Origin or VPC Origin) an access point policy (similar to bucket policy) - manage security at scale

S3 - Access Points - VPC Origin	· We can define the access
	point to be accessible
	only from within the VPC
	· You must create a VPC
	Endpoint to access the
	Access Point (Gateway
	or Interface Endpoint)
	• The VPC Endpoint Policy
	must allow access to the
	target
	Use AWS Lambda Functions to
	change the object before it is
	retrieved by the caller application
	Only one S3 bucket is needed, on
	top of which we create S3 Access
	Point and S3 Object Lambda Access
	Points.
S3 Object Lambda	· Use Cases:
33 Object Lambua	Redacting personally identifiable
	information for analytics or nonproduction
	environments.
	· Converting across data formats, such
	as converting XML to JSON.
	· Resizing and watermarking images on
	the fly using caller-specific details, such
	as the user who requested the object.
	Content Delivery Network (CDN)
	· Improves read performance,
Amazana Claudifuant	content is cached at the edge
Amazon CloudFront	DDoS protection (because
	worldwide), integration with Shield,
	AWS Web Application Firewall
	S3 bucket
	· For distributing files and caching them at the edge
	• Enhanced security with CloudFront Origin Access Control (OAC)
	• OAC is replacing Origin Access Identity (OAI)
	· CloudFront can be used as an ingress (to upload files to S3)
CloudFront - Origins	· Custom Origin (HTTP)
	Application Load Balancer
	• EC2 instance
	· S3 website (must first enable the bucket as a static S3 website)
	· Any HTTP backend you want
	CloudFront:
	Global Edge network
	 Global Edge network Files are cached for a TTL (maybe a day)
	· Files are cached for a TTL (maybe a day)
	Files are cached for a TTL (maybe a day)Great for static content that must be available everywhere
CloudFront vs S3 Cross Region Replication	 Files are cached for a TTL (maybe a day) Great for static content that must be available everywhere S3 Cross Region Replication:
CloudFront vs S3 Cross Region Replication	 Files are cached for a TTL (maybe a day) Great for static content that must be available everywhere S3 Cross Region Replication: Must be setup for each region you want replication to happen
CloudFront vs S3 Cross Region Replication	 Files are cached for a TTL (maybe a day) Great for static content that must be available everywhere S3 Cross Region Replication: Must be setup for each region you want replication to happen Files are updated in near real-time
CloudFront vs S3 Cross Region Replication	 Files are cached for a TTL (maybe a day) Great for static content that must be available everywhere S3 Cross Region Replication: Must be setup for each region you want replication to happen

 Allowlist: Allow your users to access your content only if they're in one of the countries on a list of approved countries. Blocklist: Prevent your users from accessing your content if they're in one of the countries on a list of banned countries. The "country" is determined using a 3rd party Geo-IP database Use case: Copyright Laws to control access to content CloudFront - Pricing CloudFront Edge locations are all around the world The cost of data out per edge location varies You can reduce the number of edge locations for cost reduction 	
CloudFront Geo Restriction • Blocklist: Prevent your users from accessing your content if they're in one of the countries on a list of banned countries. • The "country" is determined using a 3rd party Geo-IP database • Use case: Copyright Laws to control access to content CloudFront - Pricing CloudFront Edge locations are all around the world • The cost of data out per edge location varies	
countries on a list of banned countries. • The "country" is determined using a 3rd party Geo-IP database • Use case: Copyright Laws to control access to content CloudFront - Pricing CloudFront - Pricing • The cost of data out per edge location varies	
The "country" is determined using a 3rd party Geo-IP database Use case: Copyright Laws to control access to content CloudFront - Pricing CloudFront - Pricing The cost of data out per edge location varies	
Use case: Copyright Laws to control access to content CloudFront - Pricing CloudFront - Pricing The cost of data out per edge location varies	
CloudFront - Pricing CloudFront - Pricing CloudFront Edge locations are all around the world • The cost of data out per edge location varies	
· The cost of data out per edge location varies	
• The cost of data out per edge location varies	
You can reduce the number of edge locations for cost reduction	
•Three price classes:	
CloudFront - Price Classes 1. Price Class All: all regions - best performance	
2. Price Class 200: most regions, but excludes the most expensive regions	
3. Price Class 100: only the least expensive regions	
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	
CloudFront - Cache Invalidations · 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/*)	
Unicast IP: one server holds one IP	
address	
Unicast IP vs Anycast IP • Anycast IP: all servers hold the same	
IP address and the client is routed to	
the nearest one	
Leverage the AWS internal	
network to route to your	
application	
• 2 Anycast IP are created for your	
application	
• The Anycast IP send traffic directly	
to Edge Locations	
• The Edge locations send the traffic	
to your application	
Works with Elastic IP, EC2 instances, ALB, NLB, public or private	
AWS Global Accelerator • Consistent Performance	
 Intelligent routing to lowest latency and fast regional failover 	
 No issue with client cache (because the IP doesn't change) 	
Internal AWS network	
Health Checks	
· Global Accelerator performs a health check of your applications	
· Helps make your application global (failover less than 1 minute for unhealthy)	
Great for disaster recovery (thanks to the health checks)	
• Security	
only 2 external IP need to be whitelisted	
• DDoS protection thanks to AWS Shield	

	They both use the AWS global network and its edge locations around the world
	Both services integrate with AWS Shield for DDoS protection.
	• CloudFront
	· Improves performance for both cacheable content (such as images and videos)
	Dynamic content (such as API acceleration and dynamic site delivery)
AWS Global Accelerator vs CloudFront	· Content is served at the edge
A THE CLOSE AT A COLOR OF THE ACT	Global Accelerator
	· Improves performance for a wide range of applications over TCP or UDP
	Proxying packets at the edge to applications running in one or more AWS Regions.
	Good fit for non-HTTP use cases, such as gaming (UDP), IoT (MQTT), or Voice over IP
	Good for HTTP use cases that require static IP addresses
	Good for HTTP use cases that required deterministic, fast regional failover
	Highly-secure, portable devices to collect and process data at the edge, and migrate
	data into and out of AWS
AWS Snow Family	
	Data Migration: Snowcone, Snowball Edge, Snowmobile
	Edge Computing: Snowcone, Snowball Edge
Data Migrations with AMC Co F	AWS Snow Family: offline devices to perform data migrations
Data Migrations with AWS Snow Family	If it takes more than a week to transfer over the network, use Snowball devices!
	Physical data transport solution: move TBs or PBs of data in or out of AWS
	· Alternative to moving data over the network (and paying network fees)
	· Pay per data transfer job
	Provide block storage and Amazon S3-compatible object storage
Chauladi Edwa (fan data transfara)	Snowball Edge Storage Optimized
Snowball Edge (for data transfers)	· 80 TB of HDD capacity for block volume and S3 compatible object storage
	Snowball Edge Compute Optimized
	 42 TB of HDD or 28TB NVMe capacity for block volume and S3 compatible object
	storage
	Use cases: large data cloud migrations, DC decommission, disaster recovery
	Small, portable computing, anywhere, rugged &
	secure, withstands harsh environments
	· Light (4.5 pounds, 2.1 kg)
	Device used for edge computing, storage, and data
	transfer
AWS Snowcone & Snowcone SSD	Snowcone - 8 TB of HDD Storage
	Snowcone SSD - 14 TB of SSD Storage
	Use Snowcone where Snowball does not fit (spaceconstrained
	environment)
	Must provide your own battery / cables
	Can be sent back to AWS offline, or connect it to
	internet and use AWS DataSync to send data
	Transfer exabytes of data (1 EB = 1,000 PB = 1,000,000 TBs)
AWS Snowmobile	• Each Snowmobile has 100 PB of capacity (use multiple in parallel)
	High security: temperature controlled, GPS, 24/7 video surveillance
	Better than Snowball if you transfer more than 10 PB
	1. Request Snowball devices from the AWS console for delivery
	2. Install the snowball client / AWS OpsHub on your servers
	3. Connect the snowball to your servers and copy files using the client
Snow Family - Usage Process	4. Ship back the device when you're done (goes to the right AWS
	facility)
	5. Data will be loaded into an S3 bucket
	6. Snowball is completely wiped

What is Edge Computing?	Process data while it's being created on an edge location A truck on the road, a ship on the sea, a mining station underground These locations may have Limited / no internet access Limited / no easy access to computing power We setup a Snowball Edge / Snowcone device to do edge computing Use cases of Edge Computing: Preprocess data Machine learning at the edge Transcoding media streams Eventually (if need be) we can ship back the device to AWS (for transferring data for example)
Snow Family - Edge Computing	Snowcone & Snowcone SSD (smaller) 2 CPUs, 4 GB of memory, wired or wireless access USB-C power using a cord or the optional battery Snowball Edge - Compute Optimized 104 vCPUs, 416 GiB of RAM Optional GPU (useful for video processing or machine learning) 28 TB NVMe or 42TB HDD usable storage Storage Clustering available (up to 16 nodes) Snowball Edge - Storage Optimized Up to 40 vCPUs, 80 GiB of RAM, 80 TB storage All: Can run EC2 Instances & AWS Lambda functions (using AWS IoT Greengrass) Long-term deployment options: 1 and 3 years discounted pricing
AWS OpsHub	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))
Solution Architecture: Snowball into Glacier	Snowball cannot import to Glacier directly You must use Amazon S3 first, in combination with an S3 lifecycle policy
Amazon FSx - Overview	Launch 3rd party high-performance file systems on AWSFully managed service
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

Amazon FSx for Lustre	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)
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
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)
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)

	AWS is pushing for "hybrid cloud"
	Part of your infrastructure is on the cloud
	Part of your infrastructure is on-premises
Hybrid Cloud for Storage	• This can be due to
	Long cloud migrations
	· Security requirements
	Compliance requirements
	• IT strategy
	AWS Storage Gateway!
	Bridge between on-premises data and cloud data
	• Use cases:
	- disaster recovery
	backup & restore
	• tiered storage
AWS Storage Gateway	• on-premises cache & low-latency files access
, , , , , , , , , , , , , , , , , , , ,	• Types of Storage Gateway:
	• S3 File Gateway
	• FSx File Gateway
	· Volume Gateway
	· Tape Gateway
	Configured S3 buckets are accessible using the NFS and SMB protocol
	Most recently used data is cached in the file gateway
Amazon S3 File Gateway	Supports S3 Standard, S3 Standard IA, S3 One Zone A, S3 Intelligent Tiering
/ mazon oo rike odkemay	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
	Native access to Amazon FSx for Windows File Server
Amazon FSx File Gateway	Local cache for frequently accessed data
Amazon 13x File Galeway	· Windows native compatibility (SMB, NTFS, Active Directory)
	Useful for group file shares and home directories
	Block storage using iSCSI protocol backed by S3
Volume Gateway	 Backed by EBS snapshots which can help restore on-premises volumes!
Volume Galeway	Cached volumes: low latency access to most recent data
	Stored volumes: entire dataset is on premise, scheduled backups to S3
	Some companies have backup processes using physical tapes (!)
	· With Tape Gateway, companies use the same processes but, in the cloud
Tape Gateway	· 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
	Using Storage Gateway means you need
	on-premises virtualization
	Otherwise, you can use a Storage
	Gateway Hardware Appliance
	You can buy it on amazon.com
Storage Gateway - Hardware appliance	Works with File Gateway, Volume Gateway,
	Tape Gateway
	Has the required CPU, memory, network,
	SSD cache resources
	Helpful for daily NFS backups in small data
	centers

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,
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
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
Amazon SQS - Standard Queue	Attributes: • Unlimited throughput, unlimited number of messages in queue • Default retention of messages: 4 days, maximum of 14 days • Low latency (<10 ms on publish and receive) • Limitation of 256KB per message sent • Can have duplicate messages (at least once delivery, occasionally) • Can have out of order messages (best effort ordering)
SQS - Producing Messages	Produced to SQS using the SDK (SendMessage API) The message is persisted in SQS until a consumer deletes it Message retention: default 4 days, up to 14 days
SQS - Consuming Messages	Consumers (running on EC2 instances, servers, or AWS Lambda) Poll SQS for messages (receive up to 10 messages at a time) Process the messages (example: insert the message into an RDS database) Delete the messages using the DeleteMessage API

SQS - Multiple EC2 Instances Consumers	Consumers receive and process messages in parallel At least once delivery Best-effort message ordering Consumers delete messages after processing them We can scale consumers horizontally to improve throughput of processing
Amazon SQS - Security	Encryption: In-flight encryption using HTTPS API At-rest encryption using KMS keys Client-side encryption if the client wants to perform encryption/decryption itself Access Controls: IAM policies to regulate access to the SQS API SQS Access Policies (similar to S3 bucket policies) Useful for cross-account access to SQS queues
SQS - Message Visibility Timeout	Useful for allowing other services (SNS, S3) to write to an SQS queue After a message is polled by a consumer, it becomes invisible to other consumers By default, the "message visibility timeout" is 30 seconds That means the message has 30 seconds to be processed After the message visibility timeout is over, the message is "visible" in SQS If a message is not processed within the visibility timeout, it will be processed twice A consumer could call the ChangeMessageVisibility API to get more time If visibility timeout is high (hours), and consumer crashes, re-processing will take time If visibility timeout is too low (seconds), we may get duplicates
Amazon SQS - Long Polling	When a consumer requests messages from the queue, it can optionally "wait" for messages to arrive if there are none in the queue This is called Long Polling LongPolling decreases the number of API calls made to SQS while increasing the efficiency and reducing latency of your application The wait time can be between 1 sec to 20 sec (20 sec preferable) Long Polling is preferable to Short Polling
Amazon SQS - FIFO Queue	FIFO = First In First Out (ordering of messages in the queue) Limited throughput: 300 msg/s without batching, 3000 msg/s with Exactly-once send capability (by removing duplicates) Messages are processed in order by the consumer
Amazon SNS	The "event producer" only sends message to one SNS topic As many "event receivers" (subscriptions) as we want to listen to the SNS topic notifications Each subscriber to the topic will get all the messages (note: new feature to filter messages) Up to 12,500,000 subscriptions per topic 100,000 topics limit

Amazon SNS - How to publish	Topic Publish (using the SDK) Create a topic Create a subscription (or many) Publish to the topic Direct Publish (for mobile apps SDK) Create a platform application Create a platform endpoint Publish to the platform endpoint Works with Google GCM, Apple APNS, Amazon ADM
Amazon SNS - Security	Encryption: In-flight encryption using HTTPS API At-rest encryption using KMS keys Client-side encryption if the client wants to perform encryption/decryption itself Access Controls: IAM policies to regulate access to the SNS API SNS Access Policies (similar to S3 bucket policies) Useful for cross-account access to SNS topics Useful for allowing other services (S3) to write to an SNS topic
SNS + SQS: Fan Out	Push once in SNS, receive in all SQS queues that are subscribers Fully decoupled, no data loss SQS allows for: data persistence, delayed processing and retries of work Ability to add more SQS subscribers over time Make sure your SQS queue access policy allows for SNS to write Cross-Region Delivery: works with SQS Queues in other regions For the same combination of: event type (e.g. object create) and prefix (e.g. images/) you can only have one S3 Event rule If you want to send the same S3 event to many SQS queues, use fan-out
Amazon SNS - FIFO Topic	FIFO = First In First Out (ordering of messages in the topic) Similar features as SQS FIFO: Ordering by Message Group ID (all messages in the same group are ordered) Deduplication using a Deduplication ID or Content Based Deduplication Can only have SQS FIFO queues as subscribers Limited throughput (same throughput as SQS FIFO)
SNS FIFO + SQS FIFO: Fan Out	In case you need fan out + ordering + deduplication
SNS - Message Filtering	JSON policy used to filter messages sent to SNS topic's subscriptions • If a subscription doesn't have a filter policy, it receives every message
Kinesis Data Firehose	load data streams into AWS data stores
Kinesis Data Streams	Capture, process, and store data streams
Kinesis Data Analytics	Analyze data streams with SQL or Apache Flink
Kinesis Video Streams	Capture, process, and store video streams
Kinesis Data Streams explanation	Retention between 1 day to 365 days Ability to reprocess (replay) data Once data is inserted in Kinesis, it can't be deleted (immutability) Data that shares the same partition goes to the same shard (ordering) Producers: AWS SDK, Kinesis Producer Library (KPL), Kinesis Agent Consumers: Write your own: Kinesis Client Library (KCL), AWS SDK Managed: AWS Lambda, Kinesis Data Firehose, Kinesis Data Analytics,

Kinesis Data Streams - Capacity Modes	Provisioned mode: You choose the number of shards provisioned, scale manually or using API Each shard gets 1MB/s in (or 1000 records per second) Each shard gets 2MB/s out (classic or enhanced fan-out consumer) You pay per shard provisioned per hour On-demand mode: No need to provision or manage the capacity Default capacity provisioned (4 MB/s in or 4000 records per second) Scales automatically based on observed throughput peak during the last 30 days Pay per stream per hour & data in/out per GB
Kinesis Data Streams Security	Control access / authorization using IAM policies • Encryption in flight using HTTPS endpoints • Encryption at rest using KMS • You can implement encryption/decryption of data on client side (harder) • VPC Endpoints available for Kinesis to access within VPC • Monitor API calls using CloudTrail
Kinesis Data Firehose description	Fully Managed Service, no administration, automatic scaling, serverless - AWS: Redshift / Amazon S3 / OpenSearch - 3rd party partner: Splunk / MongoDB / DataDog / NewRelic / - Custom: send to any HTTP endpoint - Pay for data going through Firehose - Near Real Time - 60 seconds latency minimum for non full batches - Or minimum 1 MB of data at a time - Supports many data formats, conversions, transformations, compression - Supports custom data transformations using AWS Lambda - Can send failed or all data to a backup S3 bucket
Kinesis Data Streams vs Firehose	Kinesis Data Streams Streaming service for ingest at scale · Write custom code (producer / consumer) · Real-time (-200 ms) · Manage scaling (shard splitting / merging) · Data storage for 1 to 365 days · Supports replay capability Kinesis Data Firehose Load streaming data into S3 / Redshift / OpenSearch / 3rd party / custom HTTP · Fully managed · Near real-time (buffer time min. 60 sec) · Automatic scaling · No data storage · Doesn't support replay capability
Ordering data into Kinesis	Answer: send using a "Partition Key" value of the "truck_id" • The same key will always go to the same shard

	For SQS standard, there is no ordering.
Ordering data into SQS	· For SQS FIFO, if you don't use a Group ID, messages are consumed in the
	order they are sent, with only one consumer
	You want to scale the number of consumers, but you want messages to be "grouped"
	when they are related to each other
	• Then you use a Group ID (similar to Partition Key in Kinesis)
	Let's assume 100 trucks, 5 kinesis shards, 1 SQS FIFO
	Kinesis Data Streams:
	· On average you'll have 20 trucks per shard
	• Trucks will have their data ordered within each shard
	• The maximum amount of consumers in parallel we can have is 5
Kinesis vs SQS ordering	· Can receive up to 5 MB/s of data
	· SQS FIFO
	· You only have one SQS FIFO queue
	· You will have 100 Group ID
	· You can have up to 100 Consumers (due to the 100 Group ID)
	· You have up to 300 messages per second (or 3000 if using batching)
	SQS:
	• Consumer "pull data"
	Data is deleted after being
	consumed
	· Can have as many workers
	(consumers) as we want
	No need to provision
	throughput
	· Ordering guarantees only on
	FIFO queues
	· Individual message delay
	capability
	SNS:
	• Push data to many
	subscribers
	• Up to 12,500,000 subscribers
	• Data is not persisted (lost if
	Data is not persisted (tost ii
SOS vs SNS vs Kinesis	not delivered)
SQS vs SNS vs Kinesis	not delivered) - Pub/Sub
SQS vs SNS vs Kinesis	not delivered) • Pub/Sub • Up to 100,000 topics
SQS vs SNS vs Kinesis	not delivered) • Pub/Sub • Up to 100,000 topics • No need to provision
SQS vs SNS vs Kinesis	not delivered) • Pub/Sub • Up to 100,000 topics • No need to provision throughput
SQS vs SNS vs Kinesis	not delivered) Pub/Sub Up to 100,000 topics No need to provision throughput Integrates with SQS for fanout
SQS vs SNS vs Kinesis	not delivered) • Pub/Sub • Up to 100,000 topics • No need to provision throughput • Integrates with SQS for fanout architecture pattern
SQS vs SNS vs Kinesis	not delivered) • Pub/Sub • Up to 100,000 topics • No need to provision throughput • Integrates with SQS for fanout architecture pattern • FIFO capability for SQS FIFO
SQS vs SNS vs Kinesis	not delivered) • Pub/Sub • Up to 100,000 topics • No need to provision throughput • Integrates with SQS for fanout architecture pattern • FIFO capability for SQS FIFO Kinesis:
SQS vs SNS vs Kinesis	not delivered) Pub/Sub Up to 100,000 topics No need to provision throughput Integrates with SQS for fanout architecture pattern FIFO capability for SQS FIFO Kinesis: Standard: pull data
SQS vs SNS vs Kinesis	not delivered) • Pub/Sub • Up to 100,000 topics • No need to provision throughput • Integrates with SQS for fanout architecture pattern • FIFO capability for SQS FIFO Kinesis: • Standard: pull data • 2 MB per shard
SQS vs SNS vs Kinesis	not delivered) • Pub/Sub • Up to 100,000 topics • No need to provision throughput • Integrates with SQS for fanout architecture pattern • FIFO capability for SQS FIFO Kinesis: • Standard: pull data • 2 MB per shard • Enhanced-fan out: push data
SQS vs SNS vs Kinesis	not delivered) Pub/Sub Up to 100,000 topics No need to provision throughput Integrates with SQS for fanout architecture pattern FIFO capability for SQS FIFO Kinesis: Standard: pull data MB per shard Enhanced-fan out: push data MB per shard per consumer
SQS vs SNS vs Kinesis	not delivered) Pub/Sub Up to 100,000 topics No need to provision throughput Integrates with SQS for fanout architecture pattern FIFO capability for SQS FIFO Kinesis: Standard: pull data 2 MB per shard Enhanced-fan out: push data 2 MB per shard per consumer Possibility to replay data
SQS vs SNS vs Kinesis	not delivered) Pub/Sub Up to 100,000 topics No need to provision throughput Integrates with SQS for fanout architecture pattern FIFO capability for SQS FIFO Kinesis: Standard: pull data 2 MB per shard Enhanced-fan out: push data 2 MB per shard per consumer Possibility to replay data Meant for real-time big data,
SQS vs SNS vs Kinesis	not delivered) Pub/Sub Up to 100,000 topics No need to provision throughput Integrates with SQS for fanout architecture pattern FIFO capability for SQS FIFO Kinesis: Standard: pull data Me per shard Enhanced-fan out: push data Me per shard per consumer Possibility to replay data Meant for real-time big data, analytics and ETL
SQS vs SNS vs Kinesis	not delivered) Pub/Sub Up to 100,000 topics No need to provision throughput Integrates with SQS for fanout architecture pattern FIFO capability for SQS FIFO Kinesis: Standard: pull data MB per shard Enhanced-fan out: push data MB per shard per consumer Possibility to replay data Meant for real-time big data, analytics and ETL Ordering at the shard level
SQS vs SNS vs Kinesis	not delivered) Pub/Sub Up to 100,000 topics No need to provision throughput Integrates with SQS for fanout architecture pattern FIFO capability for SQS FIFO Kinesis: Standard: pull data Me per shard Enhanced-fan out: push data Me per shard per consumer Possibility to replay data Meant for real-time big data, analytics and ETL

Amazon MQ	SQS, SNS are "cloud-native" services: proprietary protocols from AWS • Traditional applications running from on-premises may use open protocols such as: MQTT, AMQP, STOMP, Openwire, WSS • When migrating to the cloud, instead of re-engineering the application to use SQS and SNS, we can use Amazon MQ • Amazon MQ is a managed message broker service for RabbitMQ and ActiveMQ Amazon MQ doesn't "scale" as much as SQS / SNS • Amazon MQ runs on servers, can run in Multi-AZ with failover • Amazon MQ has both queue feature (-SQS) and topic features (-SNS)
What is Docker?	Docker is a software development platform to deploy apps Apps are packaged in containers that can be run on any OS Apps run the same, regardless of where they're run Any machine No compatibility issues Predictable behavior Less work Easier to maintain and deploy Works with any language, any OS, any technology
Where are Docker images stored?	Docker images are stored in Docker Repositories Docker Hub (https://hub.docker.com) Public repository Find base images for many technologies or OS (e.g., Ubuntu, MySQL,) Amazon ECR (Amazon Elastic Container Registry) Private repository Public repository (Amazon ECR Public Gallery https://gallery.ecr.aws)
Docker vs. Virtual Machines	Docker is "sort of" a virtualization technology, but not exactly Resources are shared with the host => many containers on one server
Amazon Elastic Container Service (Amazon ECS)	Amazon's own container platform
Amazon Elastic Kubernetes Service (Amazon EKS)	Amazon's managed Kubernetes (open source)
AWS Fargate	Amazon's own Serverless container platform · Works with ECS and with EKS
Amazon ECR (Elastic Container Registry)	Store container images
Amazon ECS - EC2 Launch Type	Launch Docker containers on AWS = Launch ECS Tasks on ECS Clusters • EC2 Launch Type: you must provision & maintain the infrastructure (the EC2 instances) • Each EC2 Instance must run the ECS Agent to register in the ECS Cluster • AWS takes care of starting / stopping containers
Amazon ECS - Fargate Launch Type	You do not provision the infrastructure (no EC2 instances to manage) • It's all Serverless! • You just create task definitions • AWS just runs ECS Tasks for you based on the CPU / RAM you need • To scale, just increase the number of tasks. Simple - no more EC2 instances

Amazon ECS - IAM Roles for ECS	EC2 Instance Profile (EC2 Launch Type
	only):
	· Used by the ECS agent
	· Makes API calls to ECS service
	· Send container logs to CloudWatch Logs
	· Pull Docker image from ECR
	· Reference sensitive data in Secrets Manager or
	SSM Parameter Store
	• ECS Task Role:
	· Allows each task to have a specific role
	· Use different roles for the different ECS Services
	you run
	· Task Role is defined in the task definition
	Application Load Balancer supported
	and works for most use cases
500 1 10 1 11	Network Load Balancer recommended
Amazon ECS - Load Balancer Integrations	only for high throughput / high
	performance use cases, or to pair it with
	AWS Private Link
	Mount EFS file systems onto ECS tasks
	• Works for both EC2 and Fargate launch types
Amazon ECS - Data Volumes (EFS)	• Tasks running in any AZ will share the same data
,	in the EFS file system
	• Fargate + EFS = Serverless
	Automatically increase/decrease the desired number of ECS tasks
	· Amazon ECS Auto Scaling uses AWS Application Auto Scaling
	• ECS Service Average CPU Utilization
	• ECS Service Average Memory Utilization - Scale on RAM
	• ALB Request Count Per Target - metric coming from the ALB
ECS Service Auto Scaling	Target Tracking - scale based on target value for a specific CloudWatch metric
	Step Scaling - scale based on a specified CloudWatch Alarm
	Scheduled Scaling - scale based on a specified date/time (predictable changes)
	• ECS Service Auto Scaling (task level) \neq EC2 Auto Scaling (EC2 instance level)
	Fargate Auto Scaling is much easier to setup (because Serverless)
	Accommodate ECS Service Scaling by adding underlying EC2 Instances
	· Auto Scaling Group Scaling
L	Scale your ASG based on CPU Utilization
EC2 Launch Type - Auto Scaling EC2	· Add EC2 instances over time
Instances	· ECS Cluster Capacity Provider
	 Used to automatically provision and scale the infrastructure for your ECS Tasks
	· Capacity Provider paired with an Auto Scaling Group
	· Add EC2 Instances when you're missing capacity (CPU, RAM)
	ECR = Elastic Container Registry
	· Store and manage Docker images on AWS
	Private and Public repository (Amazon ECR
	Public Gallery https://gallery.ecr.aws)
Amazon ECR	· Fully integrated with ECS, backed by Amazon S3
	· Access is controlled through IAM (permission
	errors => policy)
	errors - policy)
	Supports image vulnerability scanning, versioning,

Amazon EKS Overview	Amazon EKS = Amazon Elastic Kubernetes Service
	· It is a way to launch managed Kubernetes clusters on AWS
	 Kubernetes is an open-source system for automatic deployment, scaling and
	management of containerized (usually Docker) application
	• It's an alternative to ECS, similar goal but different API
	• EKS supports EC2 if you want to deploy worker nodes or Fargate to deploy
	serverless containers
	· Use case: if your company is already using Kubernetes on-premises or in
	another cloud, and wants to migrate to AWS using Kubernetes
	· Kubernetes is cloud-agnostic (can be used in any cloud - Azure, GCP)
	· For multiple regions, deploy one EKS cluster per region
	 Collect logs and metrics using CloudWatch Container Insights
	Managed Node Groups
	· Creates and manages Nodes (EC2 instances) for you
	 Nodes are part of an ASG managed by EKS
	Supports On-Demand or Spot Instances
	Self-Managed Nodes
Amazon EKS - Node Types	 Nodes created by you and registered to the EKS cluster and managed by an ASG
	· You can use prebuilt AMI - Amazon EKS Optimized AMI
	· Supports On-Demand or Spot Instances
	• AWS Fargate
	No maintenance required; no nodes managed
	Need to specify StorageClass manifest on your EKS cluster
Amazon EKS - Data Volumes	Leverages a Container Storage Interface (CSI) compliant driver
	Fully managed service that makes it easy to deploy web
	applications and APIs at scale
	No infrastructure experience required
	Start with your source code or container image
AWS App Runner	Automatically builds and deploy the web app
	 Automatic scaling, highly available, load balancer, encryption
	· VPC access support
	· Connect to database, cache, and message queue services
	 Use cases: web apps, APIs, microservices, rapid production
	deployments
	AWS Lambda
	· DynamoDB
	• AWS Cognito
	· AWS API Gateway
	· Amazon S3
Serverless in AWS	· AWS SNS & SQS
	· AWS Kinesis Data Firehose
	Aurora Serverless
	Step Functions
	• Fargate
	Virtual functions - no servers to manage!
	Limited by time - short executions
Why AWS Lambda	• Run on-demand
	- Scaling is automated!

Benefits of AWS Lambda	Easy Pricing:
	Pay per request and compute time
	• Free tier of 1,000,000 AWS Lambda requests and 400,000 GBs of compute time
	 Integrated with the whole AWS suite of services
	Integrated with many programming languages
	Easy monitoring through AWS CloudWatch
	• Easy to get more resources per functions (up to 10GB of RAM!)
	Increasing RAM will also improve CPU and network!
	Node.js (JavaScript)
	• Python
	Java (Java 8 compatible)
	· C# (.NET Core)
AWS Lambda language support	· Golang
	· C# / Powershell
	• Ruby
	Custom Runtime API (community supported, example Rust)
Lambda Container Image	The container image must implement the Lambda Runtime API
Lambda Container image	• ECS / Fargate is preferred for running arbitrary Docker images
	- API Gateway
	- Kinesis
	- DynamoDB
	- S3
AWS Lambda Integrations	- CloudFront
3	- CloudWatch Events EventBridge
	- SNS
	- SQS
	- Cognito
	Pay per calls:
	• First 1,000,000 requests are free
	• \$0.20 per 1 million requests thereafter (\$0.0000002 per request)
AWS Lambda Pricing: evample	Pay per duration: (in increment of 1 ms)
AWS Lambda Pricing: example	• 400,000 GB-seconds of compute time per month for FREE
	• == 400,000 seconds if function is IGB RAM
	== 3,200,000 seconds if function is 128 MB RAM
	• After that \$1.00 for 600,000 GB-seconds
	Execution:
	Memory allocation: 128 MB - 10GB (1 MB increments)
	Maximum execution time: 900 seconds (15 minutes)
AWS Lambda Limits to Know - per region	• Environment variables (4 KB)
	Disk capacity in the "function container" (in /tmp): 512 MB to 10GB
	Concurrency executions: 1000 (can be increased)
ATTO Editional Elittics to Know - per region	
	Deployment: Learning of transition deployment size (compressed zip) 50 MP.
	Lambda function deployment size (compressed .zip): 50 MB
	Size of uncompressed deployment (code + dependencies): 250 MB
	Can use the /tmp directory to load other files at startup
	· Size of environment variables: 4 KB

	Many modern applications execute some form of the legis at the edge
Customization At The Edge	Many modern applications execute some form of the logic at the edge
	Edge Function: A code that you write and attach to CloudFront distributions.
	A code that you write and attach to CloudFront distributions Purps close to your ways to minimize latency.
	Runs close to your users to minimize latency
	CloudFront provides two types: CloudFront Functions &
	Lambda@Edge
	You don't have to manage any servers, deployed globally
	Use case: customize the CDN content
	Website Security and Privacy
	Dynamic Web Application at the Edge
	· Search Engine Optimization (SEO)
	 Intelligently Route Across Origins and Data Centers
CloudFront Functions & Lambda@Edge Use	Bot Mitigation at the Edge
Cases	· Real-time Image Transformation
	- A/B Testing
	User Authentication and Authorization
	User Prioritization
	User Tracking and Analytics
	Lightweight functions written in JavaScript
	For high-scale, latency-sensitive CDN customizations
	Sub-ms startup times, millions of requests/second
	· Used to change Viewer requests and responses:
	Viewer Request: after CloudFront receives a request from a
CloudFront Functions	viewer
	Viewer Response: before CloudFront forwards the response
	to the viewer
	Native feature of CloudFront (manage code entirely
	within CloudFront)
	Lambda functions written in NodeJS or Python
	Scales to 1000s of requests/second
	Used to change CloudFront requests and responses:
	Viewer Request - after CloudFront receives a request from a
	viewer
	Origin Request - before CloudFront forwards the request to the
Lambda@Edge	origin
LambadeLage	Origin Response - after CloudFront receives the response from
	the origin
	Viewer Response - before CloudFront forwards the response to
	the viewer
	Author your functions in one AWS Region (us-east-I), then CloudFront replicator to its locations.
	CloudFront replicates to its locations

	CloudFront Functions:
	- JavaScript
	- Millions of requests per second
	Triggers:
	- Viewer Request/Response
	Maximum Execution Time: < 1 ms
	Maximum Memory: 2 MB
	Total Package Size: 10 KB
	Network/File System Access: NO
	Access to the request body: NO
	Free tier available, 1/6th price of @Edge
CloudFront Functions vs. Lambda@Edge	
	Lambda@Edge:
	- Node,js, Python
	- Thousands of requests per second
	Triggers:
	- Viewer Request/Response
	- Origin Request/Response
	Maximum Execution Time: 5 - 10 seconds
	Maximum Memory: 128 MB up to 10 GB
	Total Package Size: 1 MB - 50 MB
	Network/File System Access: YES
	Access to the request body: YES
	No free tier, charged per request & duration
	CloudFront Functions
	Cache key normalization
	Transform request attributes (headers, cookies, query strings, URL) to create an
	optimal Cache Key
	· Header manipulation
	· Insert/modify/delete HTTP headers in the
	request or response
	URL rewrites or redirects
CloudFront Functions vs. Lambda@Edge -	Request authentication & authorization
Use Cases	· Create and validate user-generated tokens (e.g., JWT) to allow/deny requests
	Lambda@Edge
	· Longer execution time (several ms)
	· Adjustable CPU or memory
	· Your code depends on a 3rd libraries (e.g., AWS SDK to access other AWS services)
	Network access to use external services for processing
	• File system access or access to the body of HTTP requests
	By default, your Lambda function is
	launched outside your own VPC (in
	an AWS-owned VPC)
	· Therefore, it cannot access resources
Lambda by default	in your VPC (RDS, ElastiCache,
,	internal ELB)
	You must define the VPC ID, the
	Subnets and the Security Groups
	Lambda will create an ENI (Elastic
	Network Interface) in your subnets

	If Lambda functions directly access your
	database, they may open too many
	connections under high load
	· RDS Proxy
	Improve scalability by pooling and sharing DB
	connections
Lambda with RDS Proxy	· Improve availability by reducing by 66% the
	failover time and preserving connections
	· Improve security by enforcing IAM
	authentication and storing credentials in
	Secrets Manager
	• The Lambda function must be deployed
	in your VPC, because RDS Proxy is never
	publicly accessible
	Invoke Lambda functions from within your DB instance
	Allows you to process data events from within a
	database
	Supported for RDS for PostgreSQL and Aurora MySQL
Invoking Lambda from RDS & Aurora	Must allow outbound traffic to your Lambda function
3	from within your DB instance (Public, NAT GW, VPC
	Endpoints)
	DB instance must have the required permissions to
	invoke the Lambda function (Lambda Resource-based
	Policy & IAM Policy)
	Notifications that tells information about the DB
	instance itself (created, stopped, start,)
	· You don't have any information about the data itself
	Subscribe to the following event categories: DB
RDS Event Notifications	instance, DB snapshot, DB Parameter Group, DB
TO EVERT Notifications	Security Group, RDS Proxy, Custom Engine Version
	Near real-time events (up to 5 minutes)
	Send notifications to SNS or subscribe to events
	using EventBridge
	Using Eventuringe
	Fully managed, highly available with replication across multiple AZs
	 NoSQL database - not a relational database - with transaction support
	Scales to massive workloads, distributed database
	 Millions of requests per seconds, trillions of row, 100s of TB of storage
Amazon DynamoDB	Fast and consistent in performance (single-digit millisecond)
	 Integrated with IAM for security, authorization and administration
	Low cost and auto-scaling capabilities
	No maintenance or patching, always available
	Standard & Infrequent Access (IA) Table Cl
	DynamoDB is made of Tables
	Each table has a Primary Key (must be decided at creation time)
	• Each table can have an infinite number of items (= rows)
	Each item has attributes (can be added over time - can be null)
	Maximum size of an item is 400KB
DynamoDB - Basics	Data types supported are:
	Scalar Types - String, Number, Binary, Boolean, Null
	Document Types - List, Map Set Types - String Set Number Set Ripary Set
	Set Types - String Set, Number Set, Binary Set
	· Therefore, in DynamoDB you can rapidly evolve schemas
	Therefore, in Dynamodd ydd cantapiaty eydtye schemas

	Duradicione ad Manda (dafa dh)
DynamoDB - Read/Write Capacity Modes	Provisioned Mode (default)
	You specify the number of reads/writes per second
	· You need to plan capacity beforehand
	 Pay for provisioned Read Capacity Units (RCU) & Write Capacity Units (WCU)
	Possibility to add auto-scaling mode for RCU & WCU
	· On-Demand Mode
	 Read/writes automatically scale up/down with your workloads
	No capacity planning needed
	Pay for what you use, more expensive (\$\$\$)
	Great for unpredictable workloads, steep sudden spikes
	Fully-managed, highly available, seamless inmemory
	cache for DynamoDB
	Help solve read congestion by caching
DynamoDB Accelerator (DAX)	Microseconds latency for cached data
	Doesn't require application logic modification
	(compatible with existing DynamoDB APIs)
	• 5 minutes TTL for cache (default)
	Amazon ElastiCache: Store Aggregation Result
DynamoDB Accelerator (DAX) vs.	DynamoDB Accelerator (DAX): - Individual objects cache
ElastiCache	- Query & Scan cache
	Ordered stream of item-level modifications (create/update/delete) in a table
	· Use cases:
	React to changes in real-time (welcome email to users)
DynamoDB - Stream Processing	• Real-time usage analytics
	Insert into derivative tables
	Implement cross-region replication
	 Invoke AWS Lambda on changes to your DynamoDB table
	DynamoDB Streams
	· 24 hours retention
	• Limited # of consumers
	Process using AWS Lambda Triggers, or
	DynamoDB Stream Kinesis adapter
DynamoDB Streams vs Kinesis Data Streams	Kinesis Data Streams (newer)
(newer)	· 1 year retention
	· High # of consumers
	Process using AWS Lambda, Kinesis Data
	Analytics, Kineis Data Firehose, AWS Glue
	Streaming ETL
	Make a DynamoDB table accessible with low latency in multiple-regions
	Active-Active replication
DynamoDB Global Tables	Applications can READ and WRITE to the table in any region
	Must enable DynamoDB Streams as a pre-requisite
	Automatically delete items after an expiry
Dynama DB Time To Live (TTL)	timestamp
DynamoDB - Time To Live (TTL)	Use cases: reduce stored data by keeping only
	current items, adhere to regulatory
	obligations, web session handling

DynamoDB - Backups for disaster recovery	Continuous backups using point-in-time recovery (PITR) Optionally enabled for the last 35 days Point-in-time recovery to any time within the backup window The recovery process creates a new table On-demand backups Full backups for long-term retention, until explicitly deleted Doesn't affect performance or latency Can be configured and managed in AWS Backup (enables cross-region copy) The recovery process creates a new table
DynamoDB - Integration with Amazon S3	Export to S3 (must enable PITR) · Works for any point of time in the last 35 days · Doesn't affect the read capacity of your table · Perform data analysis on top of DynamoDB · Retain snapshots for auditing · ETL on top of S3 data before importing back into DynamoDB · Export in DynamoDB JSON or ION format · Import from S3 · Import CSV, DynamoDB JSON or ION format · Doesn't consume any write capacity · Creates a new table · Import errors are logged in CloudWatch Logs
AWS API Gateway	AWS Lambda + API Gateway: No infrastructure to manage Support for the WebSocket Protocol Handle API versioning (v1, v2) Handle different environments (dev, test, prod) Handle security (Authentication and Authorization) Create API keys, handle request throttling Swagger / Open API import to quickly define APIs Transform and validate requests and responses Generate SDK and API specifications Cache API responses
API Gateway - Integrations High Level	Lambda Function Invoke Lambda function Easy way to expose REST API backed by AWS Lambda HTTP Expose HTTP endpoints in the backend Example: internal HTTP API on premise, Application Load Balancer Why? Add rate limiting, caching, user authentications, API keys, etc AWS Service Expose any AWS API through the API Gateway Example: start an AWS Step Function workflow, post a message to SQS Why? Add authentication, deploy publicly, rate control
API Gateway - Endpoint Types	Edge-Optimized (default): For global clients Requests are routed through the CloudFront Edge locations (improves latency) The API Gateway still lives in only one region Regional: For clients within the same region Could manually combine with CloudFront (more control over the caching strategies and the distribution) Private: Can only be accessed from your VPC using an interface VPC endpoint (ENI) Use a resource policy to define access

API Gateway - Security	User Authentication through IAM Roles (useful for internal applications) Cognito (identity for external users - example mobile users) Custom Authorizer (your own logic) Custom Domain Name HTTPS security through integration with AWS Certificate Manager (ACM) If using Edge-Optimized endpoint, then the certificate must be in us-east-1 If using Regional endpoint, the certificate must be in the API Gateway region Must setup CNAME or A-alias record in Route 53
AWS Step Functions	Build serverless visual workflow to orchestrate your Lambda functions Features: sequence, parallel, conditions, timeouts, error handling, Can integrate with EC2, ECS, On-premises servers, API Gateway, SQS queues, etc Possibility of implementing human approval feature Use cases: order fulfillment, data processing, web applications, any workflow
Amazon Cognito	Give users an identity to interact with our web or mobile application Cognito User Pools: Sign in functionality for app users Integrate with API Gateway & Application Load Balancer Cognito Identity Pools (Federated Identity): Provide AWS credentials to users so they can access AWS resources directly Integrate with Cognito User Pools as an identity provider Cognito vs IAM: "hundreds of users", "mobile users", "authenticate with SAML"
Cognito User Pools (CUP) - User Features	Create a serverless database of user for your web & mobile apps · Simple login: Username (or email) / password combination · Password reset · Email & Phone Number Verification · Multi-factor authentication (MFA) · Federated Identities: users from Facebook, Google, SAML
Cognito User Pools (CUP) - Integrations	· CUP integrates with API Gateway and Application Load Balancer
Cognito Identity Pools (Federated Identities)	Get identities for "users" so they obtain temporary AWS credentials Users source can be Cognito User Pools, 3rd party logins, etc Users can then access AWS services directly or through API Gateway The IAM policies applied to the credentials are defined in Cognito They can be customized based on the user_id for fine grained control Default IAM roles for authenticated and guest users
AWS Hosted Websites	 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

Micro Services architecture	We want to switch to a micro service architecture Many services interact with each other directly using a REST API
	Each architecture for each micro service may vary in form and shape
	· We want a micro-service architecture so we can have a leaner
	development lifecycle for each service
	You are free to design each micro-service the way you want
	Synchronous patterns: API Gateway, Load Balancers
	· Asynchronous patterns: SQS, Kinesis, SNS, Lambda triggers (S3)
	Challenges with micro-services:
	· repeated overhead for creating each new microservice,
	· issues with optimizing server density/utilization
Discussions on Micro Services	· complexity of running multiple versions of multiple microservices simultaneously
	· proliferation of client-side code requirements to integrate with many separate
	services.
	Some of the challenges are solved by Serverless patterns:
	· API Gateway, Lambda scale automatically and you pay per usage
	· You can easily clone API, reproduce environments
	Generated client SDK through Swagger integration for the API Gateway
	No changes to architecture
	· Will cache software update files at the edge
	· Software update files are not dynamic, they're static (never changing)
AMb of Classes Fire and O	• Our EC2 instances aren't serverless
Why CloudFront?	But CloudFront is, and will scale for us
	· Our ASG will not scale as much, and we'll save tremendously in EC2
	· We'll also save in availability, network bandwidth cost, etc
	· Easy way to make an existing application more scalable and cheaper!
Database Types - RDBMS (= SQL / OLTP)	RDS, Aurora - great for joins
Database Types - NoSQL database - no	DynamoDB (~JSON), ElastiCache (key /
joins, no SQL	value pairs), Neptune (graphs), DocumentDB (for MongoDB), Keyspaces (for Apache
OITS, TO SQL	Cassandra)
Database Types - Object Store	S3 (for big objects) / Glacier (for backups / archives)
Database Types - Data Warehouse	(= SQL Analytics / BI): Redshift (OLAP), Athena, EMR
Database Types - Search	OpenSearch (JSON) - free text, unstructured searches
Database Types - Graphs	Amazon Neptune - displays relationships between data
Database Types - Ledger	Amazon Quantum Ledger Database
Database Types - Time series	Amazon Timestream
Amazon Aurora - Summary: Aurora	for unpredictable / intermittent workloads, no capacity planning
Serverless	
Amazon Aurora - Summary: Aurora Multi-	for continuous writes failover (high write availability)
Master	
Amazon Aurora - Summary: Aurora Global	up to 16 DB Read Instances in each region, < 1 second storage replication
Amazon Aurora - Summary: Aurora Machine	perform ML using SageMaker & Comprehend on Aurora
Amazon Aurora - Summary: Aurora Machine Learning	perform ML using SageMaker & Comprehend on Aurora
	perform ML using SageMaker & Comprehend on Aurora new cluster from existing one, faster than restoring a snapshot

	Managed Redis / Memcached (similar offering as RDS, but for caches)
	In-memory data store, sub-millisecond latency
	 Select an ElastiCache instance type (e.g., cache.móg.large)
	 Support for Clustering (Redis) and Multi AZ, Read Replicas (sharding)
	· Security through IAM, Security Groups, KMS, Redis Auth
Amazon ElastiCache	Backup / Snapshot / Point in time restore feature
	Managed and Scheduled maintenance
	Requires some application code changes to be leveraged
	Use Case: Key/Value store, Frequent reads, less writes, cache results for DB
	queries, store session data for websites, cannot use SQL.
	queries, store session data for websites, carnot use use.
	AWS proprietary technology, managed serverless NoSQL database, millisecond
	latency
	· Capacity modes: provisioned capacity with optional auto-scaling or on-demand
	capacity
	· Can replace ElastiCache as a key/value store (storing session data for example, using
	TTL feature)
	· Highly Available, Multi AZ by default, Read and Writes are decoupled, transaction
	capability
	DAX cluster for read cache, microsecond read latency
	Security, authentication and authorization is done through IAM
Amoran Dunama DR. Cumamani	
Amazon DynamoDB - Summary	Event Processing: DynamoDB Streams to integrate with AWS Lambda, or Kinesis Data
	Streams
	Global Table feature: active-active setup
	 Automated backups up to 35 days with PITR (restore to new table), or on-demand
	backups
	• Export to S3 without using RCU within the PITR window, import from S3 without using
	WCU
	· Great to rapidly evolve schemas
	· Use Case: Serverless applications development (small documents 100s KB), distributed
	serverless
	cache
Amazon S3 - Summary	S3 is a key / value store for objects
	Great for bigger objects, not so great for many small objects
	Serverless, scales infinitely, max object size is 5 TB, versioning capability
	• Tiers: S3 Standard, S3 Infrequent Access, S3 Intelligent, S3 Glacier + lifecycle policy
	• Features: Versioning, Encryption, Replication, MFA-Delete, Access Logs
	· Security: IAM, Bucket Policies, ACL, Access Points, Object Lambda, CORS,
	Object/Vault Lock
	• Encryption: SSE-S3, SSE-KMS, SSE-C, client-side, TLS in transit, default encryption
	Batch operations on objects using S3 Batch, listing files using S3 Inventory
	Performance: Multi-part upload, S3 Transfer Acceleration, S3 Select
	Automation: S3 Event Notifications (SNS, SQS, Lambda, EventBridge)
	Use Cases: static files, key value store for big files, website hosting
	- 030 Cases, static files, key value store for big files, website flosting
	 DocumentDB is the same for MongoDB (which is a NoSQL database)
	 MongoDB is used to store, query, and index JSON data
D =	· Similar "deployment concepts" as Aurora
DocumentDB	• Fully Managed, highly available with replication across 3 AZ
	DocumentDB storage automatically grows in increments of 10GB, up to 64 TB.
	Automatically scales to workloads with millions of requests per seconds
	. Elemandary deaths to membrada minimization of respectito per seconds

	Fully managed graph database
	· A popular graph dataset would be a social network
	Users have friends
	Posts have comments
	Comments have likes from users
	Users share and like posts
Amazon Neptune	Highly available across 3 AZ, with up to 15 read replicas
Amazon reptone	Build and run applications working with highly connected
	datasets - optimized for these complex and hard queries
	· Can store up to billions of relations and query the graph with
	milliseconds latency
	 Highly available with replications across multiple AZs
	· Great for knowledge graphs (Wikipedia), fraud detection,
	recommendation engines, social networking
	A managed Apache Cassandra-compatible database service
	· Serverless, Scalable, highly available, fully managed by AWS
	· Automatically scale tables up/down based on the application's traffic
	• Tables are replicated 3 times across multiple AZ
Amazon Keyspaces (for Apache Cassandra)	Using the Cassandra Query Language (CQL)
	Single-digit millisecond latency at any scale, 1000s of requests per second
	Capacity: On-demand mode or provisioned mode with auto-scaling
	Encryption, backup, Point-In-Time Recovery (PITR) up to 35 days
	· QLDB stands for "Quantum Ledger Database"
	• A ledger is a book recording financial transactions
	Fully Managed, Serverless, High available, Replication across 3 AZ
	Used to review history of all the changes made to your application data over time
Amazon QLDB	• Immutable system: no entry can be removed or modified, cryptographically verifiable
	Difference with Amazon Managed Blockchain: no decentralization component, in
	accordance with
	financial regulation rules
	Fully managed, fast, scalable, serverless time series database
	Automatically scales up/down to adjust capacity
	Store and analyze trillions of events per day
	· 1000s times faster & 1/10th the cost of relational databases
	Scheduled queries, multi-measure records, SQL compatibility
Amazon Timestream	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
	Serverless guery service to analyze data stored in Amazon S3
	Uses standard SQL language to query the files (built on Presto)
	Supports CSV, JSON, ORC, Avro, and Parquet
	Pricing: \$5.00 per TB of data scanned
Amazon Athena	Commonly used with Amazon Quicksight for
ATTIGEOTI AUTOTIA	
	reporting/dashboards
	Use cases: Business intelligence / analytics / reporting, analyze & Guarty VDC Flow Logs FLB Logs Cloud Trail trails etc. The case of the case
	query VPC Flow Logs, ELB Logs, CloudTrail trails, etc
	• Exam Tip: analyze data in S3 using serverless SQL, use Athena

	Use columnar data for cost-savings (less scan)
	Apache Parquet or ORC is recommended
	Huge performance improvement
	Use Glue to convert your data to Parquet or ORC
	· Compress data for smaller retrievals (bzip2, gzip, lz4, snappy, zlip, zstd)
	Partition datasets in S3 for easy querying on virtual columns
Amazon Athena - Performance Improvement	· s3://yourBucket/pathToTable
	/ <partition_column_name>=<value></value></partition_column_name>
	/ <partition_column_name>=<value></value></partition_column_name>
	/ <partition_column_name>=<value></value></partition_column_name>
	/etc
	Example: s3://athena-examples/flight/parquet/year=1991/month=1/day=1/
	 Use larger files (> 128 MB) to minimize overhead
	Allows you to run SQL queries across
	data stored in relational, non-relational,
	object, and custom data sources (AWS
	or on-premises)
Amazon Athena - Federated Query	· 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
	Dedebiff is based on DestaraÇOL but itle not used for OLTD
	Redshift is based on PostgreSQL, but it's not used for OLTP
	It's OLAP - online analytical processing (analytics and data
	warehousing)
De debiff Overview	• 10x better performance than other data warehouses, scale to PBs of data
Redshift Overview	Columnar storage of data (instead of row based) & parallel query engine Payanayay as based on the instances provisioned.
	Pay as you go based on the instances provisioned
	Has a SQL interface for performing the queries Place is such as Amazon Quicksight or Tableau integrate with it.
	 BI tools such as Amazon Quicksight or Tableau integrate with it vs Athena: faster queries / joins / aggregations thanks to indexes
	- vs Attiena. Taster queries / Johns / aggregations thanks to indexes
	· Leader node: for query
	planning, results
	aggregation
	Compute node: for
	performing the queries,
Redshift Cluster	send results to leader
	· You provision the node
	size in advance
	· You can used Reserved
	Instances for cost
	savings
	• Redshift has "Multi-AZ" mode for some
	clusters
	Snapshots are point-in-time backups of a cluster,
	stored internally in S3
D 11:11 C 1 : 0.55	Snapshots are incremental (only what has
Redshift - Snapshots & DR	Snapshots are incremental (only what has changed is saved)
Redshift - Snapshots & DR	
Redshift - Snapshots & DR	changed is saved)
Redshift - Snapshots & DR	changed is saved) You can restore a snapshot into a new cluster

	• Query data that is already in
	S3 without loading it
	Must have a Redshift cluster
Redshift Spectrum	available to start the query
	· The query is then submitted
	to thousands of Redshift
	Spectrum nodes
	Amazon OpenSearch is successor to Amazon ElasticSearch
	· In DynamoDB, queries only exist by primary key or indexes
	· With OpenSearch, you can search any field, even partially matches
	· It's common to use OpenSearch as a complement to another database
Amazon OpenSearch Service	• Two modes: managed cluster or serverless cluster
	· Does not natively support SQL (can be enabled via a plugin)
	Ingestion from Kinesis Data Firehose, AWS IoT, and CloudWatch Logs
	Security through Cognito & IAM, KMS encryption, TLS
	· Comes with OpenSearch Dashboards (visualization)
	• EMR stands for "Elastic MapReduce"
	• EMR helps creating Hadoop clusters (Big Data) to analyze and process vast
	amount of data
	• The clusters can be made of hundreds of EC2 instances
Amazon EMR	• EMR comes bundled with Apache Spark, HBase, Presto, Flink
	EMR takes care of all the provisioning and configuration
	Auto-scaling and integrated with Spot instances
	Use cases: data processing, machine learning, web indexing, big data
	Master Node: Manage the cluster, coordinate, manage health - long running
	· Core Node: Run tasks and store data - long running
	· Task Node (optional): Just to run tasks - usually Spot
Amazon EMR - Node types & purchasing	Purchasing options:
,, ,	· On-demand: reliable, predictable, won't be terminated
	 Reserved (min 1 year): cost savings (EMR will automatically use if available)
	Spot Instances: cheaper, can be terminated, less reliable
	· Can have long-running cluster, or transient (temporary) cluster
	· Serverless machine learning-powered business intelligence service to create
	interactive dashboards
	· Fast, automatically scalable, embeddable, with per-session pricing
	• Use cases:
	Business analytics
	Building visualizations
Amazon QuickSight	Perform ad-hoc analysis
a_on doichoight	Get business insights using data
	· Integrated with RDS, Aurora,
	Athena, Redshift, S3
	· In-memory computation using SPICE
	engine if data is imported into QuickSight
	Enterprise edition: Possibility to setup
	Column-Level security (CLS)

QuickSight - Dashboard & Analysis	 Define Users (standard versions) and Groups (enterprise version) These users & groups only exist within QuickSight, not IAM!! A dashboard is a read-only snapshot of an analysis that you can share preserves the configuration of the analysis (filtering, parameters, controls, sort) You can share the analysis or the dashboard with Users or Groups To share a dashboard, you must first publish it Users who see the dashboard can also see the underlying data
AWS Glue	 Managed extract, transform, and load (ETL) service Useful to prepare and transform data for analytics Fully serverless service
Glue - Glue Job Bookmarks	prevent re-processing old data
Glue - 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 - Glue DataBrew	clean and normalize data using pre-built transformation
Glue - Glue Studio	new GUI to create, run and monitor ETL jobs in Glue
Glue - Glue Streaming ETL	(built on Apache Spark Structured Streaming): compatible with Kinesis Data Streaming, Kafka, MSK (managed Kafka)
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 your Data Lake It automates many complex manual steps (collecting, cleansing, moving, cataloging data,) and de-duplicate (using ML Transforms) Combine structured and unstructured data in the data lake Out-of-the-box source blueprints: S3, RDS, Relational & NoSQL DB Fine-grained Access Control for your applications (row and column-level) Built on top of AWS Glue
Kinesis Data Analytics (SQL application)	 Real-time analytics on Kinesis Data Streams & Firehose using SQL Add reference data from Amazon S3 to enrich streaming data Fully managed, no servers to provision Automatic scaling Pay for actual consumption rate Output: Kinesis Data Streams: create streams out of the real-time analytics queries Kinesis Data Firehose: send analytics query results to destinations Use cases: Time-series analytics Real-time dashboards Real-time metrics
Kinesis Data Analytics for Apache Flink	 Use Flink (Java, Scala or SQL) to process and analyze streaming data Run any Apache Flink application on a managed cluster on AWS provisioning compute resources, parallel computation, automatic scaling application backups (implemented as checkpoints and snapshots) Use any Apache Flink programming features Flink does not read from Firehose (use Kinesis Analytics for SQL instead)

	Alternative to Amazon Kinesis
	• Fully managed Apache Kafka on AWS
	• Allow you to create, update, delete clusters
	MSK creates & manages Kafka brokers nodes & Zookeeper nodes for you
Amazon Managed Streaming for Apache	Deploy the MSK cluster in your VPC, multi-AZ (up to 3 for HA)
Kafka (Amazon MSK)	Automatic recovery from common Apache Kafka failures
	Data is stored on EBS volumes for as long as you want
	• MSK Serverless
	Run Apache Kafka on MSK without managing the capacity
	MSK automatically provisions resources and scales compute & storage
	Kinesis Data Streams
	· 1 MB message size limit
	• Data Streams with Shards
	Shard Splitting & Merging
	• TLS In-flight encryption
	· KMS at-rest encryption
Kinesis Data Streams vs. Amazon MSK	Amazon MSK
	· 1MB default, configure for higher (ex: 10MB)
	Kafka Topics with Partitions
	· Can only add partitions to a topic
	PLAINTEXT or TLS In-flight Encryption
	· KMS at-rest encryption
	Amazon
	- Kinesis Data Analytics for Apache Flink
	- AWS Glue - Streaming ETL Jobs - Powered by Apache Spark Streaming
Amazon MSK Consumers	- Lambda
	- Applications running on Amazon EC2, ECS, EKS
	• IoT Core allows you to harvest data from IoT devices
	Kinesis is great for real-time data collection
	Firehose helps with data delivery to S3 in near real-time (1 minute)
	Lambda can help Firehose with data transformations
Big Data Ingestion Pipeline	Amazon S3 can trigger notifications to SQS
3 3	Lambda can subscribe to SQS (we could have connecter S3 to Lambda)
	Athena is a serverless SQL service and results are stored in S3
	The reporting bucket contains analyzed data and can be used by reporting tool such
	as AWS QuickSight, Redshift, etc
	Find objects, people, text, scenes in images and videos using ML
	Facial analysis and facial search to do user verification, people counting
	racial analysis and facial search to do user verification, people coording Create a database of "familiar faces" or compare against celebrities
	Use cases:
Amazon Bakagnitisa	· Labeling
Amazon Rekognition	Content Moderation Text Detaction
	• Text Detection
	• Face Detection and Analysis (gender, age range, emotions)
	• Face Search and Verification
	· Celebrity Recognition
	Pathing (ex: for sports game analysis)

	Detect content that is inappropriate, unwanted,
	or offensive (image and videos)
	Used in social media, broadcast media,
	advertising, and e-commerce situations to create
Amazon Rekognition - Content Moderation	a safer user experience
	Set a Minimum Confidence Threshold for
	items that will be flagged
	Flag sensitive content for manual review in
	Amazon Augmented AI (A2I)
	Help comply with regulations
	· Automatically convert speech to text
	· Uses a deep learning process called automatic speech recognition (ASR) to
	convert speech to text quickly and accurately
	Automatically remove Personally Identifiable Information (PII) using Redaction
Amazon Transcribe	Supports Automatic Language Identification for multi-lingual audio
Amazon manscribe	
	· Use cases:
	transcribe customer service calls
	· automate closed captioning and subtitling
	generate metadata for media assets to create a fully searchable archive
Amazon Polly	• Turn text into lifelike speech using deep learning
Amazon Folly	Allowing you to create applications that talk
	Customize the pronunciation of words with Pronunciation lexicons
	• Stylized words: St3ph4ne => "Stephane"
	• Acronyms: AWS => "Amazon Web Services"
	Upload the lexicons and use them in the SynthesizeSpeech operation
	Generate speech from plain text or from documents marked up with Speech
Amazon Polly - Lexicon & SSML	Synthesis Markup Language (SSML) - enables more customization
	• emphasizing specific words or phrases
	- using phonetic pronunciation
	· including breathing sounds, whispering
	using the Newscaster speaking style
	Natural and accurate language translation
Amazon Translate	· Amazon Translate allows you to localize content - such as websites and applications -
	for international users, and to easily translate large volumes of text efficiently.
	· Amazon Lex: (same technology that powers Alexa)
	· Automatic Speech Recognition (ASR) to convert speech to text
	· Natural Language Understanding to recognize the intent of text, callers
	Helps build chatbots, call center bots
Amazon Lex & Connect	· Amazon Connect:
	Receive calls, create contact flows, cloud-based virtual contact center
	· Can integrate with other CRM systems or AWS
	No upfront payments, 80% cheaper than traditional contact center solutions
	 For Natural Language Processing - NLP Fully managed and serverless service
	Uses machine learning to find insights and relationships in text
Amazon Comprehend	Language of the text
	• Extracts key phrases, places, people, brands, or events
	Understands how positive or negative the text is
	Analyzes text using tokenization and parts of speech
	· Automatically organizes a collection of text files by topic

Amazon Comprehend Medical	 Amazon Comprehend Medical detects and returns useful information in unstructured clinical text: Physician's notes Discharge summaries Test results Case notes Uses NLP to detect Protected Health Information (PHI) - DetectPHI API Store your documents in Amazon 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 for developers / data scientists to build ML models Typically, difficult to do all the processes in one place + provision servers Machine learning process (simplified): predicting your exam score
Amazon Forecast	 Fully managed service that uses ML to deliver highly accurate forecasts Example: predict the future sales of a raincoat 50% more accurate than looking at the data itself Reduce forecasting time from months to hours Use cases: Product Demand Planning, Financial Planning, Resource Planning,
Amazon Kendra	 Fully managed document search service powered by Machine Learning Extract answers from within a document (text, pdf, HTML, PowerPoint, MS Word, FAQs) Natural language search capabilities Learn from user interactions/feedback to promote preferred results (Incremental Learning) Ability to manually fine-tune search results (importance of data, freshness, custom,)
Amazon Personalize	Fully managed ML-service to build apps with real-time personalized recommendations Example: personalized product recommendations/re-ranking, customized direct marketing Example: User bought gardening tools, provide recommendations on the next one to buy Same technology used by Amazon.com Integrates into existing websites, applications, SMS, email marketing systems, Implement in days, not months (you don't need to build, train, and deploy ML solutions) Use cases: retail stores, media and entertainment
Amazon Textract	 Automatically extracts text, handwriting, and data from any scanned documents using AI and ML Extract data from forms and tables Read and process any type of document (PDFs, images,) Use cases: Financial Services (e.g., invoices, financial reports) Healthcare (e.g., medical records, insurance claims) Public Sector (e.g., tax forms, ID documents, passports)
Amazon CloudWatch Metrics	 CloudWatch provides metrics for every services in AWS Metric is a variable to monitor (CPUUtilization, NetworkIn) Metrics belong to namespaces Dimension is an attribute of a metric (instance id, environment, etc) Up to 30 dimensions per metric Metrics have timestamps Can create CloudWatch dashboards of metrics Can create CloudWatch Custom Metrics (for the RAM for example)

	· Continually stream CloudWatch
	metrics to a destination of your choice,
	with near-real-time delivery and low
	latency.
	Amazon Kinesis Data Firehose (and then
CloudWatch Metric Streams	its destinations)
	· 3rd party service provider: Datadog,
	Dynatrace, New Relic, Splunk, Sumo
	Logic
	• Option to filter metrics to only stream
	a subset of them
	Log groups: arbitrary name, usually representing an application
	· Log stream: instances within application / log files / containers
	· Can define log expiration policies (never expire, 1 day to 10 years)
	· CloudWatch Logs can send logs to:
	· Amazon S3 (exports)
CloudWatch Logs	· Kinesis Data Streams
	· Kinesis Data Firehose
	· AWS Lambda
	• OpenSearch
	· Logs are encrypted by default
	· Can setup KMS-based encryption with your own keys
	SDK, CloudWatch Logs Agent, CloudWatch Unified Agent
	Elastic Beanstalk: collection of logs from application
	• ECS: collection from containers
CloudWatch Logs - Sources	AWS Lambda: collection from function logs
	· VPC Flow Logs: VPC specific logs
	· API Gateway
	· CloudTrail based on filter
	• Route53: Log DNS queries
	Search and analyze log data stored in CloudWatch Logs
	• Example: find a specific IP inside a log, count occurrences of
	"ERROR" in your logs
	Provides a purpose-built query language
CloudWatch Logs Insights	 Automatically discovers fields from AWS services and JSON log events
Stood water 2090 monghto	Fetch desired event fields, filter based on conditions, calculate
	aggregate statistics, sort events, limit number of events
	Can save queries and add them to CloudWatch Dashboards
	· Can query multiple Log Groups in different AWS accounts
	· It's a query engine, not a real-time engine
	· Log data can take up to 12 hours to
	become available for export
CloudWatch Logs - S3 Export	• The API call is CreateExportTask
	Not near-real time or real-time use
	Logs Subscriptions instead
	Get a real-time log events from CloudWatch Logs for processing and analysis
	· Send to Kinesis Data Streams, Kinesis Data Firehose, or Lambda
CloudWatch Logs Subscriptions	· Subscription Filter - filter which logs are events delivered to your destination
	· Cross-Account Subscription - send log events to resources in a different AWS account
	(KDS, KDF)

	By default, no logs from your EC2
	machine will go to CloudWatch
	· You need to run a CloudWatch
	agent on EC2 to push the log files
CloudWatch Logs for EC2	you want
	· Make sure IAM permissions are
	correct
	· The CloudWatch log agent can be
	setup on-premises too
CloudWatch Logs Agent	· Old version of the agent
etoda vateli zogo vigelit	· Can only send to CloudWatch Logs
	· Collect additional system-level metrics such as RAM, processes, etc
CloudWatch Unified Agent	· Collect logs to send to CloudWatch Logs
	Centralized configuration using SSM Parameter Store
	· CPU (active, guest, idle, system, user, steal)
	 Disk metrics (free, used, total), Disk IO (writes, reads, bytes, iops)
	· RAM (free, inactive, used, total, cached)
CloudWatch Unified Agent - Metrics	 Netstat (number of TCP and UDP connections, net packets, bytes)
	· Processes (total, dead, bloqued, idle, running, sleep)
	· Swap Space (free, used, used %)
	· Reminder: out-of-the box metrics for EC2 - disk, CPU, network (high level)
	· Alarms are used to trigger notifications for any metric
	· Various options (sampling, %, max, min, etc)
	· Alarm States:
	·OK
CloudWatch Alarms	· INSUFFICIENT_DATA
	·ALARM
	· Period:
	· Length of time in seconds to evaluate the metric
	· High resolution custom metrics: 10 sec, 30 sec or multiples of 60 sec
	· Stop, Terminate, Reboot, or Recover an EC2 Instance
CloudWatch Alarm Targets	· Trigger Auto Scaling Action
	· Send notification to SNS (from which you can do pretty much anything)
	CloudWatch Alarms are on a single metric
CloudWatch Alarms - Composite Alarms	· Composite Alarms are monitoring the states of multiple other alarms
Ctobawatch Atams - Composite Atams	· AND and OR conditions
	· Helpful to reduce "alarm noise" by creating complex composite alarms
	- Status Check:
	· Instance status = check the EC2 VM
EC2 Instance Recovery	System status = check the underlying hardware
	· Recovery: Same Private, Public, Elastic IP, metadata, placement group
	Recovery: Same Private, Public, Elastic IP, metadata, placement group Alarms can be created based on CloudWatch Logs Metrics Filters
CloudWatch Alarm: good to know	· Alarms can be created based on CloudWatch Logs Metrics Filters
CloudWatch Alarm: good to know	Alarms can be created based on CloudWatch Logs Metrics Filters To test alarms and notifications, set the alarm state to Alarm using CLI
	Alarms can be created based on CloudWatch Logs Metrics Filters To test alarms and notifications, set the alarm state to Alarm using CLI aws cloudwatch set-alarm-statealarm-name "myalarm"state-value ALARMstate-
CloudWatch Alarm: good to know Amazon EventBridge (formerly CloudWatch Events)	Alarms can be created based on CloudWatch Logs Metrics Filters To test alarms and notifications, set the alarm state to Alarm using CLI aws cloudwatch set-alarm-statealarm-name "myalarm"state-value ALARMstate-reason "testing purposes"

Amazon EventBridge - Event Buses	 Event buses can be accessed by other AWS accounts using Resource-based Policies You can archive events (all/filter) sent to an event bus (indefinitely or set period) Ability to replay archived events
Amazon EventBridge - Schema Registry	EventBridge can analyze the events in your bus and infer the schema The Schema Registry allows you to generate code for your application, that will know in advance how data is structured in the event bus Schema can be versioned
CloudWatch Container Insights	 Collect, aggregate, summarize metrics and logs from containers Available for containers on Amazon Elastic Container Service (Amazon ECS) Amazon Elastic Kubernetes Services (Amazon EKS) Kubernetes platforms on EC2 Fargate (both for ECS and EKS)
CloudWatch Container Insights - EKS and Kubernetes	In Amazon EKS and Kubernetes, CloudWatch Insights is using a containerized version of the CloudWatch Agent to discover containers
CloudWatch Lambda Insights	Monitoring and troubleshooting solution for serverless applications running on AWS 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 is provided as a Lambda Layer
CloudWatch Contributor Insights	 Analyze log data and create time series that display contributor data. See metrics about the top-N contributors The total number of unique contributors, and their usage. This helps you find top talkers and understand who or what is impacting system performance. Works for any AWS-generated logs (VPC, DNS, etc) For example, you can find bad hosts, identify the heaviest network users, or find the URLs that generate the most errors. You can build your rules from scratch, or you can also 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	 Provides automated dashboards that show potential problems with monitored applications, to help isolate ongoing issues Your applications run on Amazon EC2 Instances with select technologies only (Java, .NET, Microsoft IIS Web Server, databases) And you can use other AWS resources such as Amazon EBS, RDS, ELB, ASG, Lambda, SQS, DynamoDB, S3 bucket, ECS, EKS, SNS, API Gateway 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 Amazon EventBridge and SSM OpsCenter
CloudWatch Container Insights - Insights and Operational Visibility	ECS, EKS, Kubernetes on EC2, Fargate, needs agent for Kubernetes Metrics and logs
CloudWatch Lambda Insights - Insights and Operational Visibility	Detailed metrics to troubleshoot serverless applications
CloudWatch Contributors Insights	Find "Top-N" Contributors through CloudWatch Logs
CloudWatch Application Insights - Insights and Operational Visibility	Automatic dashboard to troubleshoot your application and related AWS services
AWS CloudTrail	 Provides governance, compliance and audit for your AWS Account CloudTrail is enabled by default! Get an history of events / API calls made within your AWS Account by: Console SDK CLI AWS Services Can put logs from CloudTrail into CloudWatch Logs or S3 A trail can be applied to All Regions (default) or a single Region. If a resource is deleted in AWS, investigate CloudTrail first!
CloudTrail Events	Management Events: Operations that are performed on resources in your AWS account Examples: Configuring security (IAM AttachRolePolicy) Configuring rules for routing data (Amazon EC2 CreateSubnet) Setting up logging (AWS CloudTrail CreateTrail) By default, trails are configured to log management events. Can separate Read Events (that don't modify resources) from Write Events (that may modify resources) Data Events: By default, data events are not logged (because high volume operations) Amazon S3 object-level activity (ex: GetObject, DeleteObject, PutObject): can separate Read and Write Events AWS Lambda function execution activity (the Invoke API)
CloudTrail Insights	 Enable CloudTrail Insights to detect unusual activity in your account: inaccurate resource provisioning hitting service limits Bursts of AWS IAM actions Gaps in periodic maintenance activity CloudTrail Insights analyzes normal management events to create a baseline And then continuously analyzes write events to detect unusual patterns Anomalies appear in the CloudTrail console Event is sent to Amazon S3 An EventBridge event is generated (for automation needs)

CloudTrail Events Retention	• Events are stored for 90 days in CloudTrail
	· To keep events beyond this period, log them to S3 and use Athena
AWS Config	· Helps with auditing and recording compliance of your AWS resources
	Helps record configurations and changes over time
	· Questions that can be solved by AWS Config:
	· Is there unrestricted SSH access to my security groups?
	• Do my buckets have any public access?
	· How has my ALB configuration changed over time?
	· You can receive alerts (SNS notifications) for any changes
	· AWS Config is a per-region service
	· Can be aggregated across regions and accounts
	Possibility of storing the configuration data into S3 (analyzed by Athena)
	· Can use AWS managed config rules (over 75)
	· Can make custom config rules (must be defined in AWS Lambda)
	• Ex: evaluate if each EBS disk is of type gp2
	• Ex: evaluate if each EC2 instance is t2.micro
Config Rules	· Rules can be evaluated / triggered:
	• For each config change
	· And / or: at regular time intervals
	AWS Config Rules does not prevent actions from happening (no deny)
	Pricing: no free tier, \$0.003 per configuration item recorded per region,
	\$0.001 per config rule evaluation per region
	View compliance of a resource over time
AWS Config Resource	· View configuration of a resource over time
	· View CloudTrail API calls of a resource over time
	Automate remediation of non-compliant resources using SSM Automation
	Documents
Config Rules - Remediations	· Use AWS-Managed Automation Documents or create custom Automation Documents
	• Tip: you can create custom Automation Documents that invokes Lambda function
	· You can set Remediation Retries if the resource is still non-compliant after
	autoremediation
	· Use EventBridge to trigger notifications when AWS resources are noncompliant
Config Rules - Notifications	· Ability to send configuration changes and compliance state notifications to SNS (all
	events - use SNS Filtering or filter at client-side)
	· CloudWatch
	· Performance monitoring (metrics, CPU, network, etc) & dashboards
	• Events & Alerting
	· Log Aggregation & Analysis
	• CloudTrail
CloudWatch vs CloudTrail vs Config	Record API calls made within your Account by everyone
Cloud watch vs Cloud I rail vs Config	· Can define trails for specific resources
	· Global Service
	·Config
	· Record configuration changes
	Evaluate resources against compliance rules
	· Get timeline of changes and compliance

	· CloudWatch:
ELB - CloudWatch vs CloudTrail vs Config	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
	• Global service
	· Allows to manage multiple AWS accounts
	• The main account is the management account
	Other accounts are member accounts
AWS Organizations	Member accounts can only be part of one organization
	Consolidated Billing across all accounts - single payment method
	Pricing benefits from aggregated usage (volume discount for EC2, S3)
	Shared reserved instances and Savings Plans discounts across accounts
	API is available to automate AWS account creation
	• Advantages
	Multi Account vs One Account Multi VPC
	Use tagging standards for billing purposes
	Enable CloudTrail on all accounts, send logs to central S3 account
AWS Organizations - Advantages and	Send CloudWatch Logs to central logging account
Security	Establish Cross Account Roles for Admin purposes
	Security: Service Control Policies (SCP)
	· IAM policies applied to OU or Accounts to restrict Users and Roles
	· They do not apply to the management account (full admin power)
	 Must have an explicit allow (does not allow anything by default - like IAM)
	aws:Sourcelp
	restrict the client IP from
	which the API calls are being made
	aws:RequestedRegion
	restrict the region the
	API calls are made to
	ec2:ResourceTag
	restrict based on tags
IAM Conditions	
	aws:MultiFactorAuthPresent
	to force MFA
	a7.LishDuglich payaissian applica to
	• s3:ListBucket permission applies to
	arn:aws:s3:::test
	• => bucket level permission
	• s3:GetObject, s3:PutObject,
	s3:DeleteObject applies to
	arn:awn:s3:::test/*
	· => object level permission
Resource Policies & aws:PrincipalOrgID	aws:PrincipalOrgID can be used in any resource policies to restrict access to accounts

IAM Roles vs Resource Based Policies	 Cross account: attaching a resource-based policy to a resource (example: S3 bucket policy) OR using a role as a proxy When you assume a role (user, application or service), you give up your original permissions and take the permissions assigned to the role When using a resource-based policy, the principal doesn't have to give up his 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. Supported by: Amazon S3 buckets, SNS topics, SQS queues, etc
Amazon EventBridge - Security	When a rule runs, it needs permissions on the target Resource-based policy: Lambda, SNS, SQS, CloudWatch Logs, API Gateway IAM role: Kinesis stream, Systems Manager Run Command, ECS task
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. Can be used in combinations of AWS Organizations SCP 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)
AWS IAM Identity Center (successor to AWS Single Sign-On)	 One login (single sign-on) for all your AWS accounts in AWS Organizations Business cloud applications (e.g., Salesforce, Box, Microsoft 365,) SAML2.0-enabled applications EC2 Windows Instances Identity providers Built-in identity store in IAM Identity Center 3rd party: Active Directory (AD), OneLogin, Okta

	. Multi-Account Parmissions
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
What is Microsoft Active Directory (AD)?	Found on any Windows Server with AD Domain Services Database of objects: User Accounts, Computers, Printers, File Shares, Security Groups Centralized security management, create account, assign permissions Objects are organized in trees A group of trees is a forest
AWS Managed Microsoft AD	 Create your own AD in AWS, manage users locally, supports MFA Establish "trust" connections with your onpremises
AD Connector	 Directory Gateway (proxy) to redirect to onpremises AD, supports MFA Users are managed on the on-premises AD
Simple AD	AD-compatible managed directory on AWS Cannot be joined with on-premises AD
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
AWS Control Tower	Easy way to set up and govern a secure and compliant multi-account AWS environment based on best practices AWS Control Tower uses AWS Organizations to create accounts Benefits: Automate the set up of your environment in a few clicks Automate ongoing policy management using guardrails Detect policy violations and remediate them Monitor compliance through an interactive dashboard
AWS Control Tower - Guardrails	Provides ongoing governance for your Control Tower environment (AWS Accounts) • Preventive Guardrail - using SCPs (e.g., Restrict Regions across all your accounts) • Detective Guardrail - using AWS Config (e.g., identify untagged resources)

Encryption in flight (SSL)	 Data is encrypted before sending and decrypted after receiving SSL certificates help with encryption (HTTPS) Encryption in flight ensures no MITM (man in the middle attack) can happen
Server side encryption at rest	 Data is encrypted after being received by the server Data is decrypted before being sent It is stored in an encrypted form thanks to a key (usually a data key) The encryption / decryption keys must be managed somewhere and the server must have access to it
Client side encryption	 Data is encrypted by the client and never decrypted by the server Data will be decrypted by a receiving client The server should not be able to decrypt the data Could leverage Envelope Encryption
AWS KMS (Key Management Service)	 Anytime you hear "encryption" for an AWS service, it's most likely KMS AWS manages encryption keys for us Fully integrated with IAM for authorization Easy way to control access to your data Able to audit KMS Key usage using CloudTrail Seamlessly integrated into most AWS services (EBS, S3, RDS, SSM) Never ever store your secrets in plaintext, especially in your code! KMS Key Encryption also available through API calls (SDK, CLI) Encrypted secrets can be stored in the code / environment variables
KMS Keys Types	 KMS Keys is the new name of KMS Customer Master Key Symmetric (AES-256 keys) Single encryption key that is used to Encrypt and Decrypt AWS services that are integrated with KMS use Symmetric CMKs You never get access to the KMS Key unencrypted (must call KMS API to use) Asymmetric (RSA & ECC key pairs) Public (Encrypt) and Private Key (Decrypt) pair Used for Encrypt/Decrypt, or Sign/Verify operations The public key is downloadable, but you can't access the Private Key unencrypted Use case: encryption outside of AWS by users who can't call the KMS API
Types of KMS Keys	AWS Owned Keys (free): SSE-S3, SSE-SQS, SSE-DDB (default key) • AWS Managed Key: free (aws/service-name, example: aws/rds or aws/ebs) • Customer managed keys created in KMS: \$1 / month • Customer managed keys imported (must be symmetric key): \$1 / month • pay for API call to KMS (\$0.03 / 10000 calls)
Automatic Key rotation	AWS-managed KMS Key: automatic every 1 year Customer-managed KMS Key: (must be enabled) automatic every 1 year Imported KMS Key: only manual rotation possible using alias
KMS Key Policies	 Control access to KMS keys, "similar" to S3 bucket policies Difference: you cannot control access without them Default KMS Key Policy: Created if you don't provide a specific KMS Key Policy Complete access to the key to the root user = entire AWS account Custom KMS Key Policy: Define users, roles that can access the KMS key Define who can administer the key Useful for cross-account access of your KMS key

Copying Snapshots across accounts	1. Create a Snapshot, encrypted with
	your own KMS Key (Customer
	Managed Key)
	2. Attach a KMS Key Policy to
	authorize cross-account access
	3. Share the encrypted snapshot
	4. (in target) Create a copy of the
	Snapshot, encrypt it with a CMK in
	your account
	5. Create a volume from the snapshot
	· Identical KMS keys in different AWS Regions that can be used interchangeably
	• Multi-Region keys have the same key ID, key material, automatic rotation
	Encrypt in one Region and decrypt in other Regions
	No need to re-encrypt or making cross-Region API calls
KMS Multi-Region Keys	10 10 14 11 D 1
	KMS Multi-Region are NOT global (Primary + Replicas)
	Each Multi-Region key is managed independently
	Use cases: global client-side encryption, encryption on Global DynamoDB, Global
	Aurora
	We can encrypt specific attributes client-side
	in our DynamoDB table using the Amazon
	DynamoDB Encryption Client
	Combined with Global Tables, the client-side
	encrypted data is replicated to other regions
DynamoDB Global Tables and KMS Multi-	If we use a multi-region key, replicated in the
Region Keys Client-Side encryption	same region as the DynamoDB Global table,
Region Reys eacht side eneryption	then clients in these regions can use lowlatency
	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
	We can encrypt specific attributes client-side
	in our Aurora table using the AWS
	Encryption SDK
	Combined with Aurora Global Tables, the
	client-side encrypted data is replicated to
	other regions
	If we use a multi-region key, replicated in the
Global Aurora and KMS Multi-Region Keys	same region as the Global Aurora DB, then
Client-Side encryption	clients in these regions can use low-latency
State State Street years	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, we can
	protect specific fields even from database
	admins

	 Unencrypted objects and objects encrypted with SSE-S3 are replicated by default
	· Objects encrypted with SSE-C (customer provided key) are never replicated
	· For objects encrypted with SSE-KMS, you need to enable the option
	· Specify which KMS Key to encrypt the objects within the target bucket
	· Adapt the KMS Key Policy for the target key
	· An IAM Role with kms:Decrypt for the source KMS Key and kms:Encrypt for the target
S3 Replication Encryption Considerations	KMS Key
	You might get KMS throttling errors, in which case you can ask for a Service Quotas
	increase
	· You can use multi-region AWS KMS Keys, but they are currently treated as
	independent keys by Amazon S3 (the object will still be decrypted and then
	encrypted)
	AMI in Source Account is encrypted with KMS Key
	from Source Account
	2. Must modify the image attribute to add a Launch
	Permission which corresponds to the specified target
	AWS account
	3. Must share the KMS Keys used to encrypted the
	snapshot the AMI references with the target account
AMI Sharing Process Encrypted via KMS	/ IAM Role
	4. The IAM Role/User in the target account must have
	the permissions to DescribeKey, ReEncrypted,
	CreateGrant, Decrypt
	5. When launching an EC2 instance from the AMI,
	optionally the target account can specify a new KMS
	key in its own account to re-encrypt the volumes
	Secure storage for configuration and secrets
	Optional Seamless Encryption using KMS
	· Serverless, scalable, durable, easy SDK
SSM Parameter Store	Version tracking of configurations / secrets
	· Security through IAM
	Notifications with Amazon EventBridge
	· Integration with CloudFormation
	GetParameters or
	GetParametersByPath API - in Lambda
	·/my-department/
	· my-app/
	· dev/
	· db-url
SSM Parameter Store Hierarchy	· db-password
	· prod/
	· db-url
	· db-password
	· other-app/
	· /other-department/
	· aws/reference/secretsmanager/secret_ID_in_Secrets_Manager
	·/aws/service/ami-amazon-linux-latest/amzn2-ami-hvm-x86_64-gp2 (public)
	2 3. 4

	Standard
	Total number of parameters allowed:
	- 10,000
	Maximum size of a parameter value:
	- 4 KB
	Parameter policies available:
	- No
	Cost:
	- No additional charge
	Storage Pricing:
	- Free
Standard and advanced parameter tiers	
otandara ana advancea parameter tiers	Advanced
	Total number of parameters allowed:
	- 100,000
	Maximum size of a parameter value:
	- 8 KB
	Parameter policies available:
	- Yes
	Cost:
	- Charges apply
	Storage Pricing:
	- 0.05\$ per advanced parameter per month
	· Allow to assign a TTL to a parameter (expiration date) to force
Parameters Policies (for advanced	updating or deleting sensitive data such as passwords
parameters)	Can assign multiple policies at a time
	· Can assign motuple policies at a time
	 Newer service, meant for storing secrets
	Capability to force rotation of secrets every X days
	· Automate generation of secrets on rotation (uses Lambda)
AWS Secrets Manager	Integration with Amazon RDS (MySQL, PostgreSQL, Aurora)
	• Secrets are encrypted using KMS
	Mostly meant for RDS integration
	Replicate Secrets across multiple AWS Regions
AWS Secrets Manager - Multi-Region	· Secrets Manager keeps read replicas in sync with the primary Secret
Secrets	Ability to promote a read replica Secret to a standalone Secret
	Use cases: multi-region apps, disaster recovery strategies, multi-region DB
	- Ose cases: motti region apps, disaster recovery strategies, motti region bb
	• Easily provision, manage, and deploy TLS Certificates
	· Provide in-flight encryption for websites (HTTPS)
	Supports both public and private TLS certificates
	• Free of charge for public TLS certificates
	Automatic TLS certificate renewal
AWS Certificate Manager (ACM)	
	• Integrations with (load TLS certificates on)
	Elastic Load Balancers (CLB, ALB, NLB)
	CloudFront Distributions
	· APIs on API Gateway
	· Cannot use ACM with EC2 (can't be extracted)

	1. List domain names to be included in the certificate
	Fully Qualified Domain Name (FQDN): corp.example.com
	Wildcard Domain: *.example.com
	2. Select Validation Method: DNS Validation or Email validation
ACM Paguasting Public Cartificator	· DNS Validation is preferred for automation purposes
ACM - Requesting Public Certificates	· Email validation will send emails to contact addresses in the WHOIS database
	· DNS Validation will leverage a CNAME record to DNS config (ex: Route 53)
	3. It will take a few hours to get verified
	4. The Public Certificate will be enrolled for automatic renewal
	· ACM automatically renews ACM-generated certificates 60 days before expiry
	· Option to generate the certificate
	outside of ACM and then import it
	No automatic renewal, must import a
	new certificate before expiry
	· ACM sends daily expiration events
	starting 45 days prior to expiration
ACM - Importing Public Certificates	• The # of days can be configured
	Events are appearing in EventBridge
	AWS Config has a managed rule
	named acm-certificate-expiration-check
	to check for expiring certificates
	(configurable number of days)
	(comigurable number of days)
	For global clients
Edge-Optimized (default) API Gateway	 Requests are routed through the CloudFront Edge locations (improves latency)
	· The API Gateway still lives in only one region
	· For clients within the same region
Regional API Gateway - Endpoint Types	· Could manually combine with CloudFront (more control over the caching strategies
	and the distribution)
Drivete ADI Cetevrey Finding int Types	· Can only be accessed from your VPC using an interface VPC endpoint (ENI)
Private API Gateway - Endpoint Types	· Use a resource policy to define access
	· Create a Custom Domain Name in API Gateway
	Edge-Optimized (default): For global clients
	Requests are routed through the CloudFront Edge locations
	(improves latency)
	The API Gateway still lives in only one region
	• The TLS Certificate must be in the same region as
ACM - Integration with API Gateway	CloudFront, in us-east-1
	• Then setup CNAME or (better) A-Alias record in Route 53
	Regional:
	• For clients within the same region
	The TLS Certificate must be imported on API Gateway, in
	the same region as the API Stage
	_
	• Then setup CNAME or (better) A-Alias record in Route 53

	Protects your web applications from common web exploits (Layer 7)
	· Layer 7 is HTTP (vs Layer 4 is TCP/UDP)
	· Deploy on
	Application Load Balancer
	· API Gateway
	· CloudFront
	AppSync GraphQL API
	· Cognito User Pool
AWS WAF - Web Application Firewall	Define Web ACL (Web Access Control List) Rules:
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1	• IP Set: up to 10,000 IP addresses - use multiple Rules for more IPs
	HTTP headers, HTTP body, or URI strings Protects from common attack - SQL
	injection and Cross-Site Scripting (XSS)
	Size constraints, geo-match (block countries)
	· Rate-based rules (to count occurrences of events) - for DDoS protection
	· Web ACL are Regional except for CloudFront
	· A rule group is a reusable set of rules that you can add to a web ACL
	WAF does not support the Network Load Balancer (Layer 4)
	• We can use Global Accelerator for fixed IP and WAF on the ALB
	· DDoS: Distributed Denial of Service - many requests at the same time
	· AWS Shield Standard:
	• Free service that is activated for every AWS customer
	Provides protection from attacks such as SYN/UDP Floods, Reflection attacks and
	other
	layer 3/layer 4 attacks
	· AWS Shield Advanced:
AWS Shield: protect from DDoS attack	• Optional DDoS mitigation service (\$3,000 per month per organization)
	Protect against more sophisticated attack on Amazon EC2, Elastic Load Balancing
	(ELB), Amazon CloudFront, AWS Global Accelerator, and Route 53
	• 24/7 access to AWS DDoS response team (DRP)
	Protect against higher fees during usage spikes due to DDoS
	Shield Advanced automatic application layer DDoS mitigation automatically creates,
	evaluates and deploys AWS WAF rules to mitigate layer 7 attacks
	<u> </u>
	Manage rules in all accounts of an AWS Organization
	• Security policy: common set of security rules
	• WAF rules (Application Load Balancer, API Gateways, CloudFront)
	AWS Shield Advanced (ALB, CLB, NLB, Elastic IP, CloudFront)
AWS Firewall Manager	Security Groups for EC2, Application Load BAlancer and ENI resources in VPC
	· AWS Network Firewall (VPC Level)
	· Amazon Route 53 Resolver DNS Firewall
	Policies are created at the region level
	· Rules are applied to new resources as they are created (good for compliance) across
	all and future accounts in your Organization
	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
WAF vs. Firewall Manager vs. Shield	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
	,

	· BP1 - CloudFront
	· Web Application delivery at
	the edge
	Protect from DDoS Common
	Attacks (SYN floods, UDP
	reflection)
	· BP1 - Global Accelerator
	· Access your application from
AWS Best Practices for DDoS Resiliency	the edge
Edge Location Mitigation (BP1, BP3)	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
	la fara-hara-hara-la rana-
	· Infrastructure layer defense
	(BP1, BP3, BP6)
	Protect Amazon EC2 against high
	traffic
	• That includes using Global
	Accelerator, Route 53,
	CloudFront, Elastic Load Balancing
AWS Best Practices for DDoS Resiliency Best	• Amazon EC2 with Auto Scaling
pratices for DDoS mitigation	(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
	Elastic Load Balancing scales with the traffic increases and will
	Elastic Load Balancing scales with

	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
AWS Best Practices for DDoS Resiliency	automatically block the IPs of bad
Application Layer Defense	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
	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
AWS Best Practices for DDoS Resiliency	filter traffic based on specific IP at the
Attack surface reduction	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
	Intelligent Threat discovery to protect your AWS Account
	· Uses Machine Learning algorithms, anomaly detection, 3rd party data
	• One click to enable (30 days trial), no need to install software
	Input data includes:
	· CloudTrail Events Logs - unusual API calls, unauthorized deployments
	· CloudTrail Management Events - create VPC subnet, create trail,
Amazon GuardDuty	· CloudTrail S3 Data Events - get object, list objects, delete object,
	VPC Flow Logs - unusual internal traffic, unusual IP address
	• DNS Logs - compromised EC2 instances sending encoded data within DNS queries
	Optional Features - EKS Audit Logs, RDS & Aurora, EBS, Lambda, S3 Data Events
	Can setup EventBridge rules to be notified in case of findings
	EventBridge rules can target AWS Lambda or SNS
	Can protect against CryptoCurrency attacks (has a dedicated "finding" for it)

	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
Amazon Inspector	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
	Remember: only for EC2 instances, Container Images & Lambda functions
	· Continuous scanning of the infrastructure, only when needed
What does Amazon Inspector evaluate?	Package vulnerabilities (EC2, ECR & Lambda) - database of CVE
	Network reachability (EC2)
	· A risk score is associated with all vulnerabilities for prioritization
	Amazon Macie is a fully managed data security and data privacy service that uses
	machine learning and pattern matching to discover and protect your sensitive data in
AWS Macie	AWS.
	Macie helps identify and alert you to sensitive data, such as personally identifiable
	information (PII)
	· Classless Inter-Domain Routing - a method for allocating IP addresses
	Used in Security Groups rules and AWS networking in general
	• They help to define an IP address range:
	• We've seen WW.XX.YY.ZZ/32 => one IP
	• We've seen 0.0.0.0/0 => all IPs
	But we can define:192.168.0.0/26 =>192.168.0.0 - 192.168.0.63 (64 IP addresses)
	• A CIDR consists of two components
	Base IP
	Represents an IP contained in the range (XX.XX.XX.XX)
Understanding CIDR - IPv4	• Example: 10.0.0.0, 192.168.0.0,
	· Subnet Mask
	Defines how many bits can change in the IP
	• Example: /0, /24, /32
	· Can take two forms:
	· /8 6 255.0.0.0
	·/16 ó 255.255.0.0
	·/24 ó 255.255.255.0
	·/32 ó 255.255.255

	The Subnet Mask basically allows part of the underlying IP to get additional next values from the base IP
Understanding CIDR - Subnet Mask	(SubnetQ23): What would be the maximum number of hosts that can be supported on each subnet if you are given the following IP address: 152.77.0.0/19
	8190
	You calculate the number of hosts based on the remaining trailing bits. Since you must use 3 bits in the third octet for the subnets, 5 trailing bits remain. In addition to these 5 bits, you have all 8 bits remaining in the fourth octet. Hence you have a total of 13 trailing bits that can be used for host addresses. In this case that equates to $2^13 = 8,192 - 2 = 8,190$ hosts available to each of the 6 subnets.
	(SubnetQ10): You are given an address of 10.10.10.8/8. You need to create 3,000 subnets. What will be your new subnet mask?
	255.255.240.0
	255.255.240.0 or /20 Ask yourself the question of what to the power of 2 gives you a value of 3,000. The answer is $2^12 = 4,096 - 2 = 4,096$ subnets. $2^11 = 2,048 - 2 = 2,046$ subnets which is quite enough. This is a Class A network, so you borrow 12 bits from the last three octets. 8 bits from the second octet and 4 bits from the third octet. Hence your mask would then be 255.255.240.0. The remaining 12 bits (4 from the third octet and 8 from the fourth octet) will form your host addresses.
	(SubnetQ26): How many subnets can an IP address of 39.0.0.0/15 support?
	126
	126 Since you are using 7 bits, the number of supported subnets is 2^7 = 128 - 2 = 126.
	(SubnetQ27): What is the maximum number of hosts that you can put on each subnet if you had an IP address of 39.0.0.0/15?
	131,070
	If you borrowed the 7 leading bits for the subnet mask from the second octet, then your would have 1 trailing bit remaining from the second octet and 8 trailing bits remaining from each of the third and fourth octets which you can use to configure host IDs. That would be $2^17 - 2 = 131,070$
192.168.0.0/24 = ?	192.168.0.0 - 192.168.0.255 (256 IPs)When in doubt, use this website https://www.ipaddressguide.com/cidr
192.168.0.0/16 = ?	192.168.0.0 - 192.168.255.255 (65,536 IPs)
134.56.78.123/32 = ?	Just 134.56.78.123
Dublic vs. Drivete ID (ID://)	The Internet Assigned Numbers Authority (IANA) established certain blocks of IPv4 addresses for the use of private (LAN) and public (Internet) addresses Private IP can only allow certain values:
Public vs. Private IP (IPv4)	 10.0.0.0 - 10.255.255.255 (10.0.0.0/8) ç in big networks 172.16.0.0 - 172.31.255.255 (172.16.0.0/12) çAWS default VPC in that range 192.168.0.0 - 192.168.255.255 (192.168.0.0/16) ç e.g., home networks All the rest of the IP addresses on the Internet are Public

172.16.0.0/12	172.16.0.0 - 172.31.255.255
192.168.0.0/16	192.168.0.0 - 192.168.255.255
VPC = Virtual Private Cloud You can have multiple VPCs in an AWS region (max. 5 per region - soft limit) Max. CIDR per VPC is 5, for each CIDR: Min. size is /28 (16 IP addresses) Max. size is /16 (65536 IP addresses)	Your VPC CIDR should NOT overlap with your other networks (e.g., corporate)
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)
Internet Gateway (IGW)	 Allows resources (e.g., EC2 instances) in a VPC connect to the Internet It scales horizontally and is highly available and redundant Must be created separately from a VPC One VPC can only be attached to one IGW and vice versa Internet Gateways on their own do not allow Internet access Route tables must also be edited!
Bastion Hosts	We can use a Bastion Host to SSH into our private EC2 instances The bastion is in the public subnet which is then connected to all other private subnets Bastion Host security group must allow inbound from the internet on port 22 from restricted CIDR, for example the public CIDR of your corporation Security Group of the EC2 Instances must allow the Security Group of the Bastion Host, or the private IP of the Bastion host

	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
NAT Instance (outdated, but still at the exam)	Pre-configured Amazon Linux AMI is available
	• Reached the end of standard support on December 31, 2020
	Not highly available / resilient setup out of the box
	· You need to create an ASG in multi-AZ + resilient user-data script
	Internet traffic bandwidth depends on EC2 instance type
	You must manage Security Groups & rules:
	· Inbound:
	· Allow HTTP / HTTPS traffic coming from Private Subnets
	Allow SSH from your home network (access is provided through Internet Gateway)
	• Outbound:
	• Allow HTTP / HTTPS traffic to the Internet
	Allow Tittly Tittle didnie to the internet
	AWS-managed NAT, higher bandwidth, high availability, no administration
	Pay per hour for usage and bandwidth
	 NATGW is created in a specific Availability Zone, uses an Elastic IP
NAT Gateway	· Can't be used by EC2 instance in the same subnet (only from other subnets)
	 Requires an IGW (Private Subnet => NATGW => IGW)
	· 5 Gbps of bandwidth with automatic scaling up to 100 Gbps
	No Security Groups to manage / required
	NAT Gateway is resilient within
	a single Availability Zone
	Must create multiple NAT
	Gateways in multiple AZs for
NAT Gateway with High Availability	fault-tolerance
	• There is no cross-AZ failover
	needed because if an AZ goes
	Š
	down it doesn't need NAT
	10.34.23.0
	This is a class A network because of the value of 10 in the first octet. But the subnet
	mask is the key to this question. The 1's or the 255's in the subnet mask dictate the
(CulomatOO()) What is theturnely	network address. The default mask for a Class A network is 255.0.0.0. However, in this
(SubnetQ04): What is the network number	question, the /24 means that 24 bits are being applied to the subnet mask. In this case
for this host: 10.34.23.5/24?	the subnet mask would be 255.255.255.0. The 255's are telling you what part of the
	given IP address is network, and what part is host. Since the mask numbers are equal
	and fall right at the octet breaks, the network is 10.34.23.0. The host is 0.0.0.5. With the
	/24 you have actually created 2^16-2 = 65,534 subnets on which you can put 254 hosts
	on each of the subnets.

NAT Gateway vs NAT Instance	NAT Gateway: - Highly available within AZ (create in another AZ) - Up to 100 Gbps Bandwidth - Managed by AWS - Cost is per hour and amount of data transferred - Has Public and Private IPv4 NAT Instance: - Use a script to manage failover between instances - Depends on EC2 instance type - Managed by you (e.g., software, OS patches,) - Per hour, EC2 instance type and size, + network \$ - Has Security Groups and can be used as a Bastion Host
Network Access Control List (NACL)	 NACL are like a firewall which control traffic from and to subnets One NACL per subnet, new subnets are assigned the Default NACL You define NACL Rules: Rules have a number (1-32766), higher precedence with a lower number First rule match will drive the decision Example: if you define #100 ALLOW 10.0.010/32 and #200 DENY 10.0.010/32, the IP address will be allowed because 100 has a higher precedence over 200 The last rule is an asterisk (*) and denies a request in case of no rule match AWS recommends adding rules by increment of 100 Newly created NACLs will deny everything NACL are a great way of blocking a specific IP address at the subnet level
Default NACL	 Accepts everything inbound/outbound with the subnets it's associated with Do NOT modify the Default NACL, instead create custom NACLs
Ephemeral Ports	 For any two endpoints to establish a connection, they must use ports Clients connect to a defined port, and expect a response on an ephemeral port Different Operating Systems use different port ranges, examples: IANA & MS Windows 10 è 49152 - 65535 Many Linux Kernels è 32768 - 60999
Security Group vs NACLs	Security Group Operates at the instance level Supports allow rules only Stateful: return traffic is automatically allowed, regardless of any rules All rules are evaluated before deciding whether to allow traffic Applies to an EC2 instance when specified by someone NACL Operates at the subnet level Supports allow rules and deny rules Stateless: return traffic must be explicitly allowed by rules (think of ephemeral ports) Rules are evaluated in order (lowest to highest) when deciding whether to allow traffic, first match wins Automatically applies to all EC2 instances in the subnet that it's associated with

VPC Peering	Privately connect two VPCs using AWS' network
	Make them behave as if they were in the
	same network
	Must not have overlapping CIDRs
	VPC Peering connection is NOT transitive
	(must be established for each VPC that
	need to communicate with one another)
	You must update route tables in each
	VPC's subnets to ensure EC2 instances
	can communicate with each other
	Can communicate with each other
	· You can create VPC Peering connection between VPCs in different AWS
VPC Peering - Good to know	accounts/regions
J	 You can reference a security group in a peered VPC (works cross accounts - same
	region)
	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
\(\text{PO.5.}\)	using the public Internet
VPC Endpoints (AWS PrivateLink)	• 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
	Provisions an ENI (private IP address) as an entry
Interface Endpoints (powered by	point (must attach a Security Group)
PrivateLink)	Supports most AWS services
	• \$ per hour + \$ per GB of data processed
	Provisions a gateway and must be used as a
	target in a route table (does not use security
Gateway Endpoints	groups)
	· Supports both S3 and DynamoDB
	· Free
	Gateway is most likely going to be preferred all the time at the exam
	· Cost: free for Gateway, \$ for interface endpoint
Gateway or Interface Endpoint for S3?	• Interface Endpoint is preferred access is required from onpremises (Site to Site VPN
	or Direct Connect), a different VPC or a different region
	DynamoDB is a public service
	from AWS
	Option 1: Access from the public
	internet
	Because Lambda is in a VPC, it
	needs a NAT Gateway in a public
Lambda in VPC accessing DynamoDB	subnet and an internet gateway
	Option 2 (better & free): Access
	from the private VPC network
	Deploy a VPC Gateway endpoint
	for DynamoDB
	· Change the Route Tables
	- Gridinge the Noote rapies

	· Capture information about IP traffic going into your interfaces:
VPC Flow Logs	· VPC Flow Logs
	· Subnet Flow Logs
	• Elastic Network Interface (ENI) Flow Logs
	Helps to monitor & troubleshoot connectivity issues
	• Flow logs data can go to S3, CloudWatch Logs, and Kinesis Data Firehose
	• Captures network information from AWS managed interfaces too: ELB, RDS,
	ElastiCache, Redshift, WorkSpaces, NATGW, Transit Gateway
	· srcaddr & dstaddr - help identify problematic IP
	· srcport & dstport - help identity problematic ports
	· Action - success or failure of the request due to Security Group / NACL
VPC Flow Logs Syntax	· 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-
	logsrecords-examples.html
	Look at the "ACTION" field
	Incoming Requests
	· Inbound REJECT => NACL or SG
VPC Flow Logs - Troubleshoot SG & NACL	· Inbound ACCEPT, Outbound REJECT =>
issues	NACL
	· Outbound REJECT => NACL or SG
	· Outbound ACCEPT, Inbound REJECT =>
	NACL
	· Virtual Private Gateway (VGW)
	· VPN concentrator on the AWS side of the VPN connection
	· VGW is created and attached to the VPC from which you want to create the
1146 611 1 611 14611	Site-to-Site VPN connection
AWS Site-to-Site VPN	Possibility to customize the ASN (Autonomous System Number)
	· Customer Gateway (CGW)
	· Software application or physical device on customer side of the VPN connection
	https://docs.aws.amazon.com/vpn/latest/s2svpn/your-cgw.html#DevicesTested
	Customer Gateway Device (On-premises)
	· What IP address to use?
	Public Internet-routable IP address for your Customer
	Gateway device
	· If it's behind a NAT device that's enabled for NAT
Site-to-Site VPN Connections	traversal (NAT-T), use the public IP address of the NAT
	device
	· Important step: enable Route Propagation for
Site-to-site VEIN Connections	· Important step. enable Route Propagation for
Site-to-Site YFIN Connections	the Virtual Private Gateway in the route table
Site-to-site yfin Coillections	
Site-to-Site YFIN Connections	the Virtual Private Gateway in the route table
Site-to-Site YFIN Connections	the Virtual Private Gateway in the route table that is associated with your subnets
Site-to-Site YFIN Connections	the Virtual Private Gateway in the route table that is associated with your subnets If you need to ping your EC2 instances from

	· Provide secure communication between
	multiple sites, if you have multiple VPN
	connections
	· Low-cost hub-and-spoke model for
	primary or secondary network connectivity
AWS VPN CloudHub	between different locations (VPN only)
	· It's a VPN connection so it goes over the
	public Internet
	• To set it up, connect multiple VPN
	connections on the same VGW, setup
	dynamic routing and configure route tables
	aynumb roomig and comiger or roots tubics
	 Provides a dedicated private connection from a remote network to your
	VPC
	 Dedicated connection must be setup between your DC and AWS Direct
	Connect locations
	· You need to setup a Virtual Private Gateway on your VPC
Direct Connect (DX)	· Access public resources (S3) and private (EC2) on same connection
	· Use Cases:
	· Increase bandwidth throughput - working with large data sets - lower cost
	· More consistent network experience - applications using real-time data feeds
	Hybrid Environments (on prem + cloud)
	· Supports both IPv4 and IPv6
	If you want to gatus a Direct Connect to one or more VDC in many different regions
Direct Connect Gateway	 If you want to setup a Direct Connect to one or more VPC in many different regions (same account), you must use a Direct Connect Gateway
	 Dedicated Connections: 1Gbps,10 Gbps and 100 Gbps capacity
	Physical ethernet port dedicated to a customer
	 Request made to AWS first, then completed by AWS Direct Connect Partners
Direct Connect - Connection Types	 Hosted Connections: 50Mbps, 500 Mbps, to 10 Gbps
2	Connection requests are made via AWS Direct Connect Partners
	Capacity can be added or removed on demand
	· 1, 2, 5, 10 Gbps available at select AWS Direct Connect Partners
	· Lead times are often longer than 1 month to establish a new connection
	Data in transit is not encrypted but is
	private
	· AWS Direct Connect + VPN
	provides an IPsec-encrypted private
Direct Connect - Encryption	connection
	Good for an extra level of security,
	but slightly more complex to put in
	place
	· .
	High Resiliency for Critical Workloads - One connection at multiple locations
Direct Connect - Resiliency	Maximum Resiliency for Critical Workloads - Maximum resilience is achieved by
	separate connections terminating on separate devices in more than one location.
Site-to-Site VPN connection as a backup	

Transit Gateway	· For having transitive peering between thousands of VPC
	and on-premises, hub-and-spoke (star) connection
	• Regional resource, can work cross-region
	Share cross-account using Resource Access Manager (RAM)
	You can peer Transit Gateways across regions
	Route Tables: limit which VPC can talk with other VPC
	Works with Direct Connect Gateway, VPN connections
	Supports IP Multicast (not supported by any other AWS
	service)
	• ECMP = Equal-cost multi-path
	routing
	Routing strategy to allow to
	forward a packet over multiple
Transit Gateway: Site-to-Site VPN ECMP	best path
	Use case: create multiple Siteto-
	Site VPN connections to
	increase the bandwidth of
	your connection to AWS
	• Allows you to capture and inspect network
	traffic in your VPC
	Route the traffic to security appliances that
	you manage
	· Capture the traffic
	• From (Source) - ENIs
VPC - Traffic Mirroring	• To (Targets) - an ENI or a Network Load
j 	Balancer
	· Capture all packets or capture the packets of
	your interest (optionally, truncate packets)
	Source and Target can be in the same VPC or
	different VPCs (VPC Peering)
	Use cases: content inspection, threat
	monitoring, troubleshooting,
	• IPv4 designed to provide 4.3 Billion addresses (they'll be exhausted soon)
	• IPv6 is the successor of IPv4
	\cdot IPv6 is designed to provide 3.4 × 10!" unique IP addresses
	• Every IPv6 address is public and Internet-routable (no private range)
	• Format è x.x.x.x.x.x.x (x is hexadecimal, range can be from 0000 to ffff)
What is IPv6?	• Examples:
Trinacio il vo.	· 2001:db8:3333:4444:5555:6666:7777:8888
	· 2001:db8:3333:4444:cccc:dddd:eeee:ffff
	∙ :: è all 8 segments are zero
	• 2001:db8:: è the last 6 segments are zero
	∙ ::1234:5678 è the first 6 segments are zero
	• 2001:db8::1234:5678 è the middle 4 segments are zero
	• IPv4 cannot be disabled for your VPC and
	subnets
	· You can enable IPv6 (they're public IP addresses)
IPv6 in VPC	to operate in dual-stack mode
	Your EC2 instances will get at least a private
	internal IPv4 and a public IPv6
	•They can communicate using either IPv4 or IPv6
	to the internet through an Internet Gateway

IPv6 Troubleshooting	IPv4 cannot be disabled for your VPC and subnets So, if you cannot launch an EC2 instance in your subnet It's not because it cannot acquire an IPv6 (the space is very large) It's because there are no available IPv4 in your subnet Solution: create a new IPv4 CIDR in your subnet
Egress-only Internet Gateway	Used for IPv6 only • (similar to a NAT Gateway but for IPv6) • Allows instances in your VPC outbound connections over IPv6 while preventing the internet to initiate an IPv6 connection to your instances • You must update the Route Tables
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 - VPC Section	managed by AWS, provides scalable Internet access to private EC2 instances, IPv4 only
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)
Minimizing egress traffic network cost	Egress traffic: outbound traffic (from AWS to outside) Ingress traffic: inbound traffic - from outside to AWS (typically free) Try to keep as much internet traffic within AWS to minimize costs Direct Connect location that are co-located in the same AWS Region result in lower cost for egress network
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

Pricing: NAT Gateway vs Gateway VPC Endpoint	\$0.045 NAT Gateway / hour \$0.045 NAT Gateway data processed / GB \$0.09 Data transfer out to S3 (cross-region) \$0.00 Data transfer out to S3 (same-region) No cost for using Gateway Endpoint. \$0.01 Data transfer in/out (sameregion)
AWS Network Firewall	 Protect your entire Amazon VPC From Layer 3 to Layer 7 protection Any direction, you can inspect VPC to VPC traffic Outbound to internet Inbound from internet To / from Direct Connect & Site-to-Site VPN Internally, the AWS Network Firewall uses the AWS Gateway Load Balancer Rules can be centrally managed crossaccount by AWS Firewall Manager to apply to many VPCs
Network Firewall - Fine Grained Controls	 Supports 1000s of rules IP & port - example: 10,000s of IPs filtering Protocol - example: block the SMB protocol for outbound communications Stateful domain list rule groups: only allow outbound traffic to *.mycorp.com or third-party software repo General pattern matching using regex Traffic filtering: Allow, drop, or alert for the traffic that matches the rules Active flow inspection to protect against network threats with intrusionprevention capabilities (like Gateway Load Balancer, but all managed by AWS) Send logs of rule matches to Amazon S3, CloudWatch Logs, Kinesis Data Firehose
RPO vs RTO	RPO is the amount of data loss a business is willing to lose, measured in minutes to hours (whatever backups are); RTO is how much downtime (in hours) that a business can have and still survive.
Disaster Recovery - Pilot Light	 A small version of the app is always running in the cloud Useful for the critical core (pilot light) Very similar to Backup and Restore Faster than Backup and Restore as critical systems are already up
Warm Standby	Full system is up and running, but at minimum sizeUpon disaster, we can scale to production load
Multi Site / Hot Site Approach	Very low RTO (minutes or seconds) - very expensive • Full Production Scale is running AWS and On Premise
Backup - Disaster Recovery	 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
Backup - 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
Backup - Replication	 RDS Replication (Cross Region), AWS Aurora + Global Databases Database replication from on-premises to RDS Storage Gateway

· CloudFormation / Elastic Beanstalk to re-create a whole new environment	
Backup - Automation	Recover / Reboot EC2 instances with CloudWatch if alarms fail
	· AWS Lambda functions for customized automations
'	
	· Quickly and securely migrate databases to
	AWS, resilient, self healing
	• The source database remains available
	during the migration
	· Supports:
DMS - Database Migration Service	Homogeneous migrations: ex Oracle to
	Oracle
	Heterogeneous migrations: ex Microsoft SQL
	Server to Aurora
	Continuous Data Replication using CDC
	· You must create an EC2 instance to
	perform the replication tasks
	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 \$3
	DocumentDB
DMS Sources and Targets	
DMS Sources and Targets	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
	Convert your Database's Schema from one engine to another
	• Example OLTP: (SQL Server or Oracle) to MySQL, PostgreSQL, Aurora
	• Example OLAP: (Teradata or Oracle) to Amazon Redshift
AWS Schema Conversion Tool (SCT)	Prefer compute-intensive instances to optimize data conversions
	You do not need to use SCT if you are migrating the same DB engine
	• Ex: On-Premise PostgreSQL => RDS PostgreSQL
	• The DB engine is still PostgreSQL (RDS is the platform)
	When Multi-AZ Enabled, DMS provisions and maintains a
	synchronously stand replica in a different AZ
AWS DMS - Multi-AZ Deployment	· Advantages:
	Provides Data Redundancy
	- Eliminates I/O freezes
	Minimizes latency spikes

RDS MySQL to Aurora MySQL	Option 1: DB Snapshots from RDS MySQL restored as MySQL Aurora DB 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 \$) Use DMS if both databases are up and running
External MySQL to Aurora MySQL	Option 1: - Use Percona XtraBackup to create a file backup in Amazon S3 - Create an Aurora MySQL DB from Amazon S3 - Option 2: - Create an Aurora MySQL DB - Use the mysqldump utility to migrate MySQL into Aurora (slower than S3 method) Use DMS if both databases are up and running
RDS PostgreSQL to Aurora PostgreSQL	Option 1: DB Snapshots from RDS PostgreSQL restored as PostgreSQL Aurora DB Option 2: Create an Aurora Read Replica from your RDS PostgreSQL, and when the replication lag is 0, promote it as its own DB cluster (can take time and cost \$) Use DMS if both databases are up and running
External PostgreSQL to Aurora PostgreSQL	 Create a backup and put it in Amazon S3 Import it using the aws_s3 Aurora extension Use DMS if both databases are 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

	• Fully managed service
	Centrally manage and automate backups across AWS services
	No need to create custom scripts and manual processes
	- Supported services:
	- Amazon EC2 / Amazon EBS
	- Amazon S3
	- Amazon RDS (all DBs engines) / Amazon Aurora / Amazon DynamoDB
	- Amazon DocumentDB / Amazon Neptune
	- Amazon EFS / Amazon FSx (Lustre & Windows File Server)
AWS Backup	· AWS Storage Gateway (Volume Gateway)
·	Supports cross-region backups
	Supports cross-account backups
	Supports PITR for supported services
	· On-Demand and Scheduled backups
	· Tag-based backup policies
	· You create backup policies known as Backup Plans
	Backup frequency (every 12 hours, daily, weekly, monthly, cron expression)
	Backup window
	· Transition to Cold Storage (Never, Days, Weeks, Months, Years)
	· Retention Period (Always, Days, Weeks, Months, Years)
	Enforce a WORM (Write Once Read Many)
	state for all the backups that you store in
	your AWS Backup Vault
	Additional layer of defense to protect your
AWS Backup Vault Lock	backups against:
·	Inadvertent or malicious delete operations
	· Updates that shorten or alter retention periods
	• Even the root user cannot delete backups
	when enabled
	Plan migration projects by gathering information about on-premises data centers
	Server utilization data and dependency mapping are important for migrations
	Agentless Discovery (AWS Agentless Discovery Connector)
	VM inventory, configuration, and performance history such as CPU, memory, and disk
AWS Application Discovery Service	usage - Agent-based Discovery (AWS Application Discovery Agent)
	System configuration, system performance, running processes, and details of the
	network
	connections between systems
	•
	Resulting data can be viewed within AWS Migration Hub
	• The "AWS evolution" of CloudEndure Migration, replacing AWS Server Migration
	Service (SMS)
AWS Application Migration Service (MGN)	 Lift-and-shift (rehost) solution which simplify migrating applications to AWS
,	· Converts your physical, virtual, and cloud-based servers to run natively on AWS
	Supports wide range of platforms, Operating Systems, and databases
	Minimal downtime, reduced costs
	Some customers use VMware Cloud to manage their on-premises Data Center
	· They want to extend the Data Center capacity to AWS, but keep using the VMware
	Cloud software
	·Enter VMware Cloud on AWS
VMware Cloud on AWS	• Use cases
	Migrate your VMware vSphere-based workloads to AWS
	Run your production workloads across VMware vSphere-based private, public, and
	hybrid cloud environments
	Have a disaster recover strategy
	• nave a disaster recover strategy

Transferring large amount of data into AWS - Over the internet / Site-to-Site VPN	 Immediate to setup Will take 200(TB)1000(GB)1000(MB)*8(Mb)/100 Mbps = 16,000,000s = 185d
Transferring large amount of data into AWS -	- Long for the one-time setup (over a month)
Over direct connect 1Gbps	• Will take 200(TB) 1000(GB) 8(Gb)/1 Gbps = 1,600,000s = 18.5d
Transferring large energy of data into AMC	• Will take 2 to 3 snowballs in parallel
Transferring large amount of data into AWS -	• Takes about 1 week for the end-to-end transfer
Over Snowball	- Can be combined with DMS
Transferring large amount of data into AWS - For on-going replication / transfers	Site-to-Site VPN or DX with DMS or DataSync
	S3:ObjectCreated, S3:ObjectRemoved,
	S3:ObjectRestore, S3:Replication
	Object name filtering possible (*.jpg)
	Use case: generate thumbnails of images
S3 Event Notifications	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
	Advanced filtering options with JSON rules (metadata, object size, name)
S3 Event Notifications with Amazon	Multiple Destinations - ex Step Functions, Kinesis Streams / Firehose
EventBridge - duplicate	EventBridge Capabilities - Archive, Replay Events, Reliable delivery
	• The cloud is the perfect place to perform HPC
	You can create a very high number of resources in no time
	You can speed up time to results by adding more resources
High Performance Computing (HPC)	You can pay only for the systems you have used
	Perform genomics, computational chemistry, financial risk modeling,
	weather prediction, machine learning, deep learning, autonomous driving
	· Which services help perform HPC?
	· AWS Direct Connect:
	Move GB/s of data to the cloud, over a private secure network
	· Snowball & Snowmobile
Data Management & Transfer	• Move PB of data to the cloud
	· AWS DataSync
	· Move large amount of data between on-premises and S3, EFS, FSx for Windows
	Higher bandwidth, higher PPS (packet per second), lower latency
EC2 Enhanced Networking (SR-IOV)	• Option 1: Elastic Network Adapter (ENA) up to 100 Gbps
	• Option 2: Intel 82599 VF up to 10 Gbps - LEGACY
	· Improved ENA for HPC, only works for Linux
	Great for inter-node communications, tightly coupled workloads
Elastic Fabric Adapter (EFA)	· Leverages Message Passing Interface (MPI) standard
	Bypasses the underlying Linux OS to provide low-latency, reliable transport
	AWS Batch supports multi-node parallel jobs, which enables you to run single
AWS Batch	jobs that span multiple EC2 instances.
	Easily schedule jobs and launch EC2 instances accordingly
	Open-source cluster management tool to deploy HPC on AWS
	· Configure with text files
AWS ParallelCluster	-
AWS ParallelCluster	Automate creation of VPC, Subnet, cluster type and instance types

	CloudFormation is a declarative way of outlining your AWS
	Infrastructure, for any resources (most of them are supported).
	· For example, within a CloudFormation template, you say:
	· I want a security group
	• I want two EC2 instances using this security group
	· I want an S3 bucket
	· I want a load balancer (ELB) in front of these machines
	· Then CloudFormation creates those for you, in the right order, with the exact
	configuration that you specify
	• Infrastructure as code
	· No resources are manually created, which is excellent for control
	· Changes to the infrastructure are reviewed through code
	· Cost
	• Each resources within the stack is tagged with an identifier so you can easily see how
What is CloudFormation	much a stack costs you
	You can estimate the costs of your resources using the CloudFormation template
	Savings strategy: In Dev, you could automation deletion of templates at 5 PM and
	recreated at 8 AM, safely
	· Productivity
	Ability to destroy and re-create an infrastructure on the cloud on the fly
	Automated generation of Diagram for your templates! Deplace the generation of a great to find the second and a generation.
	Declarative programming (no need to figure out ordering and orchestration)
	Don't re-invent the wheel
	Leverage existing templates on the web!
	Leverage the documentation
	Supports (almost) all AWS resources:
	Everything we'll see in this course is supported
	· You can use "custom resources" for resources that are not supported
	• Example: WordPress CloudFormation Stack
CloudFormation Stack Designer	· We can see all the resources
	· We can see the relations between the components
	Fully managed service to send emails securely, globally and at scale
	· Allows inbound/outbound emails
	Reputation dashboard, performance insights, anti-spam feedback
Amazon Simple Email Service (Amazon SES)	Provides statistics such as email deliveries, bounces, feedback loop
	results, email open
	Supports DomainKeys Identified Mail (DKIM) and Sender Policy
	Framework (SPF)
	Flexible IP deployment: shared, dedicated, and customer-owned IPs
	Send emails using your application using AWS Console, APIs, or SMTP
	Use cases: transactional, marketing and bulk email communications
	· OSE CASES: HAIISACHOHAL, MAI KEHING AND DUK EMIAH COMMINUMICATIONS

	Scalable 2-way (outbound/inbound) marketing
	communications service
	Supports email, SMS, push, voice, and in-app messaging
	 Ability to segment and personalize messages with the
	right content to customers
	Possibility to receive replies
	Scales to billions of messages per day
Amazon Pinpoint	Use cases: run campaigns by sending marketing, bulk,
	transactional SMS messages
	· Versus Amazon SNS or Amazon SES
	 In SNS & SES you managed each message's audience,
	content, and delivery schedule
	· In Amazon Pinpoint, you create message templates,
	delivery schedules, highly-targeted segments, and full
	campaigns
	· Allows you to start a secure shell on your EC2 and
	on-premises servers
	No SSH access, bastion hosts, or SSH keys
Systems Manager - SSM Session Manager	needed
-,	No port 22 needed (better security)
	Supports Linux, macOS, and Windows
	Send session log data to S3 or CloudWatch Logs
	Execute a document (= script) or just run a
	command
	Run command across multiple instances
	(using resource groups)
	· No need for SSH
Systems Manager - Run Command	Command Output can be shown in the AWS
,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	Console, sent to S3 bucket or CloudWatch
	Logs
	Send notifications to SNS about command
	status (In progress, Success, Failed,)
	Integrated with IAM & CloudTrail
	Can be invoked using EventBridge
	Automates the process of patching managed
	instances
	· OS updates, applications updates, security
	updates
	Supports EC2 instances and on-premises
Systems Manager - Patch Manager	servers
	Supports Linux, macOS, and Windows
	Patch on-demand or on a schedule using
	Maintenance Windows
	Scan instances and generate patch compliance
	report (missing patches)
	Defines a schedule for when to perform actions on your instances Defines a schedule for when to perform actions on your instances
Systems Manager - Maintenance Windows	Example: OS patching, updating drivers, installing software, Maintenance Windows partains
	Maintenance Window contains
	• Schedule
	Duration
	Set of registered instances
	Set of registered tasks

	Simplifies common maintenance and deployment tasks of EC3 instances and other
	deployment tasks of EC2 instances and other AWS resources
	Examples: restart instances, create an AMI,
	EBS snapshot
	Automation Runbook - SSM Documents to
	define actions preformed on your EC2
Systems Manager - Automation	instances or AWS resources (pre-defined or
	custom)
	· Can be triggered using:
	Manually using AWS Console, AWS CLI or SDK
	· Amazon EventBridge
	On a schedule using Maintenance Windows
	By AWS Config for rules remediations
	Visualize, understand, and manage your AWS costs and usage over time
	Create custom reports that analyze cost and usage data.
	Analyze your data at a high level: total costs and usage across all accounts
AWS Cost Explorer	· Or Monthly, hourly, resource level granularity
	· Choose an optimal Savings Plan (to lower prices on your bill)
	• Forecast usage up to 12 months based on previous usage
Cost Explorer - Savings Plan Alternative to	Cost Explorer - Savings Plan Alternative to Reserved Instances
Reserved Instances	
	Elastic Transcoder is used to convert media files stored in S3 into media files in the
	formats required by consumer playback devices (phones etc)
	· Benefits:
Amazon Elastic Transcoder	• Easy to use
	Highly scalable - can handle large volumes of media files and large file sizes
	· Cost effective - duration-based pricing model
	Fully managed & secure, pay for what you use
	• Lambda:
	• Time limit
	Limited runtimes
	Limited temporary disk space
Batch vs Lambda	• Serverless
baten vs tambua	- Batch:
	No time limit
	· Any runtime as long as it's packaged as a Docker image
	• Rely on EBS / instance store for disk space
	• Relies on EC2 (can be managed by AWS)
	Fully managed integration service that enables you to securely transfer data between
	Software-as-a-Service (SaaS) applications and AWS
Amazon AppFlow	Sources: Salesforce, SAP, Zendesk, Slack, and ServiceNow
	Destinations: AWS services like Amazon S3, Amazon Redshift or non-AWS such as
	SnowFlake and Salesforce
	Frequency: on a schedule, in response to events, or on demand
	Data transformation capabilities like filtering and validation
	 Encrypted over the public internet or privately over AWS PrivateLink Don't spend time writing the integrations and leverage APIs immediately

Framework and adopt architectural best practices How does it work? AWS Well-Architected Tool Select your workload and answer questions Review your answers against the 6 pillars	AWS Amplify - web and mobile applications	 A set of tools and services that helps you develop and deploy scalable full stack web and mobile applications Authentication, Storage, API (REST, GraphQL), CI/CD, PubSub, Analytics, AI/ML Predictions, Monitoring, Connect your source code from GitHub, AWS CodeCommit, Bitbucket, GitLab, or
- Test systems at production scale - Automate to make architectural experimentation easier - Automate to make architectural experimentation easier - Automate to make architectures - Automate to make architectures - Automate to revolutionary architectures - Design based on changing requirements - Drive architectures using data - Improve through game days - Simulate applications for flash sale days - Department of the special experiments - Drive architectures - Design based on changing requirements - Drive architectures using data - Improve through game days - Simulate applications for flash sale days - Department of the special experiments - Design based on changing requirements - Design based on changing rethietcures - Design based days - Design based days - Design based days - Design based and sale architectural based days - Design based on changing requirements - Design based and provides -		upload directly
. 2) Security . 3) Reliability . 4) Performance Efficiency . 5) Cost Optimization . 6) Sustainability . Free tool to review your architectures against the 6 pillars Well-Architected Framework and adopt architectural best practices . How does it work? AWS Well-Architected Tool AWS Well-Architected Tool Review your answers against the 6 pillars . Obtain advice: get videos and documentations, generate a report, see the dashboard - Analyze your AWS accounts and provides recommendation on 5 categories . Cost optimization . Performance . Security . Fault tolerance . Service limits 7 CORE CHECKS Basic & Developer Support plan . S3 Bucket Permissions . Security Groups - Specific Ports Unrestricted . IAM Use (one IAM user minimum) . MFA on Root Account . EBS Public Snapshots . RDS Public Snapshots . Service Limits FULL CHECKS Business & Enterprise Support plan . Full Checks available on the 5 categories		 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
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 1 dashboard Analyze your AWS accounts and provides recommendation on 5 categories Cost optimization Performance Security Fault tolerance Service limits 7 CORE CHECKS Basic & Developer Support plan Sa Bucket Permissions Security Groups - Specific Ports Unrestricted IAM Use (one IAM user minimum) MFA on Root Account EBS Public Snapshots RDS Public Snapshots Service Limits FULL CHECKS Business & Enterprise Support plan Full Checks available on the 5 categories	Well Architected Framework 6 Pillars	 2) Security 3) Reliability 4) Performance Efficiency 5) Cost Optimization
recommendation on 5 categories Cost optimization Performance Security Fault tolerance Service limits 7 CORE CHECKS Basic & Developer Support plan S3 Bucket Permissions Security Groups - Specific Ports Unrestricted IAM Use (one IAM user minimum) MFA on Root Account EBS Public Snapshots RDS Public Snapshots Service Limits FULL CHECKS Business & Enterprise Support plan Full Checks available on the 5 categories	AWS Well-Architected Tool	 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
- S3 Bucket Permissions - Security Groups - Specific Ports Unrestricted - IAM Use (one IAM user minimum) - MFA on Root Account - EBS Public Snapshots - RDS Public Snapshots - Service Limits FULL CHECKS Business & Enterprise Support plan - Full Checks available on the 5 categories	Trusted Advisor	recommendation on 5 categories Cost optimization Performance Security Fault tolerance
reaching limits • Programmatic Access using AWS Support	Trusted Advisor - Support Plans	 S3 Bucket Permissions Security Groups - Specific Ports Unrestricted IAM Use (one IAM user minimum) MFA on Root Account EBS Public Snapshots RDS Public Snapshots Service Limits FULL CHECKS Business & Enterprise Support plan Full Checks available on the 5 categories Ability to set CloudWatch alarms when reaching limits