

LEARNING AWS

A CLOUD GURU TRAINING SESSIONS

MISCELLANEOUS NOTES

- AWS US EAST (N. VIRGINIA)
 - THIS REGION IS THE OLDEST AND IT IS WHERE MOST OF THE NEW SERVICES ARE DEPLOYED FIRST FOR CONSUMPTION BY THE PUBLIC
 - IS A REGION THAT TENDS TO GO DOWN AT LEAST ONCE A YEAR.
- AWS SUPPORT LEVELS
 - BASIC PLAN
 - DEVELOPER PLAN
 - BUSINESS PLAN

AWS GLOBAL INFRASTRUCTURE (2019)

- REGION
 - A GEOGRAPHICAL AREA THAT MAY CONTAIN 2 OR MORE AVAILABILITY ZONES
 - 24 REGIONS GLOBALLY
- AVAILABILITY ZONE
 - AN AZ IS A FACILITY SIMILAR TO A DATA CENTER
 - HAVE COMPUTERS, SAN, REDUNDANT POWER, REDUNDANT BACKBONE TO INTERNET.
 - FACILITIES WHICH ARE CLOSE TOGETHER ARE COUNTED AS BEING IN THE SAME AVAILABILITY ZONE
 - 72 AVAILABILITY ZONES GLOBALLY
- EDGE LOCATION
 - THESE ARE ENDPOINTS FOR AWS WHICH ARE USED FOR CACHING CONTENT. EXAMPLES ARE
 - AWS **CLOUDFRONT** WHICH IS AMAZON's CDN
 - THERE ARE OVER 150 EDGE LOCATIONS GLOBALLY

CORE AWS SERVICES FOR ARCHITECT CERT

- AWS GLOBAL INFRASTRUCTURE
 - REGION > AVAILABILITY ZONE > EDGE LOCATION
- COMPUTE
 - EC2, LAMBDA
- STORAGE
 - EFS, S3
- DATABASES
 - RDS, DYNAMO
- NETWORK & CONTENT DELIVERY
 - KNOWING THIS AREA INSIDE-OUT IS KEY TO PASSING THE EXAM
 - KNOW VPC & ROUTE 53
 - **CLOUDFRONT** (AMAZON'S CDN)
- SECURITY, IDENTITY & COMPLIANCE
 - KNOWING THIS AREA INSIDE-OUT IS KEY TO PASSING THE EXAM
 - KNOW IAM VERY WELL

IAM 101 WITH AWS

- CENTRALIZED CONTROL OF YOUR AWS ACCOUNT
- SHARED ACCESS TO YOUR AWS ACCOUNT
- GRANULAR PERMISSIONS
- IDENTITY FEDERATION (AD, FB, LINKEDIN, ETC.)
- MFA
- PROVIDE TEMPORARY ACCESS FOR USERS/DEVICES AND SERVICES WHERE NECESSARY
- ALLOWS SETUP OF PASSWORD ROTATION POLICIES
- INTEGRATES WITH MANY DIFFERENT AWS SERVICES
- SUPPORTS PCI DSS COMPLIANCE
- IAM IS THE MOST IMPORTANT SERVICE IN AWS BECAUSE IT CONTROLS ACCESS TO ALL OTHER AWS SERVICES.
- POWER USER ACCESS – ALLOWS ACCESS TO ALL AWS SERVICES EXCEPT THE MANAGEMENT OF GROUPS AND USERS WITHIN IAM
- READ THIS WHITE PAPER: [HTTPS://AWS.AMAZON.COM/WHITEPAPERS/](https://aws.amazon.com/whitepapers/)

AWS IAM – KEY TERMINOLOGY

- USERS
 - REPRESENTS END USERS (I.E. PEOPLE IN AN ORG)
 - ACCOUNTS FOR INDIVIDUAL USERS
 - CAN HAVE PERMISSIONS ASSIGNED TO THEM
- GROUPS
 - CAN HAVE PERMISSIONS ASSIGNED TO THEM
 - USERS CAN BE MEMBERS OF A GROUP
 - USERS WHO ARE MEMBERS OF A GROUP WILL INHERIT THE PERMISSIONS OF SAID GROUP
- ROLES
 - ARE ASSIGNED TO AWS RESOURCES
- POLICIES
 - JSON DOCUMENTS THAT DESCRIBE PERMISSIONS THAT A USER/GROUP/ROLE IS ABLE TO DO

IAM – KEY CONCEPTS

- **IMPORTANT:**
 - **ACTIVATE MFA ON THE ROOT ACCOUNT.** KEEP YOUR ROOT ACCOUNT QR CODE IN A SECURE LOCATION TO BE ABLE TO REACTIVATE IN CASE YOU LOSE YOUR PHONE.
 - **OPERATIONS ON IAM ARE GLOBAL.** WE CANNOT CREATE A USER FOR A SPECIFIC REGION OR AVAILABILITY ZONE.
- YOU CAN CREATE USERS OF 2 TYPES
 - **PROGRAMMATIC ACCESS** – TO ACCESS AWS SERVICES PROGRAMMATICALLY (FOR DEVS)
 - **CONSOLE ACCESS** – TO ACCESS AWS WEB-BASED CONSOLE TO MANAGE YOUR AWS ACCOUNT SETTINGS
- **USERS** – REPRESENT THE USERS OF AWS SERVICES. THEY CAN BE GRANTED **ACTIONS** TO **RESOURCES** (NOT RECOMMENDED)
- **GROUPS** – CONTAIN USERS AND ARE GOVERNED BY **POLICIES**
- **POLICY** – DETERMINES WHAT **ACTIONS** CAN BE TAKEN ON WHAT **RESOURCES**
- **ROLES** – DETERMINES A SECURE WAY TO GRANT PERMISSIONS TO ENTITIES WE TRUST

IAM – CREATE A BILLING ALARM

- A BILLING ALARM HELPS KEEP WATCH OVER AN AWS ACCOUNT.
- **CLOUDWATCH** – IT IS THE AWS SERVICE THAT CAN BE USED TO CREATE BILLING ALARMS
- **SNS TOPIC** – IT IS AN AWS SERVICE CALLED “SIMPLE NOTIFICATION SERVICE” (SNS) THAT ALLOWS SERVICES TO SEND MESSAGES TO OTHER SERVICES.

AWS SERVICE – S3

- S3 – SIMPLE STORAGE SERVICE
- S3 ALLOWS YOU TO CREATE FOLDERS (AWS CALLS THEM **BUCKETS**) WHERE YOU CAN STORE OBJECTS (FILES)
- FILES CAN HAVE 0 BYTES TO 5 TERABYTES
- BY DEFAULT, AN ACCOUNT CAN HAVE 100 BUCKETS
- **BUCKETS** HAVE A GLOBALLY UNIQUE WEB ADDRESS
- UPLOADING LARGE OBJECTS:
<https://docs.aws.amazon.com/AmazonS3/latest/dev/uploadingobjects.html>
- WHEN UPLOADING FILES TO A BUCKET LOCATION YOU, IF IT WAS SUCCESSFUL, YOU WILL GET AN **HTTP 200** RESPONSE CODE.
- EACH OBJECT (FILE) HAS CERTAIN ATTRIBUTES
 - KEY – THE NAME OF THE FILE. IT MUST BE UNIQUE WITHIN A **BUCKET**
 - VALUE – THE CONTENT OF THE FILE (BYTES)
 - VERSION – THE VERSION OF THE FILE. IMPORTANT WHEN YOU ARE UPLOADING OBJECTS WITH THE SAME KEY TO THE BUCKET LOCATION.
 - METADATA – ADDITIONAL DATA ABOUT THE **S3** OBJECT THAT IS STORED IN A BUCKET
 - SUBRESOURCES
 - ACCESS CONTROL LISTS
 - TORRENTS

AWS SERVICE – S3

- DATA CONSISTENCY
 - READ AFTER WRITE CONSISTENCY FOR 1ST PUT OF A NEW OBJECTS
 - EVENTUAL CONSISTENCY FOR OVERWRITE PUTs AND DELETES FOR EXISTING OBJECT (CAN TAKE SOME TIME TO PROPAGATE)
- S3 GUARANTEES
 - BUILT FOR 99.99% AVAILABILITY – 52.6 SECONDS IN A YEAR OF DOWNTIME
 - AMAZON GUARANTEES 99.9% - 8.76 HOURS IN A YEAR OF DOWNTIME
 - AMAZON GUARANTEES 99.999999999% OF DURABILITY –
- FEATURES
 - TIERED STORAGE
 - LIFECYCLE MANAGEMENT
 - VERSIONING
 - ENCRYPTION AT REST
 - MFA DELETE
 - SECURITY WITH ACL AND BUCKET POLICIES

AWS SERVICES – S3

- STORAGE CLASSES (TIERS)

- **S3 STANDARD**

- 99.99% AVAILABILITY
- 99.999999999% DURABILITY (11 x 9's)
- STORED REDUNDANTLY ACROSS MULTIPLE DEVICES IN MULTIPLE FACILITIES
- DESIGNED TO SUSTAIN THE LOSS OF 2 FACILITIES CONCURRENTLY

- **S3 INFREQUENT ACCESS (IA)** – CHEAPER THAN S3 STANDARD BUT DATA RETRIEVAL FEES APPLY

- AVAILABILITY ??

- **S3 ONE ZONE IA** – LOWER COST OPTION WHEN IA APPLIES BUT DATA RESILIENCY IS NOT REQUIRED

- AVAILABILITY IS 99.50%

- **S3 INTELLIGENT TIER** – DESIGNED TO OPTIMIZE STORAGE COSTS BY AUTOMATICALLY MOVING DATA TO THE MOST COST EFFECTIVE TIER WITHOUT PERFORMANCE IMPACT OR OPERATIONAL OVERHEAD

- **S3 GLACIER** – CHEAPER OPTION THAN PREVIOUS. COMPETITIVE WITH ON-PREM DATA STORAGE. RETRIEVAL TIMES ARE CONFIGURABLE FROM MINUTES TO HOURS.

- **S3 GLACIER DEEP ARCHIVE** – LOWEST COST STORAGE CLASS FOR FILE STORAGE CASES WHERE A RETRIEVAL TIME OF 12 HOURS IS ACCEPTABLE.

- READ THROUGH THE S3 FAQ'S BEFORE TAKING THE EXAM

S3 – ENCRYPTION FOR BUCKETS

- BY DEFAULT NEWLY CREATED BUCKETS ARE PRIVATE
- ACCESS CONTROL CAN BE SET FOR A BUCKET WITH
 - BUCKET POLICIES
 - ACCESS CONTROL LISTS
- S3 BUCKETS CAN LOG ACCESS TO IT
 - LOGS CAN BE SENT TO ANOTHER BUCKET ON THE SAME ACCOUNT
 - LOGS CAN BE SENT TO A BUCKET OWNED BY ANOTHER ACCOUNT
- BUCKETS CAN BE ENCRYPTED
 - IN TRANSIT – USE HTTPS WITH SSL/TSL
 - AT REST – ENCRYPTING DATA THAT IS STORED
 - SERVER SIDE – AWS ENCRYPTS THE OBJECT WHEN IT IS STORED
 - S3 MANAGED KEYS – AWS MANAGES THE ENCRYPTION KEYS FOR YOU WITH SERVER SIDE ENCRYPTION (**SSE-S3**)
 - AWS KEY MANAGEMENT SERVICE – YOU AND AWS MANAGE THE ENCRYPTION KEYS (**SSE-KMS**)
 - CUSTOMER PROVIDED THE KEYS (**SSE-C**)
 - CLIENT SIDE – YOU ENCRYPT THE OBJECT BEFORE YOU UPLOAD IT TO AWS S3 BUCKET

S3 - VERSIONING

- STORES ALL VERSIONS OF AN OBJECT INCLUDING WRITES AND EVEN DELETES
- REMOVING A DELETE MARKER WILL RE-ENABLE THE FILE THAT WAS DELETED. IT IS LIKE A SOFT DELETE.
- CAN BE USED AS A BACKUP TOOL
- ONCE ENABLED IT CANNOT BE REMOVED FROM THE BUCKET. IT CAN ONLY BE SUSPENDED. IF YOU NEED TO REMOVE YOU'LL NEED TO
- INTEGRATES WITH LIFECYCLE RULES
- CAN ACTIVATE MFA DELETION FOR AN ADDITIONAL LAYER OF SECURITY

S3 – LIFECYCLE MANAGEMENT

- AUTOMATES MOVING OBJECTS BETWEEN THE DIFFERENT STORAGE TIERS
- CAN BE USED IN CONJUNCTION WITH VERSIONING (VERSIONING DOES NOT HAVE TO BE TURNED ON TO USE LIFECYCLE MANAGEMENT)
- CAN BE APPLIED TO CURRENT VERSIONS OR PREVIOUS VERSIONS

S3 – CROSS-REGION REPLICATION

- VERSIONING MUST BE ENABLED IN BOTH THE SOURCE AND DESTINATION BUCKETS
- REGIONS HAVE TO BE UNIQUE WHEN USING CROSS-REGION REPLICATION.
- SAME-REGION REPLICATION IS POSSIBLE (AS OF RECENT IN 2019?) TO A BUCKET IN THE SAME REGION – BUCKET NAMES HAVE TO BE UNIQUE.
- FILES IN AN EXISTING BUCKET ARE NOT REPLICATED AUTOMATICALLY. ONLY FILES THAT ARE ADDED TO A BUCKET AFTER **CRR** IS TURNED ON ARE REPLICATED.
- ALL SUBSEQUENT UPDATED FILES WILL BE REPLICATED AUTOMATICALLY
- DELETE MARKERS ARE NOT REPLICATED
- DELETING INDIVIDUAL VERSIONS ARE NOT REPLICATED

S3 – TRANSFER ACCELERATION

- USES THE **CLOUDFRONT EDGE NETWORK** TO SPEED UP YOUR UPLOADS TO S3 BUCKETS
- YOU CAN UPLOAD TO CFE LOCATION WHICH THEN IT WILL TRANSFER TO AN S3 BUCKET
- USE THE TESTING TOOL TO MEASURE TIMES AND FIGURE OUT BETTER EDGE LOCATIONS TO USE FOR YOUR APPLICATION
 - <HTTP://S3-ACCELERATE-SPEEDTEST.S3-ACCELERATE.AMAZONAWS.COM/EN/ACCELERATE-SPEED-COMPERSION.HTML>

AWS CLOUDFRONT

- EDGE LOCATION: AN AWS NETWORK LOCATION WHERE CONTENT IS CACHED. THIS IS PHYSICALLY SEPARATE FROM A AWS REGION OR AWS AVAILABILITY ZONE (AZ)
- ORIGIN: IS THE LOCATION WHERE THE ORIGINAL CONTENT RESIDES. IT CAN BE AN S3 BUCKET, AN EC2 INSTANCE, AND ELB OR ROUTE 53.
- DISTRIBUTION: IS THE NAME GIVEN TO A CDN WHICH IS A COLLECTION OF EDGE LOCATIONS
- TYPES OF CF
 - WEB DISTRIBUTION – USED FOR WEBSITES
 - RTMP – USED FOR MEDIA STREAMING. RTMP IS ADOBE'S REAL-TIME MESSAGING PROTOCOL
- EDGE LOCATIONS ARE NOT JUST READ-ONLY. YOU CAN WRITE TO THEM TOO.
- OBJECTS ARE CACHED AT AN EDGE LOCATION BASED ON A TIME-TO-LIVE (TTL) VALUE
- YOU CAN CLEAR (INVALIDATE) CACHED OBJECTS BEFORE THEIR TTL VALUE. HOWEVER, YOU WILL BE CHARGED EXTRA FOR THAT ACTION

AWS – CLOUDFRONT LAB

- CLOUDFRONT IS A GLOBAL SERVICE
- ACCESS TO CACHED CONTENT CAN BE RESTRICTED TO SIGNED URLs/COOKIES
- TO ACTIVATE, YOU NEED TO CREATE/DEFINE A CLOUDFRONT DISTRIBUTION
- THIS IS NOT COVERED BY AWS' FREE TIER – YOU WILL INCUR IN SOME COSTS
- TO DELETE A CLOUDFRONT DISTRIBUTION YOU NEED TO FIRST MARK IT AS “DISABLED” THEN IT CAN BE DELETED

SNOWBALL

- IT IS A SERVICE TO HELP MOVE MASSIVE AMOUNTS OF DATA
 - FROM ON-PREM TO AWS
 - FROM AWS TO ON-PREM
- NEED TO ISSUE A REQUEST. AWS SENDS YOU A DEVICE THAT YOU NEED TO ACTIVATE AND CONNECT TO YOUR LOCAL NETWORK
- SNOWBALL - IT IS SECURE, FULL CHAIN OF CUSTODY, 256 BIT ENCRYPTION. CAN TRANSFER 50TB OR 80TB. AWS PERFORMS ERASURE
- SNOWBALL EDGE – 100TB. IT IS LIKE HAVING A MINI-AWS AT YOUR DISPOSAL. IT CAN BE CLUSTERED WITH OTHER SE.
- SNOWMOBILE – 100PB. IT IS A TRUCK THAT HAS SERVERAL SNOWBALL EDGES NETWORKED
- USED SNOWBALL WHEN YOU HAVE
 - 44 MBIT/S TO INTERNET AND YOU NEED TO TRANSFER MORE THAN 2TB OF DATA
 - 100 MBIT/S TO INTERNET AND YOU NEED TO TRANSFER MORE THAN 5TB OF DATA
 - 1000 MBIT/S TO INTERNET AND YOU NEED TO TRANSFER MORE THAN 60TB OF DATA

AWS – STORAGE GATEWAY

- USED TO MOVE DATA INTO AWS. IT IS A SERVICE THAT CONNECTS ON-PREM SOFTWARE APPLIANCE TO AWS CLOUD
- IT CAN BE A VIRTUAL (DOWNLOAD A VM IMAGE) OR PHYSICAL DEVICE
- OPTIONS
 - **FILE GATEWAY** – FILES ARE STORED AS OBJECTS ON S3 BUCKETS
 - IDEAL FOR EXISTING OPERATIONS WITH LARGE BODIES OF DOCUMENTS
 - **VOLUME GATEWAY** – iSCSI BLOCK PROTOCOL - COPIES OF DATA IN HD FOR
 - STORED VOLUMES – ALL DATA ON THE HD WILL BE COPIED/BACKED-UP TO AN S3 BUCKET
 - CACHED VOLUMES – LET'S YOU USE S3 AS YOUR PRIMARY DATA STORAGE WHILE RETAINING FREQUENTLY USED DATA ON-PREM STORAGE INFRASTRUCTURE
 - **TAPE GATEWAY** – A VIRTUAL TAPE LIBRARY
 - ALLOWS TO LEVERAGE YOUR EXISTING TAPE INFRA AND MOVING YOUR BACKUPS TO CLOUD
 - USED FOR ARCHIVAL PURPOSES
 - YOU CAN USE GLACIER SERVICE TO REDUCE COSTS OF ARCHIVAL OF INFREQUENTLY USED DATA

EC2 - 101

- ELASTIC COMPUTE CLOUD
- IS A WEB SERVICE FOR RESIZABLE COMPUTE IN THE CLOUD
- REDUCES THE TIME IT TAKES TO OBTAIN AND REBOOT SERVERS
- ALLOWS FOR QUICK SCALING (UP/DOWN) OF COMPUTE REQUIREMENTS
- VIRTUAL MACHINES IN THE CLOUD
- RUN ON XEN AND NITRO HYPERVISOR SOFTWARE
- PRICING
 - PAY AS YOU GO
 - PAY FOR WHAT YOU USE
 - PAY LESS AS YOU USE MORE EVEN AS CAPACITY IS RESERVED
 - TIERS:
 - **R**ESERVED – CAPACITY RESERVATION FOR 1 YR OR 3 YR
 - **O**N-DEMAND
 - **D**EDICATED HOST – EXCLUSIVE FOR AN ACCOUNT
 - **S**POT – BID YOUR PRICE FOR AVAILABLE CAPACITY AT LOWER RATE

EC2 TIER – ON DEMAND

- PAY A FIXED RATE
- PRICING IS BY THE HOUR OR BY THE SECOND
- NO COMMITMENTS

EC2 TIER – RESERVED INSTANCES

- CONTRACTS FOR 1 YR OR 3 YR WITH DEEPLY DISCOUNTED PRICES FROM THE ON-DEMAND TIER
- STANDARD RESERVED INSTANCES
 - UP TO 75% DISCOUNT.
 - THE MORE YOU PAY UPFRONT OR THE LONGER THE CONTRACT THE MORE YOU SAVE
 - ONCE THE CONTRACT STARTS YOU CANNOT CHANGE THE COMPUTING CHOICES YOU MAKE
- CONVERTIBLE RESERVED INSTANCES
 - UP TO 45% DISCOUNT
 - ALLOWS YOU TO CHANGE THE COMPUTING CHOICES WITHIN THE CONTRACT TO EQUAL OR GREATER VALUE
- SCHEDULED RESERVED INSTANCES
 - THESE INSTANCES ALLOW YOU TO USE COMPUTE CAPACITY DURING A SPECIFIC TIME WINDOW (FRACTION OF DAY, WEEK, A MONTH)

EC2 TIER – SPOT & DEDICATED HOST PRICING

- SPOT
 - PRICE MOVES UP/DOWN DEPENDING ON AMAZON'S DEMAND AND AVAILABLE CAPACITY
 - DESIGNED FOR APPS THAT HAVE FLEXIBLE START/STOP LIFECYCLE
 - IDEAL FOR APPS THAT ARE ONLY FEASIBLE AT VERY LOW COMPUTE PRICES
 - ACCOMMODATES URGENT COMPUTING NEEDS AND LARGE AMOUNT OF ADDITIONAL CAPACITY
 - IF THE SYSTEM TERMINATES YOUR WORK LOAD YOU WILL NOT BE CHARGED FOR FRACTION OF HOUR.
IF YOU TERMINATE THE INSTANCE, YOU WILL BE CHARGED FOR WHATEVER FRACTIONS YOU USED.
- DEDICATED Host
 - WHEN REGULATION DOES NOT ALLOW FOR SHARING COMPUTE ENVIRONMENT (I.E MULTITENANT VIRTUALIZATION)
 - WHEN SOFTWARE LICENSING DOES NOT ALLOW MULTITENANT VIRTUALIZATION

EC2 INSTANCE TYPES

- F - FPGA
- I - IOPS
- G - GRAPHICS
- H – HIGH DISK I/O
- T – CHEAP T2 MICRO
- D – DENSE STORAGE
- R - RAM
- M- MAIN CHOICE (GEN)
- C - COMPUTE
- P – PROC PIX/BIT COIN
- X – EXTREME MEMORY
- Z – EXTREME MEM & CPU
- A – ARM BASED LOADS
- U – BARE METAL

Family	Specialty	Use case
F1	Field Programmable Gate Array	Genomics research, financial analytics, real-time video processing, big data etc
I3	High Speed Storage	NoSQL DBs, Data Warehousing etc
G3	Graphics Intensive	Video Encoding/ 3D Application Streaming
H1	High Disk Throughput	MapReduce-based workloads, distributed file systems such as HDFS and MapR-FS
T3	Lowest Cost, General Purpose	Web Servers/Small DBs
D2	Dense Storage	Fileservers/Data Warehousing/Hadoop
R5	Memory Optimized	Memory Intensive Apps/DBs
M5	General Purpose	Application Servers
C5	Compute Optimized	CPU Intensive Apps/DBs
P3	Graphics/General Purpose GPU	Machine Learning, Bit Coin Mining etc
X1	Memory Optimized	SAP HANA/Apache Spark etc
Z1D	High compute capacity and a high memory footprint.	Ideal for electronic design automation (EDA) and certain relational database workloads with high per-core licensing costs.
A1	Arm-based workloads	Scale-out workloads such as web servers
U-6tb1	Bare Metal	Bare metal capabilities that eliminate virtualization overhead



EC2 – LAB NOTES

- CREATING A SIMPLE LINUX / APACHE HTTP SERVICE

- SPIN UP AN EC2 INSTANCE

- SELECT A TYPE (ELIGIBLE FOR FREE TIER)
 - TAKE ALL DEFAULTS
 - STORAGE
 - YOU CAN ADD ADDITIONAL STORAGE SPACE BY ADDING VOLUMES (ON TOP OF OS REQUIREMENTS)
 - ROOT DEVICE VOLUMES CAN BE ENCRYPTED (IT WAS NOT ALWAYS THE CASE)

- TERMINATION PROTECTION

- TICK THE BOX FOR WARNING FOR ACCIDENTAL TERMINATION. IT IS TURNED OFF BY DEFAULT
 - PROVIDE SOME TAG ATTRIBUTES
 - GENERATE KEYS PAIR (PUB/PRIVATE)
 - SSH TO THE NEW INSTANCE WITH IP ADDRESS (USER IS EC2-USER)
 - DO "SUDO SU" TO BE ROOT
 - INSTALL APACHE HTTPD SERVICE
 - SETUP A BASIC INDEX.HTML FILE
 - TURN ON HTTPD SERVICE
 - VOILA!

EXAM TIPS:

- TERMINATION PROTECTION IS TURNED “ON” BY DEFAULT
- ON EBS-BACKED INSTANCE, THE DEFAULT ACTION IS FOR THE ROOT EBS VOLUME TO BE DELETED WHEN INSTANCE IS TERMINATED
- EBS ROOT VOLUMES CAN NOW BE ENCRYPTED
- YOU CAN ENCRYPT ROOT VOLUMES WITH 3RD PARTY TOOLS LIKE BITLOCKER
- YOU CAN ADD MORE VOLUMES AND ENCRYPT THOSE AS WELL

SECURITY GROUPS

- ALL IN-BOUND TRAFFIC IS BLOCKED BY DEFAULT
- ALL OUT-BOUND TRAFFIC IS ALLOWED
- CHANGES TO SECURITY GROUPS TAKE PLACE IMMEDIATELY
- YOU CAN HAVE ANY NUMBER OF EC2 INSTANCES IN A SECURITY GROUP
- YOU CAN HAVE MULTIPLE SECURITY GROUPS ATTACHED TO AN EC2 INSTANCE
- SEC GROUPS ARE **STATEFUL**
 - WHEN OPENING PORTS, IT OPENS FOR BOTH IN-BOUND AND OUT-BOUND TRAFFIC
- NETWORK ACL ARE STATELESS
 - WHEN OPENING PORTS, YOU HAVE TO SPECIFICALLY DECLARE THE IN-BOUND AND OUT-BOUND RULES
- IF YOU CREATE A RULE ALLOWING TRAFFIC IN, AUTOMATICALLY THERE WILL BE A RULE CREATED THAT ALLOWS TRAFFIC BACK OUT
- IP ADDRESSES CANNOT BE SPECIFICALLY BLOCKED WITH SECURITY GROUPS
- NETWORK ACL's ALLOW YOU TO DEFINE BLOCKED IP ADDRESSES

EBS 101 – ELASTIC BLOCK STORE

- PROVIDES PERSISTENT BLOCK STORAGE VOLUMES FOR USE WITH EC2 INSTANCES
- THINK OF IT AS HARD DRIVE IN THE CLOUD FOR A GIVEN VIRTUAL MACHINE
- THIS SERVICE IS BASED ON BLOCK STORAGE (EFS AND FSX ARE ALSO BASED ON BLOCK STORAGE)
- READ THIS TO LEARN ABOUT OPTIMIZING EBS VOLUME PERFORMANCE
- EACH EBS IS AUTOMATICALLY REPLICATED WITHIN IT'S AZ TO INCREASE PERFORMANCE. THIS DONE TO REDUCE LATENCY BETWEEN THE VIRTUAL MOTHER-BOARD AND THE VIRTUAL HDD
 - DURABILITY
 - AVAILABILITY
 - PREVENT SINGLE POINT OF FAILURE
- THERE IS A RESTRICTION OF 50:1 RATIO BETWEEN IOPS AND VOLUME SIZE
- THERE ARE 5 TYPES OF EBS
 - GENERAL PURPOSE SSD
 - PROVISIONED IOPS SSD
 - THROUGHPUT OPTIMIZED HDD
 - COLD HDD
 - EBS MAGNETIC

ELASTIC BLOCK STORAGE – TYPES

G.P. T.E.C.

SOLID STATE DRIVES (SSD)

- **GENERAL PURPOSE SSD**

- USE CASE: MOST WORKLOADS
- BALANCES PRICES V PERFORMANCE
- SIZE: 1GB – 16TB
- MAX IOPS 16,000
- API NAME GP2

- **PROVISIONED IOPS SSD**

- USE CASE: DATABASES
- HIGHEST PERFORMANCE FOR MISSION CRITICAL CASES
- SIZE: 4GB – 16TB
- MAX IOPS 64,000
- API NAME IO1

HARD DISK DRIVES (MAGNETIC)

- **THROUGHPUT OPTIMIZED HDD**

- USE CASE: BIG DATA/DATAWAREHOUSE
- LOW COST FREQ. ACCESS HIGH THROUGHPUT WORKLOADS
- API NAME: ST1
- SIZE: 500 GB – 16 TB
- MAX IOPS: 500

- **EBS MAGNETIC**

- USE CASE: INFREQUENTLY ACCESSED DATA AND YOU ARE NOT USING GLACIER
- PREVIOUS GENERATION HDD
- API NAME: STANDARD
- SIZE: 1 GB – 1 TB
- MAX IOPS: 40 – 200

- **COLD HDD**

- USE CASE: FILE SERVERS
- LOWEST Cost HDD. LESS FREQUENTLY ACCESSED WORKLOADS
- API NAME: SC1
- SIZE: 500 GB – 16 TB
- MAX IOPS: 250

ELASTIC BLOCK STORE – EXAM TIPS

- VOLUMES EXIST ON EBS. EBS CAN BE THOUGHT OF AS VIRTUAL HDD ATTACHED TO A VM
- CREATING AN EC2 INSTANCE WILL AUTOMATICALLY CREATE A ROOT VOLUME
- SNAPSHOTS (A “PICTURE” OF A VOLUME) ARE STORED IN S3
- SNAPSHOTS ARE A POINT-IN-TIME COPY OF A VOLUME
- SNAPSHOTS ARE INCREMENTAL. SUBSEQUENT “PICTURES” OF THE SAME VOLUME WILL ONLY COPY DIFF BLOCKS BETWEEN THE ORIGINAL PICTURE AND THE SUBSEQUENT PICTURES.
- THE FIRST SNAPSHOT TAKES THE LONGEST TO CREATE.
- IN A PRODUCTION CASE, YOU SHOULD STOP THE EC2 INSTANCE BEFORE CREATING A SNAPSHOT. THIS WILL ENSURE THAT THE SNAPSHOT IS PRISTINE.
- SNAPSHOTS CAN BE CREATED ON A RUNNING INSTANCE.
- YOU CAN CREATE AN **AMI** (AMAZON MACHINE IMAGE) FROM EITHER VOLUMES OR SNAPSHOTS
- VOLUME SIZES AND TYPE CAN BE CHANGED ON THE FLY. YOU MAY NEED TO ACCESS THE SPECIFIC OS TO ENSURE THAT THE ADDED VOLUMES ARE VISIBLE TO THE OS.
- **AEV** - YOU CANNOT DECREASE THE SIZE OF A VOLUME, ONLY INCREASE.
- EBS VOLUMES WILL ALWAYS BE ON THE SAME AZ AS THE EC2 INSTANCE THAT OWNS THEM
- TO MOVE AN EC2 INSTANCE FROM ONE AZ TO ANOTHER, YOU NEED TO FIRST CREATE A SNAPSHOT OF IT'S ROOT VOLUME, THEN CREATE AN AMI OFF OF IT AND THEN USE AMI TO LAUNCH THE EC2 INSTANCE IN ANOTHER AZ.
- TO MOVE AN EC2 INSTANCE FROM ONE REGION TO ANOTHER, YOU NEED TO SNAPSHOT > AMI AND COPY THE AMI TO THE DESIRED REGION. THEN YOU CAN LAUNCH THE EC2 INSTANCE IN THE DESIRED REGION/AZ.
- YOU CAN PERFORM ACTIONS ON SNAPSHOTS WITH AWS API, CLI OR AWS CONSOLE

AMI TYPES – EBS **vs** INSTANCE STORE

DIFFERENCES BETWEEN EBS AND INSTANCE STORE

- ALL AMI'S ARE CATEGORIZED BY THE TYPE OF BACKING OF THE VOLUME OF THE ROOT DEVICE (WHERE THE **OS** IS INSTALLED)
- **EBS** - BACKED
 - THE ROOT DEVICE IS A VOLUME CREATED FROM AN EBS SNAPSHOT
 - IT PERSISTS THROUGH TIME AFTER THE EC2 INSTANCE IS STOPPED
 - IF THE HYPERVISOR RUNNING AN EBS-BACKED VM STOPS WORKING YOU CAN SIMPLY RE-START IT
- **INSTANCE STORE** – BACKED
 - THE ROOT DEVICE IS A VOLUME CREATED FROM A TEMPLATE STORED IN AMAZON S3
 - IS-BACKED EC2 INSTANCES CAN ONLY BE REBOOTED OR TERMINATED. CANNOT BE STOPPED
 - IF THE HYPERVISOR RUNNING AN IS-BACKED VM CRASHES YOU WILL LOSE ALL OF YOUR DATA/INFO ON THAT INSTANCE

S. O. L. A. R. AMI

- SELECTION OF **AMI** CAN BE BASED ON
 - **S**TORAGE OF THE ROOT DEVICE VOLUME TYPE
 - EBS BACKED VOLUMES
 - INSTANCE STORE, ALSO KNOWN AS Ephemeral (Temporary) STORAGE
 - **O**PERATING SYSTEM
 - **L**AUNCH PERMISSIONS
 - **A**rchitecture (32 BIT OR 64 BIT)
 - **R**EGION OR AZ

ENCRYPTED ROOT DEVICE VOLUMES AND SNAPSHOTS

- SNAPSHOTS OF ENCRYPTED VOLUMES ARE ENCRYPTED AUTOMATICALLY
- VOLUMES RESTORED FROM ENCRYPTED SNAPSHOTS ARE ENCRYPTED AUTOMATICALLY
- YOU CAN SHARE SNAPSHOTS ONLY IF THEY ARE UNENCRYPTED
- SNAPSHOTS CAN BE SHARED WITH OTHER AWS ACCOUNTS OR MADE PUBLIC
- YOU CAN NOW ENCRYPT VOLUMES UPON **CREATION** OF THE EC2 INSTANCE

TO ENCRYPT A VOLUME AFTER IT HAS BEEN CREATED:

- CREATE A SNAPSHOT OF THE UNENCRYPTED VOLUME
- CREATE A COPY OF THE SNAPSHOT AND SELECT AN ENCRYPTION OPTION
- CREATE AN AMI FROM THE ENCRYPTED SNAPSHOT
- LAUNCH NEW EC2 INSTANCE OFF OF THE ENCRYPTED AMI

CLOUDWATCH 101

MONITORS VIRTUAL THINGS LIKE:

- COMPUTE
 - EC2 INSTANCES
 - AUTOSCALING GROUPS
 - ELASTIC LOAD BALANCERS
 - ROUTE53 HEALTHCHECKS
- STORAGE & CONTENT DELIVERY
 - EBS VOLUMES
 - STORAGE GATEWAYS
 - CLOUDFRONT

MONITORS HOST-LEVEL EVENTS LIKE:

- CPU
- NETWORK
- DISK
- STATUS CHECK
 - SYSTEM – UNDERLYING HYPERVISOR (PHYSICAL COMPUTER)
 - INSTANCE – EC2 INSTANCE (VM)

EXAM TIPS:

- **CLOUDWATCH** IS USED FOR MONITORING PERFORMANCE
- IT MONITORS MOST OF AWS INFRA AS WELL AS VMs AND APPS RUNNING ON AWS INFRA
- **CLOUDWATCH** WITH EC2 WILL MONITOR EVENTS EVERY 5 MINUTES
- WITH DETAILED MONITORING YOU CAN EVALUATE IN 1 MINUTE INTERVALS
- YOU CAN CREATE CLOUDWATCH ALARMS THAT CAN TRIGGER NOTIFICATIONS
- **CLOUDWATCH** IS ABOUT PERFORMANCE (HOW ARE THE HOST MACHINE AND VM MACHINE OPERATING).
- **Do NOT CONFUSE** WITH **CLOUDTRAIL**, WHICH MONITORS AUDITING (WHO IS DOING WHAT ON WHICH RESOURCE)
- YOU CAN CREATE DASHBOARDS TO VISUALIZE:
 - ALARMS – TO NOTIFY OF THRESHOLDS ARE HIT
 - EVENTS – TO RESPOND TO STATE CHANGES IN AWS RESOURCES
 - LOGS – AGGREGATE, MONITOR AND STORE LOGS

AWS COMMAND LINE INTERFACE (CLI)

- YOU CAN INTERACT WITH AWS FROM ANYWHERE IN THE WORLD JUST BY USING THE COMMAND LINE
- YOU WILL NEED TO SETUP A USER WITH PROGRAMMATIC CREDENTIALS IN IAM
- COMMANDS ARE NOT IN THE EXAM BUT IT IS USEFUL TO KNOW BASIC CLI COMMANDS. FOR EXAMPLE
 - **AWS CONFIGURE** – TO SETUP YOUR CREDENTIALS ON AN EC2 INSTANCE. HOWEVER THIS IS NOT RECOMMENDED SINCE IF YOUR EC2 INSTANCE IS COMPROMISED, THEN YOUR CREDENTIALS ARE TOO. USE ROLES TO AVOID THESE ISSUES.
 - **AWS S3 LS** – THIS LISTS ALL OF THE S3 BUCKETS
 - **AWS S3 MB <BUCKET_NAME>** – THIS ALLOWS YOU TO CREATE A NEW BUCKET. THE BUCKET_NAME MUST BE UNIQUE ACROSS ALL EXISTING BUCKET NAMES IN AMAZON S3

IAM - ROLES

- ROLES ARE MORE SECURE THAN STORING ACCESS KEYS AND SECRET ACCESS KEY ON INDIVIDUAL EC2 INSTANCES
- ROLES ARE EASIER TO MANAGE (I.E. IMAGINE YOU NEED TO MODIFY ACCESS KEYS ON 1000 EC2 INSTANCES)
- ROLES CAN BE ASSIGNED TO AN EC2 INSTANCE EVEN AFTER IT IS CREATED WITH EITHER THE WEB CONSOLE APP OR CLI TOOL
- ROLES ARE UNIVERSAL. YOU CAN USE THEM IN ANY REGION

USING BOOTSTRAP SCRIPTS

- USEFUL TOOL TO ADD PROGRAMMATIC BEHAVIORS WHEN LAUNCHING AN EC2 INSTANCE
- THESE REQUIRE A ROLE WITH SYS-ADMIN PRIVILEGE IS ATTACHED TO THE INSTANCE ESPECIALLY IF YOU NEED TO RUN AWS CLI COMMANDS DURING THE BOOTSTRAPPING PROCESS

EC2 INSTANCE METADATA

- THIS CONTAINS SPECIFIC INFORMATION ABOUT A GIVEN EC2 INSTANCE
- TO GET TO THIS INFORMATION YOU NEED TO EXECUTE A CURL COMMAND TO
 - <HTTP://169.254.169.254/LATEST/META-DATA>
- TO GET THE DATA OF THE BOOTSTRAP SCRIPT, YOU NEED TO EXECUTE A CURL COMMAND TO
 - <HTTP://169.254.169.254/LATEST/USER-DATA>

ELASTIC FILE SYSTEM (EFS)

- EFS
 - SUPPORTS THE NETWORK FILE SYSTEM v4 (NFSv4) PROTOCOL
 - YOU PAY ONLY FOR THE STORAGE YOU USE (NO PRE PROVISIONING IS REQUIRED)
 - CAN SCALE UP TO PETABYTES
 - CAN SUPPORT THOUSANDS OF CONCURRENT CONNECTIONS
 - DATA IS STORED ACROSS MULTIPLE AZ WITHIN A REGION
 - PROVIDES READ AFTER WRITE CONSISTENCY
 - MICROSOFT WINDOWS IS NOT SUPPORTED
- **NOT IN ACG LECTURE**
 - CONFIGURING CLIENT ACCESS ON EFS
 - CONFIGURE FILE SYSTEM POLICY
 - CONFIGURE ACCESS POINTS
 - **NEED TO DISCUSS THIS WITH THE COLLEAGUES PREPARING FOR THE EXAM**

EC2 PLACEMENT GROUPS

- RELATES TO THE PLACEMENT OF A SET OF EC2 INSTANCES WITHIN A AWS AVAILABILITY ZONE
- CLUSTERED PLACEMENT GROUP
 - EC2 INSTANCES ARE PLACED AS CLOSE TOGETHER AS POSSIBLE
 - IDEAL IF YOU ARE LOOKING TO HAVE LOW NETWORK LATENCY, HIGH THROUGH PUT OR BOTH
 - CANNOT SPAN MULTIPLE AZ
 - AWS RECOMMENDS HOMOGENEOUS INSTANCES
- SPREAD PLACEMENT GROUP
 - EACH EC2 INSTANCE IS CREATED ON A SEPARATE HARDWARE RACK WITH SEPARATE POWER AND PING
 - IDEAL WHEN YOU HAVE TO ENSURE THAT EC2 INSTANCES WILL NOT BE SUSCEPTIBLE TO HARDWARE RACK FAILURE WITHIN THE AVAILABILITY ZONE
 - CAN SPAN AZ WITHIN A REGION
 - THERE IS MAX OF 7 EC2 INSTANCES PER AZ
- PARTITIONED PLACEMENT GROUP
 - LOGICAL SEGMENTS (PARTITIONS) ARE CREATED AND EC2 INSTANCES ARE ADDED TO EACH SEGMENT
 - LATENCY/THROUGHPUT IS LOW WITHIN A PARTITION, BUT SUSCEPTIBLE TO HARDWARE FAILURE
 - HARDWARE FAILURES OF ONE PARTITION DO NOT AFFECT INSTANCES RUNNING IN OTHER PARTITIONS
 - IDEAL FOR IMPLEMENTATIONS LIKE HADOOP (HDFS), HBASE AND CASSANDRA
 - CAN SPAN AZ WITHIN A REGION

PLACEMENT GROUP RESTRICTIONS:

- THE NAME OF PLACEMENT GROUP HAS TO BE UNIQUE WITHIN AN AWS ACCOUNT
- ONLY CERTAIN EC2 INSTANCE TYPES CAN BE LAUNCHED WITHIN A PLACEMENT GROUP
 - COMPUTE OPTIMIZED
 - GPU
 - MEMORY OPTIMIZED
 - STORAGE OPTIMIZED
- GROUPS CANNOT BE MERGED
- YOU CANNOT MOVE AN EXISTING INSTANCE INTO A PLACEMENT GROUP. YOU HAVE TO CREATE A AMI FROM AN EXISTING INSTANCE AND THEN LAUNCH THAT AMI INTO THE PLACEMENT GROUP

DATABASES 101

UNDERSTAND DIFFERENCES BETWEEN RELATIONAL DATABASES (RDBMS) AND NO-SQL DBs

- RDBMS FLAVORS AVAILABLE AT AWS
 - MS-SQL
 - ORACLE
 - MYSQL
 - POSTGRESQL
 - AURORA
 - MARIADB
- NO-SQL FLAVORS AVAILABLE AT AWS
 - DYNAMODB (JSON DOC OBJECT STORE / KEY-VALUE PAIRS)
- UNDERSTAND ARCHITECTURE CONFIGURATION AND WHICH YOU WOULD USE FOR WHAT SCENARIO
 - MULTI-AZ – FOR DISASTER RECOVERY
 - READ REPLICAS – FOR PERFORMANCE

UNDERSTAND DIFFERENCES BETWEEN OLTP AND OLAP

- DB SOLUTIONS FOR OLTP
 - SEE ALL SQL & NO-SQL SOLUTIONS
- DB SOLUTIONS FOR OLAP
 - REDSHIFT
 - USED FOR BUSINESS INTELLIGENCE OR DATAWAREHOUSING

UNDERSTAND WHEN TO USE CACHING MECHANISM

- AWS OFFERS ELASTICACHE SERVICE WITH 2 OPTIONS
 - MEMCACHED
 - REDIS CACHE
- USED TO SPEED UP THE PERFORMANCE OF EXISTING DATABASES

RDS – RELATIONAL DATABASE SERVICES LAB

EXAM TIPS

- RDS RUNS ON VIRTUAL MACHINES
- YOU CANNOT LOGIN TO THE OS HOSTING RDS INSTANCES (ONLY AWS CAN)
- PATCHING OF THE OS AND DB IS AMAZON'S RESPONSIBILITY (NOT YOURS). WHEN YOU PROVISION YOUR OWN EC2 INSTANCES PATCHING IS YOUR RESPONSIBILITY.
- RDS IS NOT SERVERLESS (THERE IS ONE EXCEPTION)
- AURORA SERVERLESS DB IS THE ONLY SERVERLESS OFFERING IN RDS
- LEARN ABOUT AMAZON'S ATHENA SERVICE (WTF?)

RDS – BACKUPS, MULTI-AZ & READ REPLICAS

- RDS BACKUPS
 - AUTOMATED BACKUPS
 - ALLOW TO RECOVER DB AT ANY POINT IN TIME WITHIN A RETENTION PERIOD
 - RETENTION PERIOD COULD BE BETWEEN 1 AND 35 DAYS
 - IT TAKES A FULL DAILY SNAPSHOT AND STORE TRANSACTION LOGS DURING THE DAY
 - WHEN A RECOVERY IS NEEDED, AWS WILL TAKE THE MOST RECENT SNAPSHOT AND APPLY ALL TRANSACTION LOGS THAT OCCURRED FROM THE TIME THE SNAPSHOT WAS TAKEN
 - THIS ALLOWS A RECOVERY POINT TO THE MOST RECENT SECOND WITHIN THE RETENTION PERIOD
 - ARE ENABLED BY DEFAULT
 - BACKUP DATA IS STORED IN S3 EQUAL TO THE SIZE OF YOUR DB
 - YOU MAY EXPERIENCE SOME LATENCY DURING THE BACKUP WINDOW
 - WHEN RDS INSTANCE IS DELETED AUTO BACKUPS ARE DELETED AS WELL
 - RDS BACKUPS
 - DB SNAPSHOTS
 - THESE ARE USE-INITIATED
 - THESE ARE STORED EVEN AFTER YOU DELETE THE RDS INSTANCE
 - RESTORING RDS INSTANCES FROM BACKUPS
 - IT RESULTS IN A BRAND NEW RDS INSTANCE WITH A NEW DNS ENDPOINT
 - ENCRYPTION AT REST
 - SUPPORTED FOR MS-SQL, ORACLE, MYSQL, MARIADB, POSTGRESQL AND AURORA
 - ENCRYPTION IS DONE VIA AMAZON'S KEY MANAGEMENT SERVICE (KMS)
 - ONCE ENCRYPTION AT REST IS TURNED ON ALSO BACKUPS, READ REPLICAS AND SNAPSHOTS ARE ENCRYPTED

RDS – BACKUPS, MULTI-AZ & READ REPLICAS

- RDS MULTI-AZ
 - THIS IS FOR DISASTER RECOVERY PURPOSES ONLY (NOT FOR PERFORMANCE – SEE READ REPLICAS)
 - ALLOWS YOU TO HAVE AN EXACT COPY OF YOUR DB IN ANOTHER AZ
 - AWS HANDLES THE REPLICATION FOR YOU
 - IN THE EVENT OF PLANNED DB MAINTENANCE, DB INSTANCE FAILURE OR AZ FAILURE, AMAZON RDS WILL AUTOMATICALLY FAIL OVER TO THE STAND-BY SO THAT DB OPERATIONS CAN RESUME QUICKLY WITHOUT ADMIN INTERVENTION
 - MULTI-AZ IS AVAILABLE FOR MS-SQL, POSTGRESQL, MySQL, ORACLE, AND MARIADB (AURORA IS BY DESIGN RESILIENT TO DISASTERS)
- READ REPLICAS
 - ALLOWS YOU TO HAVE READ-ONLY REPLICAS OF YOUR PRIMARY DATABASE
 - USES ASYNCHRONOUS REPLICATION FROM PRIMARY TO THE REPLICA
 - EFFECTIVE IN SITUATIONS WHERE AN APPLICATION HAS READ-INTENSIVE WORKLOADS
 - THIS IS AVAILABLE FOR MySQL, Oracle, PostgreSQL, MariaDB and Aurora
 - USED FOR SCALING AN APPLICATION (NOT FOR DISASTER RECOVERY – DR)
 - YOU CAN HAVE UP TO 5 READ-REPLICA COPIES OF YOUR DB
 - YOU CAN HAVE READ-REPLICAS OF READ-REPLICAS (LATENCY MAY BE AN ISSUE IN THIS SCENARIO)
 - EACH RR WILL HAVE ITS OWN DNS END POINT
 - RR CAN HAVE MULTI-AZ
 - YOU CAN CREATE RR OF MULTI-AZ DATABASES
 - RR CAN BE PROMOTED TO BE A PRIMARY DB (THIS BREAKS REPLICATION)
 - RR CAN BE IN A DIFFERENT REGION

REDSHIFT – OLAP DB

FEATURES

- USED FOR BUSINESS INTELLIGENCE AND DATA WAREHOUSING
- FAST, POWERFUL, FULLY MANAGED
- PETABYTE SCALE
- AS LOW AS \$0.25 X HOUR. \$1000 PER YR PER 1TB
- CHEAPER THAN MOST OTHER DATA WAREHOUSING SOLUTIONS
- **PRICING:** CHARGES ARE BASED ON
 - COMPUTE-NODE-HOURS
 - BACKUP
 - DATA TRANSFERS WITHIN A VPC
- WHEN IN MULTI-NODE, LEADER NODE DOES NOT INCUR IN CHARGES.
- KNOW THE DIFFERENCE BETWEEN OLTP/OLAP
- CONFIGURATION
 - SINGLE-NODE
 - 160 GB STORAGE SIZE
 - MULTI-NODE
 - 1 LEADER NODE
 - N COMPUTE NODES
 - UP TO 128 COMPUTE NODES
- USES ADVANCE COMPRESSION BY COLUMN BECAUSE SIMILARITY IN THE DATA TYPE
- USES LESS SPACE THAN TRADITIONAL OLTP DBs
- USES MASSIVE PARALLELISM
- EASY TO SCALE OUT BY ADDING NODES
- BACKUPS ARE ENABLED BY DEFAULT BY 1 DAY RETENTION (UP TO 35)
 - UP TO 3 COPIES STORED IN S3
 - ASYNCHRONOUSLY REPLICATES SNAPSHTOS TO ANOTHER REGION FOR DR PURPOSES
- SECURITY
 - ENCRYPTED IN TRANSIT WITH SSL
 - ENCRYPTED AT-REST AES-256
 - KEY MANAGEMENT THROUGH
 - HSM (YOU MANAGE)
 - AWS KMS (AWS MANAGES)
- AVAILABILITY
 - 24/7 WITHIN ONE AZ (NOT MULTI AZ)
 - YOU CAN MOVE BACKUPS TO ANOTHER AZ AND RESTORE THERE FOR DR PURPOSES

AWS AURORA

- IT IS AMAZON'S PROPRIETARY SQL DATABASE COMPATIBLE WITH MYSQL AND POSTGRESQL
- 2 COPIES OF YOUR DATA IS CONTAINED IN EACH AZ WITH A MINIMUM OF 3 AZ'S (6 COPIES OF YOUR DATA)
- IT IS ONLY AVAILABLE IN REGIONS WITH 3 OR MORE AZS
- YOU CAN SHARE AURORA SNAPSHOTS WITH OTHER AWS ACCOUNTS
- OFFERS 2 TYPES OF REPLICAS
 - AURORA REPLICAS
 - MYSQL REPLICAS
- OFFERS AUTOMATED FAILOVER IS ONLY AVAILABLE WITH AURORA REPLICAS
- AUTOMATED BACKUP IS TURNED ON BY DEFAULT

ELASTICACHE

- IT IS A SERVICE TO ASSIST WITH WEB APP PERFORMANCE IMPROVEMENT
- IT COMES WITH 2 OPTIONS
 - MEMCACHED
 - VERY SIMPLE, HIGH-PERFORMING OPTION
 - LACKS SOME FEATURES
 - REDIS CACHE
 - OFFERS MORE FEATURES THAN MEMCACHED
 - COMPLEX OBJECTS
 - AVAILABLE IN MULTIPLE AZ
 - YOU CAN DO BACKUPS AND RESTORE
 - AND MUCH MORE - - -> SEE PIC

Requirement	Memcached	Redis
Simple Cache to offload DB	Yes	Yes
Ability to scale horizontally	Yes	Yes
Multi-threaded performance	Yes	No
Advanced data types	No	Yes
Ranking/Sorting data sets	No	Yes
Pub/Sub capabilities	No	Yes
Persistence	No	Yes
Multi-AZ	No	Yes
Backup & Restore Capabilities	No	Yes

AWS – ROUTE 53 – A DOMAIN NAME SERVICE (DNS 101)

- ELASTIC LOAD BALANCER SERVICE (ELBs) DO NOT HAVE PRE-DEFINED IPv4 ADDRESSES
- YOU RESOLVE TO ELBs USING A DNS NAME
- UNDERSTAND THE DIFFERENCE BETWEEN
 - CNAME
 - ALIAS RECORD
- ON THE EXAM, WHEN GIVEN THE CHOICE BETWEEN A CNAME OR AN ALIAS RECORD, ALWAYS PICK THE ALIAS RECORD (IT'D BE NICE TO KNOW WHY THO)
- UNDERSTAND COMMON DNS TYPES
 - SOA RECORDS – START OF AUTHORITY RECORD
 - NS RECORDS – NAME SERVER – INDICATES WHICH DNS SERVER IS AUTHORITATIVE FOR A DOMAIN
 - A RECORDS – MOST BASIC DNS RECORDS. POINTS A DOMAIN/SUB-DOMAIN NAME TO AN IPv4 ADDRESS
 - CNAMEs – CANONICAL NAME RECORD
 - MX RECORDS – ELECTRONIC MAIL
 - PTR RECORDS – ALLOW LOOKUP OF DNS RECORDS GIVEN AN IP ADDRESS

AWS – ROUTE 53, REGISTER A DOMAIN NAME

- YOU CAN BUY DOMAIN NAMES FROM AMAZON ROUTE 53
- IT CAN TAKE UP TO 3 DAYS TO REGISTER AND PROPAGATE THROUGH THE INTERNET

AWS – ROUTE 53, ROUTING OPTIONS

- SIMPLE ROUTING (ROUND ROBIN)
- WEIGHTED ROUTING (%-BASED ROUTING)
- LATENCY-BASED ROUTING (PING-BASED ROUTING)
- FAILOVER ROUTING (ACTIVE/PASSIVE ROUTING)
- GEOLOCATION ROUTING (REGION BASED)
- GEO-PROXIMITY ROUTING (TRAFFIC FLOW ONLY)
- MULTI-VALUE ANSWER ROUTING (SIMPLE W/ HEALTH CHECK)
- RESOURCES:
 - [HTTPS://DOCS.AWS.AMAZON.COM/ROUTE53/LATEST/DEVELOPERGUIDE/RESOURCERECORDTYPES.HTML](https://docs.aws.amazon.com/Route53/latest/DeveloperGuide/ResourceRecordTypes.html)

AWS – ROUTE 53, ROUTING OPTIONS

- SIMPLE ROUTING
 - AS THE NAME INDICATES, IT IS THE MOST BASIC ROUTING OPTION
 - YOU CAN HAVE A SINGLE DNS RECORD WITH MULTIPLE IP ADDRESSES ASSOCIATED TO IT
 - WHEN A BROWSER HITS THE DNS RECORD AND IF MULTIPLE IP ADDRESSES HAVE BEEN DEFINED, ROUTE 53 WILL RETURN A RANDOM IP ADDRESS FROM THE LIST

AWS – ROUTE 53, ROUTING OPTIONS

- WEIGHTED ROUTING POLICY
 - R53 USES A % WITH WHICH TO SPLIT THE TRAFFIC COMING IN TO A GIVEN DNS
 - THE % SPLITS HAVE TO ADD TO 100
- HEALTHCHECKS (THIS IS FOR ALL ROUTING POLICIES)
 - YOU CAN SET HEALTH CHECKS ON INDIVIDUAL RECORD SETS
 - IF A RECORD FAILS THE HEALTH CHECK, IT WILL BE REMOVED FROM R53 UNTIL IT PASSES THE HEALTHCHECK
 - YOU CAN SET SNS (SIMPLE NOTIFICATION SERVICE) ALERTS TO NOTIFY YOU THAT A HEALTH CHECK FAILED

AWS – ROUTE 53, ROUTING OPTIONS

- LATENCY-BASED ROUTING POLICY
 - ROUTE 53 DETECTS THE LATENCY FOR A SPECIFIC REQUEST
 - IT THEN ROUTES THE USER TO THE AWS THAT IS EXPERIENCING THE LEAST LATENCY FROM THE PERSPECTIVE OF THE USER
 - IF YOU ARE IN JAPAN AND HIT A DNS RECORD AND THERE IS LOW LATENCY TO SYDNEY THAN TO US-WEST, THE TRAFFIC WILL BE ROUTED TO SYDNEY

AWS – ROUTE 53, ROUTING OPTIONS

- FAIL-OVER ROUTING POLICY
 - YOU NEED TO CONFIGURE THE PASSIVE AND ACTIVE IP ADDRESSES
 - YOU ALSO NEED TO DEFINE HEALTH CHECKS FOR THE EC2 INSTANCES BEHIND YOUR IP ADDRESSES
 - ROUTE 53 DETECTS A REQUEST TO A HOSTED DOMAIN
 - ALL REQUESTS WILL GO TO THE ACTIVE IP ADDRESS AS LONG AS IT HAS A PASSING HEALTH-CHECK
 - IN CASE THE HEALTH-CHECK FAILS, THE R53 WILL RE-ROUTE ALL REQUESTS TO THE PASSIVE IP ADDRESS REGISTERED IN THE ROUTING POLICY

AWS – ROUTE 53, ROUTING OPTIONS

- GEO-LOCATION ROUTING POLICY
 - YOU DEFINE WHICH ARE THE GEOGRAPHICAL LOCATIONS THAT A WEB REQUEST SHOULD GO TO
 - IF YOU DEFINE THAT USERS IN EUROPE SHOULD GO TO IP ADDRESS “A” AND USA USERS SHOULD GO TO IP ADDRESS “B” R53 WILL FORWARD REQUESTS ACCORDING TO SETTINGS
 - **QUESTIONS:**
 - WHAT HAPPENS WHEN USERS COME IN FROM REGIONS NOT DEFINED?
 - WHAT HAPPENS WHEN ONE OF THE REGIONS FAILS HEALTH-CHECKS?

AWS – ROUTE 53, ROUTING OPTIONS

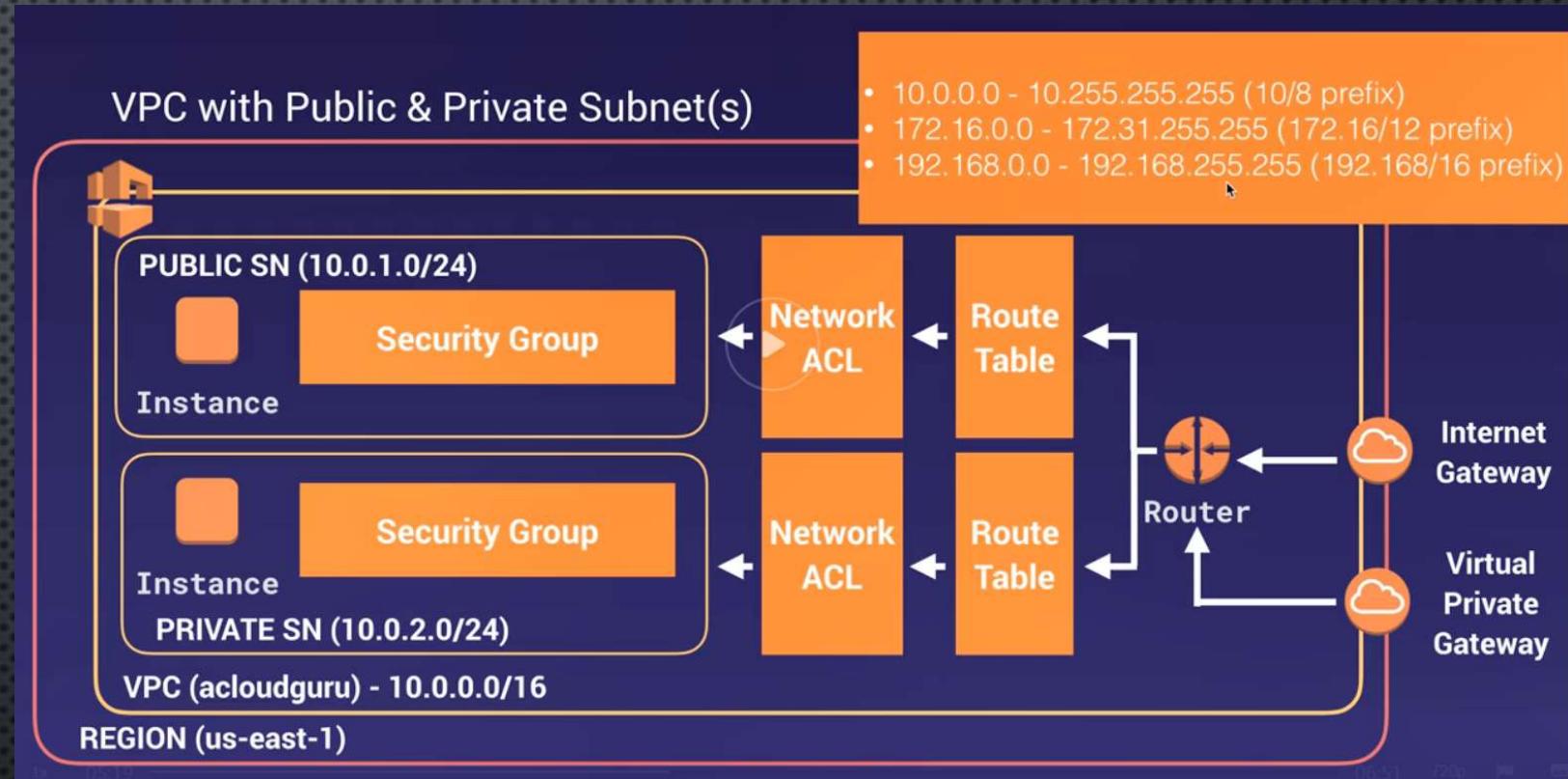
- GEO-PROXIMITY ROUTING POLICY
 - LETS TRAFFIC TO YOUR AWS RESOURCES BASED ON
 - GEOGRAPHIC LOCATION OF YOUR USER (LAT/LON)
 - GEOGRAPHIC LOCATION OF YOUR RESOURCES (LAT/LON)
 - YOU CAN ADJUST (INCREASE OR DECREASE) TRAFFIC FLOW BASED ON A BIAS
 - TO USE THIS FEATURE, YOU MUST USE R53 TRAFFIC FLOW TOOL

AWS – ROUTE 53, ROUTING OPTIONS

- MULTI-VALUE ANSWER ROUTING POLICY
 - THIS IS BASICALLY SIMPLE ROUTING BUT FOR EACH IP ADDRESS DEFINED AS AN INDEPENDENT RECORD
 - HEALTH-CHECKS NEED TO BE ENABLED FOR THIS TO WORK

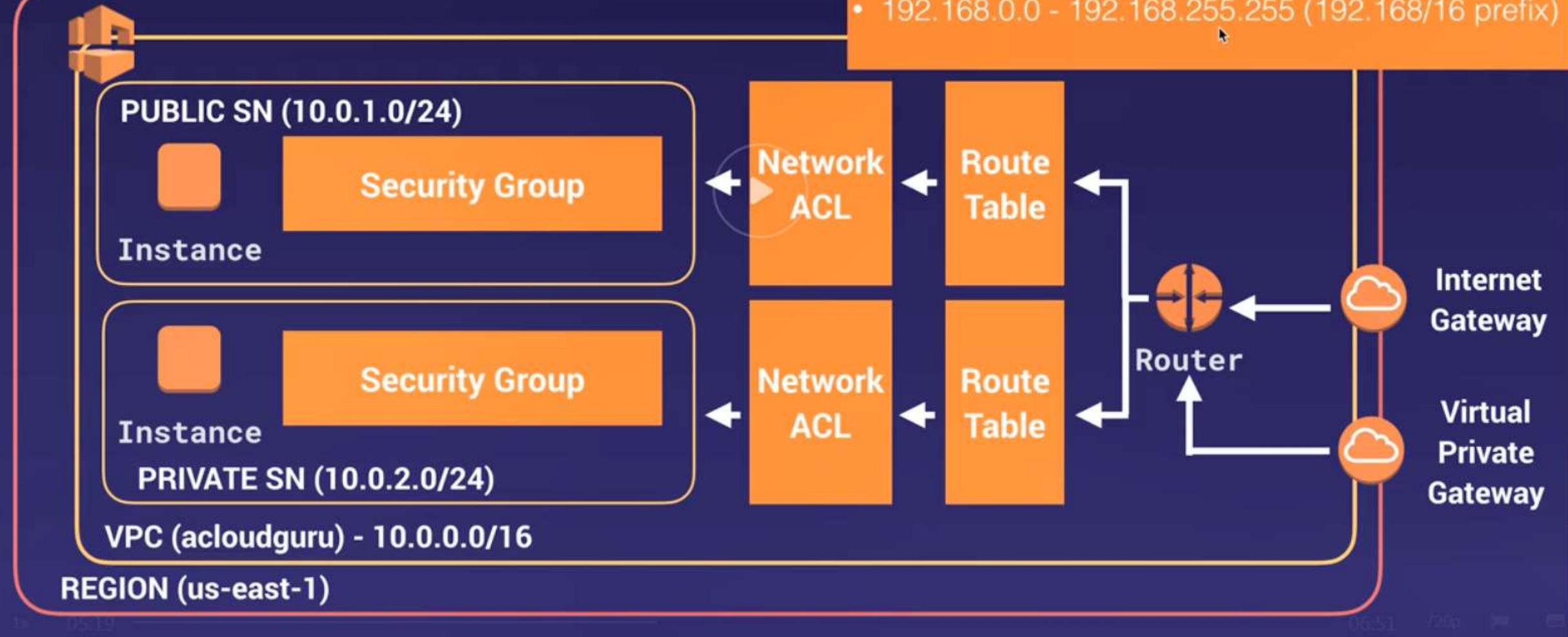
AWS - VPC

- THINK OF A VPC AS A LOGICAL DATACENTER IN AWS
- A VPC CONSISTS OF
 - INTERNET GATEWAYS & VIRTUAL PRIVATE GATEWAYS
 - ROUTE TABLES
 - NETWORK ACLs
 - SUBNETS
 - SECURITY GROUPS
- 1 SUBNET = 1 AZ, MEANING SUBNETS CANNOT SPAN OVER MULTIPLE AZS
- 1 AZ CAN HAVE MULTIPLE SUBNETS
- SECURITY GROUPS ARE STATEFUL
 - YOU CAN ADD ALLOW & DENY RULES
 - IF YOU DEFINE ALLOW OR DENY, YOU DON'T NEED TO DEFINE THE RETURN RULE
- NETWORK ACLs ARE STATELESS
 - YOU CAN ADD ALLOW & DENY RULES
 - OPENING IN-BOUND PORT DOES NOT AUTOMATICALLY ADD OUT-BOUND TRAFFIC. YOU HAVE TO ADD MANUALLY
- NO TRANSITIVE PEERING IS ALLOWED
 - GIVEN: VPC A < - > VPC B & VPC B < - > VPC D
 - THEN NO TRAFFIC CAN GO FROM VPC A TO VPC D UNLESS YOU DEFINE PEERING BETWEEN VPC A AND VPC D



VPC with Public & Private Subnet(s)

- 10.0.0.0 - 10.255.255.255 (10/8 prefix)
- 172.16.0.0 - 172.31.255.255 (172.16/12 prefix)
- 192.168.0.0 - 192.168.255.255 (192.168/16 prefix)



VPC – EXAM TIPS

- WHEN YOU CREATE A VPC, THE FOLLOWING ARE CREATED BY DEFAULT
 - DEFAULT ROUTE TABLE
 - NETWORK ACL
 - DEFAULT SECURITY GROUP
- WHEN YOU CREATE A VPC, YOU HAVE TO DEFINE
 - SUBNETS
 - INTERNET GATEWAY
- REMEMBER
 - FOR 2 DIFFERENT AWS ACCOUNTS THE SAME AZ NAME COULD REFER TO DIFFERENT, PHYSICAL DATACENTERS. THE AZS ARE RANDOMIZED WHEN ASSIGNED.
 - AMAZON ALWAYS RESERVES 5 IP ADDRESSES WITHIN YOUR SUBNETS
 - X.X.X.0 – NETWORK ADDRESS
 - X.X.X.1 – RESERVED FOR VPC ROUTER
 - X.X.X.2 – RESERVED FOR MAPPING TO THE AWS-PROVIDER DNS
 - X.X.X.3 – FUTURE USE (NOT REALLY DISCLOSED WHY)
 - X.X.X.255 – STANDARD BROADCAST ADDRESS
- YOU CAN ONLY HAVE 1 INTERNET GATEWAY PER VPC
- SECURITY GROUPS CANNOT SPAN VPCs

VPC – NAT INSTANCES VS NAT GATEWAYS

- **RESOURCE:**

<https://docs.aws.amazon.com/vpc/latest/ug/vpc-nat-instance.html>

- **NAT INSTANCE TIPS**

- THESE ARE NOT POPULAR ANY MORE
 - NOW WE HAVE **NAT GATEWAYS**
- MUST PROVIDE AN EC2 INSTANCE BE IN A PUBLIC SUBNET AND **DISABLE SOURCE/DESTINATION CHECKS OF THAT INSTANCE**
- THERE MUST BE A ROUTE (0.0.0.0/0) FROM THE PRIVATE SUBNET TO THE NAT INSTANCE FOR IT TO WORK
- THE SIZE/RESOURCES OF THE NAT INSTANCE WILL DICTATE ITS ABILITY TO SUPPORT TRAFFIC
- IF YOUR NAT INSTANCE IS BOTTLENECKING, YOU CAN INCREASE THE INSTANCE SIZE
- YOU CAN CREATE HIGH-AVAILABILITY USING AUTO SCALING GROUPS, MULTIPLE SUBNETS IN DIFFERENT AZS AND A SCRIPT TO AUTO-FAIL-OVER, BUT THIS IS VERY DIFFICULT TO SCRIPT AND MAINTAIN
- NAT INSTANCES SHOULD BE WITHIN A SECURITY GROUP (TYPICALLY A WEBDMZ SG)
- YOU ARE RESPONSIBLE FOR PATCHING THE OS OF YOUR NAT INSTANCE

- **NAT GATEWAYS TIPS**

- NAT-GW ARE REDUNDANT WITHIN AN AZ
- PREFERRED ENTERPRISE WAY TO ACCESS INSTANCES IN PRIVATE SUBNETS
- SPEED RATES START AT 5 GB/S AND SCALES UP TO 45 GB/S
- NO NEED TO PATCH THE OS
- NO NEED TO ASSOCIATE THEM TO A SECURITY GROUP
- AUTOMATICALLY GET ASSIGNED AN IP ADDRESS
- REMEMBER TO UPDATE YOUR DEFAULT VPC ROUTE TABLE AND ATTACH TRAFFIC TO YOUR NAT GW
- NO NEED TO DISABLE SOURCE/DESTINATION CHECKS
- MULTIPLE AZS CAN SHARE A SINGLE NAT GW BUT THIS CREATES A SINGLE POINT OF FAILURE CONCERN
- IF YOU HAVE RESOURCES IN PRIVATE SUBNETS IN MULTIPLE AZS THAT NEED INTERNET ACCESS, CONFIGURE A NAT GW FOR EACH OF THE AZ AND UPDATE THE DEFAULT ROUTE TABLE TO ALLOW 0.0.0.0/0 TRAFFIC

VPC – NETWORK ACL vs SECURITY GROUPS

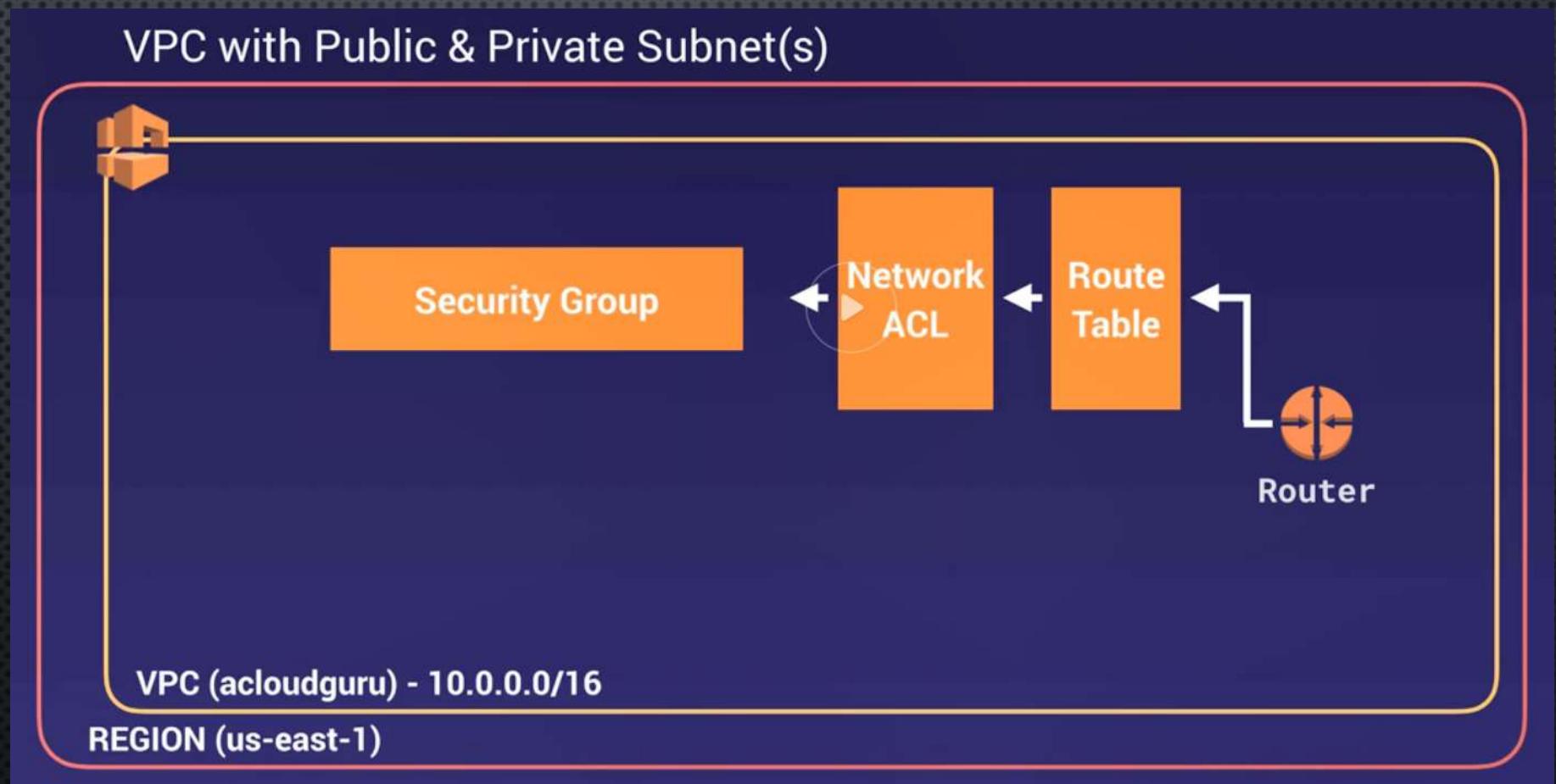
- TRAFFIC NETWORK ACL RULES ARE EVALUATED **BEFORE** TRAFFIC RULES IN SECURITY GROUPS
- UPON CREATION OF A VPC, A DEFAULT NACL IS CREATED
- A DEFAULT NACL ALLOWS ALL IN-BOUND/OUT-BOUND TRAFFIC BY DEFAULT
- YOU CAN CREATE CUSTOM NACL
- BY DEFAULT, CUSTOM NACL DENIES ALL IN-BOUND/OUT-BOUND TRAFFIC UNTIL YOU ADD SPECIFIC IN-BOUND/OUT-BOUND RULES
- EACH SUBNET IN YOUR VPC MUST BE ASSOCIATED TO A NACL
- WHEN A SUBNET IS CREATED, IT WILL BE ASSIGNED TO THE DEFAULT NACL. YOU CAN CHANGE THAT ASSOCIATION LATER TO A CUSTOM NACL
- TO BLOCK SPECIFIC IP ADDRESSES USE A NACL RULE DENY ACCESS TO THAT IP. YOU CANNOT USE SECURITY GROUPS TO BLOCK IP ADDRESSES
- A NACL CAN BE ASSOCIATED TO MULTIPLE SUBNETS, BUT A SUBNET CAN ONLY BE ASSOCIATED TO 1 NACL AT A TIME
- WHEN YOU ASSOCIATE A SUBNET TO A NACL, THE PREVIOUS NACL ASSOCIATION WILL BE REMOVED
- A NACL CONTAINS A LIST OF ALLOW/DENY RULES THAT ARE ORDERED SEQUENTIALLY. STARTING FROM THE LOWEST NUMBER, THE SEQUENCE NUMBER INDICATES THE ORDER IN WHICH THE RULE WILL BE EVALUATED BY NETWORKING HW/SW.
- NACLs HAVE SEPARATE IN-BOUND/OUT-BOUND RULES WHICH CAN EITHER ALLOW OR DENY TRAFFIC
- NACLs ARE STATELESS. RESPONSES TO ALLOWED IN-BOUND TRAFFIC ARE SUBJECT TO THE RULES FOR OUT-BOUND TRAFFIC (AND VICE VERSA)
 - **NOTE:** WITH SECURITY GROUPS (WHICH ARE STATEFUL) THAT CONSIDERATION IS NOT NEED. ALLOWED IN-BOUND TRAFFIC DOES NOT REQUIRE TO DEFINE THE ALLOWED OUT-BOUND TRAFFIC.

VPC – STEPS TO CREATE VPC WITH PUB/PRIVATE SUBNETS

- Create VPC (10.0.0.0/16 – 65K+ ip addresses)
 - Default SG
 - Default NACL
 - Default RT
- Create Internet Gateway
 - Attach IG to VPC
- Create a subnets
 - Public (10.0.1.0/24 – 256 ip addr)
 - Private (10.0.2.0/24 – 256 ip addr)
- Create custom RT for public-exposed traffic
 - In bound traffic
 - Allow ports 22, 80, and 443 (SSH, HTTP, HTTPS)
 - Out bound traffic
 - Allow all (0.0.0.0/0)

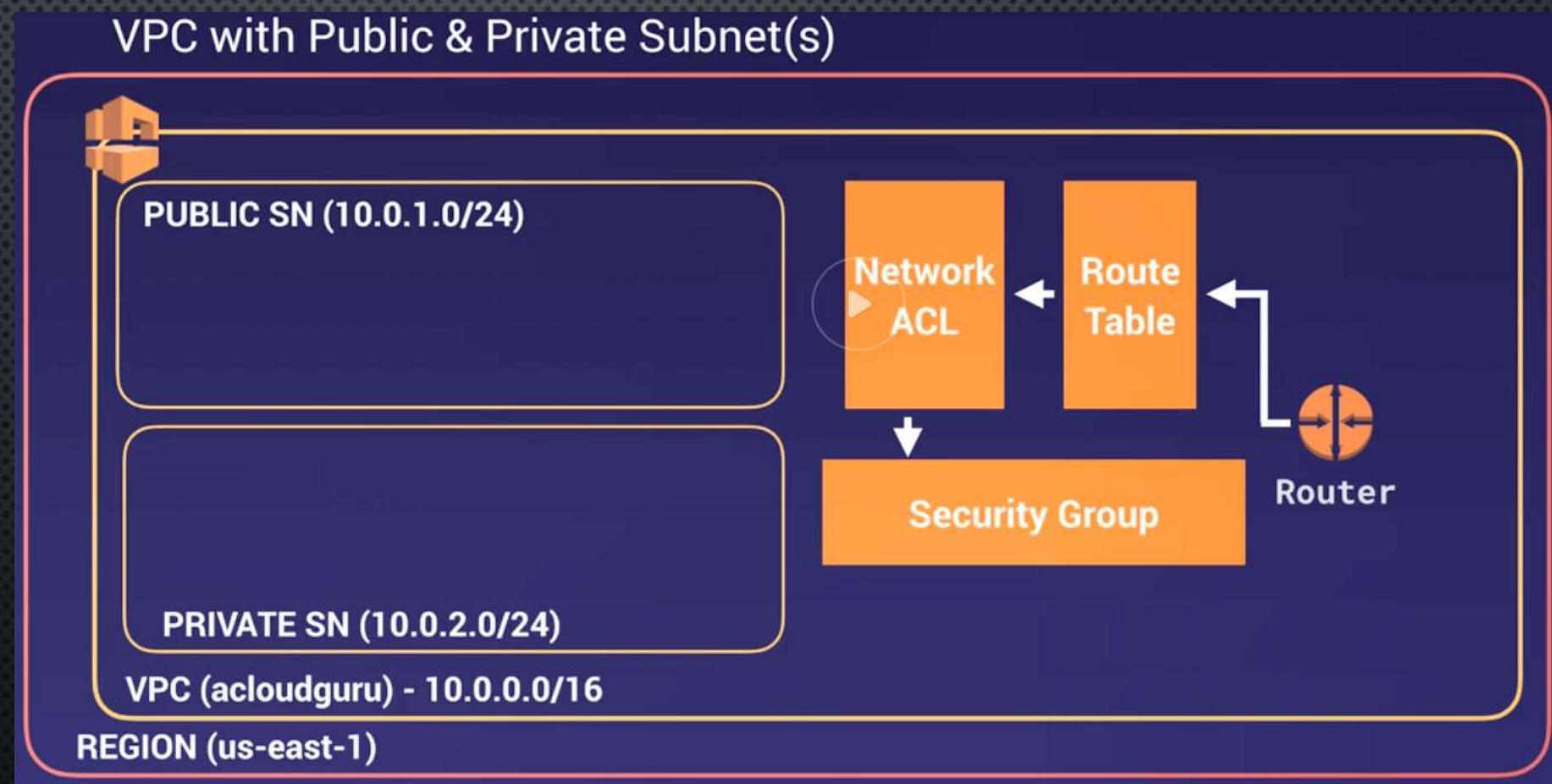
VPC – STEPS TO CREATE VPC WITH PUB/PRIVATE SUBNETS

- Select a Region
- Create VPC in that region
 - CIDR 10.0.0.0/16 for biggest address possible – 65K+ ip addresses
 - Give a name
- By default, the following are created
 - Default vpc Routing Table
 - Default vpc NACL
 - Default vpc Security Group



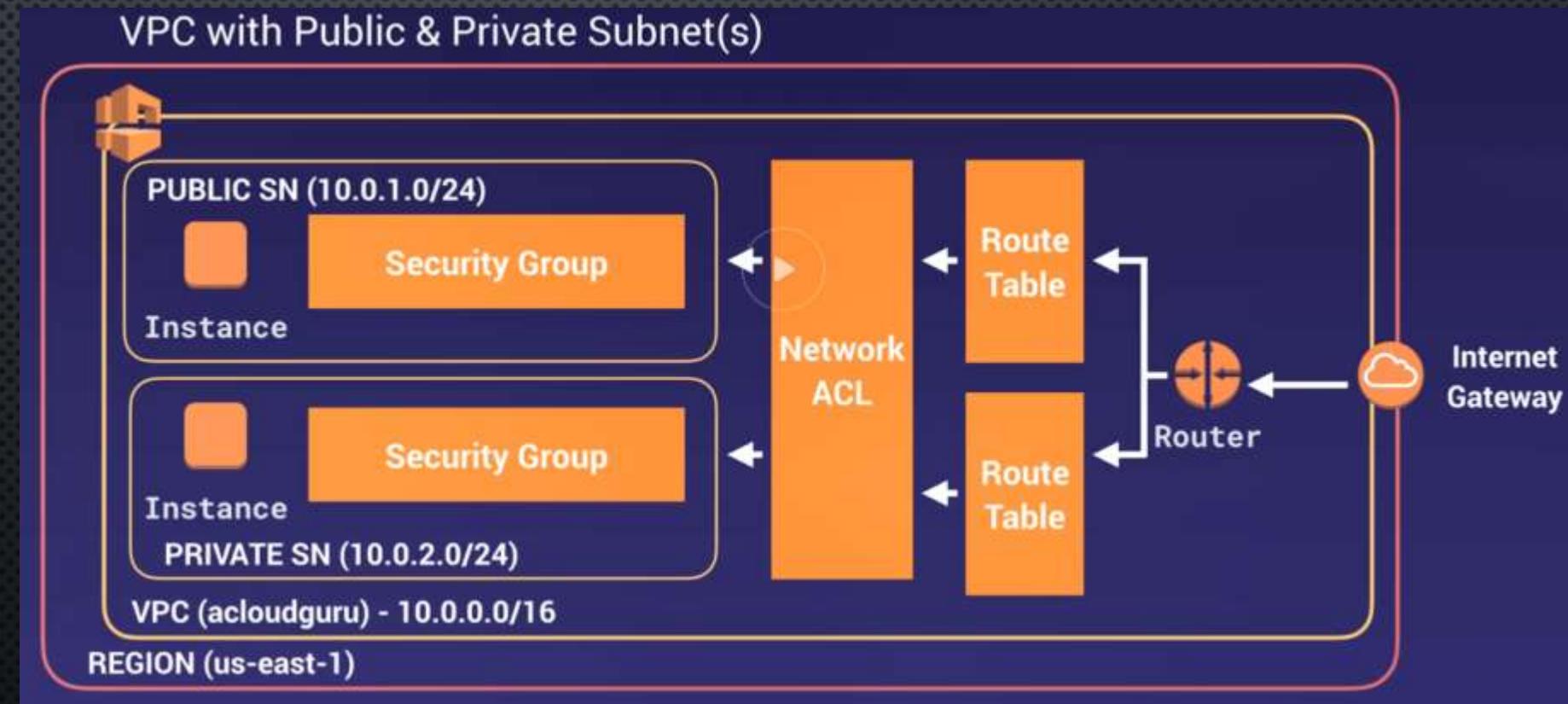
VPC – STEPS TO CREATE VPC WITH PUB/PRIVATE SUBNETS

- Create **Public** Subnet
 - CIDR **10.0.1.0/24** – 256 IP addresses
 - Amazon reserves 5 of those for their use cases
 - Select AZ for subnet
 - For name, consider using “**<CIDR> – <AZ>**”
 - Set the subnet with auto-assign public IP addresses
- Create **Private** Subnet
 - Use a different AZ
 - CIDR **10.0.2.0/24**
- Associate all subnets to custom VPC



VPC – STEPS TO CREATE VPC WITH PUB/PRIVATE SUBNETS

- Add Internet Gateway and attach it with the VPC
- Configure the Main Route table
 - Always keep it PRIVATE without access to the internet
- Create a new Public Route Table and give it access to the internet
 - Add a route for IPv4 to the Internet Gateway associated to the VPC: 0.0.0.0/0 (all traffic)
 - If needed add a route for IPv6 (::/0)
- Associate the Public subnet to this route table
- Provision an EC2 Instance within our Custom VPC with the Public subnet.
- Create a WebDMZ Security Group to firewall the publicly exposed EC2 instances
- Provision an EC2 instance within Custom VPC with the Private subnet
- Use the default Security group for the VPC



VPC – STEPS TO CREATE VPC WITH PUB/PRIVATE SUBNETS

- Create a new security group for your assets deployed in your Private subnet
- Give it a Name and Description
- Restrict in-bound traffic that only comes from Public subnet
- In-bound traffic – allow
 - Ping – All ICMP
 - HTTP - 80
 - HTTPS - 443
 - SSH - 22
 - DB – your DB's port
- Change the security group of your private EC2 instances
- For private instances you need to provide a safe route out to the Internet. For that you can use
 - NAT Instances (EC2 instance)
 - NAT Gateways – a more modern way to provide access to internet resources

VPC FLOW LOGS

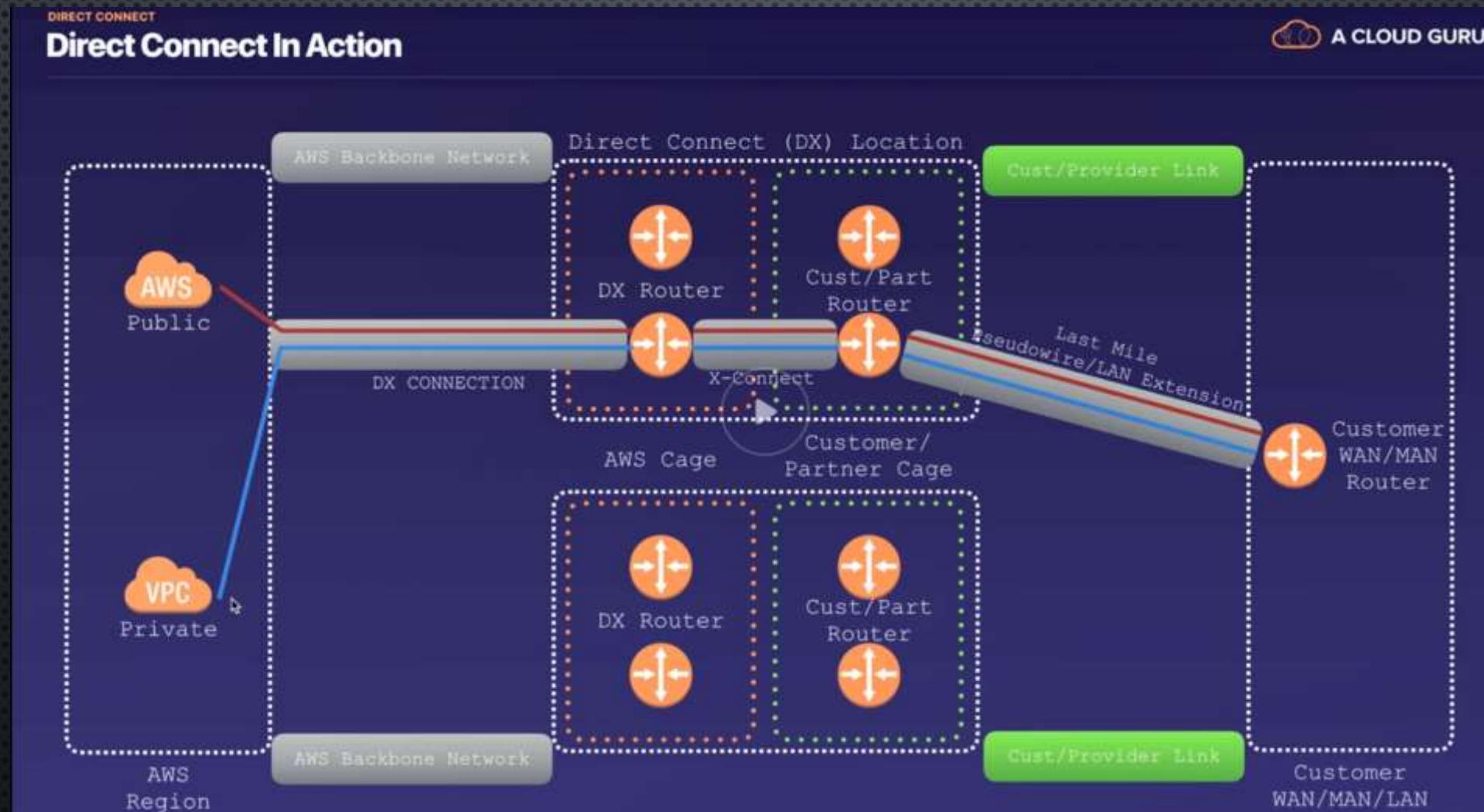
- Flow Logs:
 - ALLOW YOU TO CAPTURE INFORMATION ABOUT IP TRAFFIC GOING THROUGH NETWORK INTERFACES IN YOUR CUSTOM VPC
 - FLOW LOG DATA IS STORED IN CLOUD WATCH OR IN S3 BUCKET
 - WHEN USING CLOUD WATCH YOU CAN VIEW LOGS THERE
 - CAN BE CREATED AT THE LEVEL OF: **VPC > SUBNET > NETWORK INTERFACE** (ELASTIC NETWORK INTERFACE - ENI)
- YOU CANNOT ENABLE FLOW LOGS FOR VPCs THAT ARE PEERED WITH YOUR VPC UNLESS THE PEER-VPC IS IN YOUR ACCOUNT.
- YOU CANNOT ADD TAGS FOR A FLOW LOG (NAME IT WISELY)
- AFTER YOU HAVE CREATED A FLOW LOG, YOU CANNOT CHANGE ITS CONFIGURATION.
 - FOR EXAMPLE: YOU CANNOT ASSOCIATE A DIFFERENT IAM ROLE FOR A FLOW LOG AFTER ITS CREATION
- NOT ALL IP TRAFFIC IS MONITORED
 - TRAFFIC GENERATED BY INSTANCES WHEN THEY CONTACT AMAZON'S DNS SERVER. HOWEVER, IF YOU USE YOUR OWN DNS SERVER, THEN ALL TRAFFIC TO THAT DNS SERVER IS LOGGED.
 - TRAFFIC GENERATED BY A WINDOWS INSTANCE FOR AMAZON WINDOWS LICENSE ACTIVATION
 - TRAFFIC TO/FROM 169.254.169.254 FOR INSTANCE META DATA
 - DHCP TRAFFIC
 - TRAFFIC TO THE RESERVED IP ADDRESS FOR THE DEFAULT VPC ROUTER

VPC – BASTION HOSTS

- IS AN EC2 INSTANCE THAT LIVES IN A PUBLIC SUB NET IN CUSTOM VPC TO ALLOW REMOTE ACCESS TO RESOURCES IN PRIVATE SUBNETS
- A NAT GATEWAY OR NAT INSTANCE IS USED TO PROVIDE INTERNET TRAFFIC TO EC2 INSTANCES IN PRIVATE SUBNETS
- A BASTION IS USED TO SECURELY ADMINISTER EC2 INSTANCES, VIA SSH OR RDP (WINDOWS MACHINES). BASTIONS ARE ALSO KNOWN AS JUMP-BOXES ELSEWHERE.
- YOU CANNOT USE A NAT GATEWAY AS A BASTION HOST

DIRECT CONNECT

- DIRECT CONNECTION FROM YOUR DATA CENTER TO AWS
- USEFUL FOR HIGH-THROUGHPUT WORKLOADS (LOTS OF NETWORK TRAFFIC)
- USEFUL FOR ESTABLISHING A STABLE, RELIABLE, AND SECURE CONNECTION TO AWS



VPC ENDPOINTS

- ENABLES YOU TO PRIVATELY CONNECT YOUR VPC TO
 - AWS SUPPORTED SERVICES
 - VPC ENDPOINT SERVICES POWERED BY PRIVATE LINK
- IT DOES NOT REQUIRE YOU TO HAVE CONNECTIVITY VIA
 - INTERNET GATEWAY
 - NAT DEVICE / NAT GATEWAY
 - VPN CONNECTION
 - AWS DIRECT CONNECT
- EC2 INSTANCES DO NOT REQUIRE A PUBLIC IP ADDRESS
- TRAFFIC BETWEEN YOUR VPC AND AWS SERVICES DO NOT HAVE TO LEAVE THE AWS NETWORK
 - EXAMPLE – YOU HAVE AN EC2 INSTANCE THAT NEEDS TO DROP OBJECTS INTO AN S3 BUCKET
- VPC ENDPOINTS ARE VIRTUAL DEVICES
 - SCALE HORIZONTALLY
 - ARE REDUNDANT
 - HIGHLY AVAILABLE
 - DO NOT IMPOSE
 - AVAILABILITY RISKS
 - BANDWIDTH CONSTRAINTS ON YOUR NETWORK TRAFFIC
- TYPES OF VPC ENDPOINTS
 - INTERFACE ENDPOINTS
 - GATEWAY ENDPOINTS
- CURRENTLY GW ENDPOINTS SUPPORT
 - AMAZON S3
 - DYNAMO DB

ELASTIC LOAD BALANCERS (ELBs)

- APPLICATION LOAD BALANCER (LAYER 7)
 - HTTP/HTTPS
 - INTELLIGENT
 - ADVANCE / INTELLIGENT ROUTING
- NETWORK LOAD BALANCER (LAYER 4)
 - TCP TRAFFIC BALANCER
 - USE IF YOU WANT
 - MAXIMUM PERFORMANCE
 - FIXED IP ADDRESS
 - MILLIONS OF REQUESTS PER SECOND
 - HIGHEST COST
- CLASSIC LOAD BALANCER
 - HTTP/HTTPS
 - NOT APP AWARE
 - STICKY SESSION (X-FORWARDED-FOR HEADER, IPv4)
 - LOWEST COST
- ERRORS
 - 504 ERRORS – GATEWAY ERROR ON THE APP AT THE WEB SERVER OR BACK-END, DB
- INSTANCES MONITORED BY ELBS ARE REPORTED AS
 - IN SERVICE
 - OUT OF SERVICE
- HEALTH CHECKS ENSURE THAT AN INSTANCE UNDER ELB IS HEALTHY BY HITTING AN HTTP PATH TO THAT INSTANCE
- LOAD BALANCERS HAVE THEIR OWN DNS NAME
- YOU WILL NEVER GET AN IP ADDRESS FOR A LOAD BALANCER
- SLAS GUARANTEE (ALL 3 TYPES OF ELBs)
 - MONTHLY AVAILABILITY OF 99.99%
- READ THE ELB FAQ:
 - <https://aws.amazon.com/elasticloadbalancing/faqs/>

ADVANCED LOAD BALANCER THEORY

- STICKY SESSIONS
 - USERS CAN STICK TO THE SAME EC2
 - USEFUL WHEN THERE IS APPLICATION STATE FOR GIVEN USER STORED IN THE EC2
- CROSS ZONE LB
 - ALLOWS AN ELB TO FORWARD TRAFFIC TO A DIFFERENT AZ OTHER THAN THE ONE IT HAS BEEN DEFINED FOR
- PATH PATTERNS
 - ALLOWS AN ELB TO FORWARD TRAFFIC TO A DIFFERENT AZ BASED ON THE URL PATH COMING IN THE HTTP REQUEST

HAA – HIGH AVAILABILITY ARCHITECTURE

- ALWAYS DESIGN FOR FAILURE
- USE MULTIPLE AZ IN MULTIPLE REGIONS WHENEVER YOU CAN
- KNOW THE DIFFERENCE BETWEEN
 - RDS MULTI AZ
 - RDS READ REPLICAS
- KNOW THE DIFFERENCE BETWEEN
 - SCALING OUT – PROCURING MORE SIMILAR INSTANCES (LAUNCHING MORE T2 MICRO EC2's)
 - SCALING UP – INCREASING THE RESOURCES OF AN INSTANCE (FROM T2 – > XL)
 - S3 STORAGE CLASSES
 - HA S3 – STANDARD S3, INFREQUENT ACCESS S3
 - NOT HA S3 – INFREQUENT ACCESS 1 ZONE OR REDUCED REDUNDANCY STORAGE ([RRS](#))
- FOR THE EXAM, READ QUESTIONS CAREFULLY AND ALWAYS CONSIDER THE COST FACTOR

HA – EXAM TIPS

- LOAD BALANCER TYPES
 - APPLICATION LB
 - LAYER-7 AWARE. THESE LBs KNOW MORE INTIMATE DETAILS ABOUT THE APPLICATION METADATA
 - WITH THESE LB's YOU GET A DNS / URL – NEVER AN IP ADDRESS
 - NETWORK LB
 - EXTREME PERFORMANCE AT THE HARDWARE LEVEL (LAYER-4 AWARE)
 - CAN HARD-CODE IP ADDRESS
 - HIGHEST COST
 - CLASSIC LB
 - LOWEST COST
 - NO “INTELLIGENT” ROUTING
 - X-FORWARDED-FOR HEADER TELLS YOU THE IP ADDRESS OF THE CLIENT CALLING THE APP
 - WITH THESE LB's YOU GET A DNS / URL – NEVER AN IP ADDRESS
 - ELB REPORTING
 - ERRORS
 - 504 ERROR CODE – GATEWAY TIMEOUT
 - ERROR IS NOT WITH THE LB. IT IS PROBABLY WITH THE APP OR THE BACKEND (APP SERVER OR DB)
 - EC2's MONITORED BY ELB PRESENT EITHER
 - IN-SERVICE
 - OUT-OF-SERVICE
 - YOU CAN TRACK INSTANCE'S HEALTH WITH HEALTH CHECKS (GETTING A 200 HTTP CODE AFTER CALLING GET HTTP VERB)
 - READ ELB FAQs
- LB THEORY
 - STICKY SESSIONS
 - ALLOWS USERS TO “STICK” TO THE SAME INSTANCE DURING THEIR SESSION
 - CROSS-ZONE LB
 - ENABLES YOU TO LOAD-BALANCE TRAFFIC BETWEEN AVAILABILITY ZONES
 - PATH PATTERNS
 - ENABLES YOU TO LOAD-BALANCE TRAFFIC BETWEEN DIFFERENT EC2 INSTANCES BASED ON URL PATTERNS
 - CLOUDFORMATION
 - A WAY TO SCRIPT DEFINITION OF A CLOUD ENVIRONMENT – LOTS OF GRANULAR CONTROL: VPC, SUBNETS, NACL, EC2, ETC.
 - QUICKSTART TEMPLATES – PREDEFINED CLOUDFORMATION CLOUD ENVIRONMENTS TARGETING SPECIFIC TECH STACKS
 - TARGETED TO DEVOPS TEAMS WHO WANT FULL CONTROL OF THEIR ENVIRONMENTS
 - ELASTIC BEANSTALK
 - DESIGNED TO QUICKLY DEPLOY AND MANAGE APPLICATIONS WITHOUT WORRYING ABOUT INFRA SETTINGS
 - USER DEPLOYS THE APP AND EBS TAKES IT FROM THERE (SIMILAR TO CLOUD FOUNDRY)
 - TARGETED TO DEVELOPERS WHO DO NOT HAVE/WANT/NEED TO DEAL WITH INFRASTRUCTURE SETTINGS

AWS SQS – SIMPLE QUEUE SERVICE

- THE FIRST AWS SERVICE – HAS LOTS OF LEGACY
 - MESSAGING / QUEUING PLATFORM
 - DECOUPLES COMPONENTS OF AN APPLICATION
 - MESSAGES ARE STORED IN A FAIL-SAFE QUEUE
 - SQS IS “PULL” BASED NOT “PUSH” BASED. A CALLER NEEDS TO ACCESS THE QUEUE TO EITHER ADD OR GET MESSAGES FROM THE QUEUE
 - MESSAGES CAN ONLY BE UP TO 256 KB OF TEXT IN ANY FORMAT
 - MESSAGES CAN BE BIGGER BUT IT USES S3 TO STORE THOSE – UP TO 2 GB
 - MESSAGES CAN BE RETRIEVED PROGRAMMATICALLY USING THE SQS API
 - AUTO-SCALING CAN BE TRIGGERED BY NUMBER OF MESSAGES IN THE QUEUE
 - MESSAGES CAN BE KEPT FROM 1 MIN – 14 DAYS – DEFAULT IS 4 DAYS
 - MSGS ARE SUBJECT TO **VISIBILITY TIMEOUT** – A PERIOD OF TIME (12 HOURS MAX) THAT A MSG THAT IS ACCESSED BY AN EC2 IS INVISIBLE TO ANY OTHER EC2 WHICH IS GETTING MESSAGES – THIS COULD RESULT IN THE SAME MESSAGE BEING DELIVERED TWICE
 - LONG POLLING – A REQUEST FOR A MESSAGE DOES NOT RETURN RESPONSE IF Q IS EMPTY OR LONG POLL TIMES OUT
 - SHORT POLLING – A REQUEST FOR MSG RETURNS IMMEDIATELY EVEN IF Q IS EMPTY
- **STANDARD QUEUE TYPE**
 - IS THE DEFAULT TYPE
 - NEARLY UNLIMITED NUMBER OF TPS
 - GUARANTEE THAT MESSAGES ARE DELIVERED AT LEAST ONCE
 - PROVIDE “BEST-EFFORT” ORDERING: MESSAGES ARE GENERALLY DELIVERED IN THE ORDER THEY WERE SENT
 - HOWEVER, MULTIPLE COPIES OF A MESSAGE MAY BE DELIVERED OUT OF ORDER MORE THAN ONCE
 - **FIFO QUEUE TYPE**
 - FIRST-IN FIRST-OUT
 - UP TO 300 TPS
 - COMPLEMENTS THE STANDARD QUEUE
 - ORDER OF MESSAGES IS STRICTLY PRESERVED
 - MESSAGE DELIVERY IS GUARANTEED TO BE ONLY ONCE AND REMAINS AVAILABLE UNTIL A CONSUMER DELETES IT
 - DUPLICATES ARE NOT INTRODUCED INTO THE QUEUE

AWS SIMPLE WORKFLOW SERVICE (SWF)

- KNOW THE DIFFERENCE BETWEEN SWF vs SQS
- RETENTION
 - SQS HAS A RETENTION PERIOD OF 14 DAYS
 - SWF – WORKFLOW EXECUTIONS CAN PERSIST FOR UP TO 1 YEAR
- ORIENTATION OF USE CASE
 - SQS – OFFERS A “MESSAGE-ORIENTED” API SOLUTION
 - SWF – OFFERS A “TASK-ORIENTED” API SOLUTION
- DUPLICATION
 - SQS – HAS FLEXIBILITY OF HANDLING DUPLICATED MESSAGES OR ENFORCING FIFO
 - SWF – ENSURES THAT TASKS ARE ASSIGNED ONLY ONCE AND ARE NEVER DUPLICATED
- TRACKING
 - SQS – YOUR APPLICATION NEEDS TO IMPLEMENT APP-LEVEL TRACKING OF MESSAGES
 - SWF – KEEPS TRACK OF ALL TASKS AND EVENTS IN AN APPLICATION
- ACTORS (APPLIES ONLY TO SWF)
 - WORKFLOW STARTERS – APPLICATIONS THAT CAN INITIATE AN SWF WORKFLOW
 - DECIDERS – DECIDE WHAT TO DO NEXT AFTER A TASK FINISHES IN A WORKFLOW, HANDLE ERROR STATE
 - ACTIVITY WORKERS – CARRY OUT LOGIC WORK WITHIN A STEP IN THE WORKFLOW

AWS – SIMPLE NOTIFICATION SERVICE (SNS)

- SERVICE THAT ALLOWS SETTING UP, OPERATE AND SEND MESSAGES FROM THE CLOUD
- USABLE FOR PUSH NOTIFICATIONS FOR MOBILE APPLICATIONS
- ABLE TO DELIVER MESSAGES BY SQS, SMS TEXT OR EMAIL
- USES TOPICS TO ASSOCIATE PUBLISHERS AND SUBSCRIBERS FOR MESSAGE DISTRIBUTION
- MESSAGES PUBLISHED TO SNS ARE STORED REDUNDANTLY ACROSS MULTIPLE AZS
- BENEFITS
 - INSTANTANEOUS PUSH-BASED DELIVERY
 - SIMPLE API AND SIMPLE INTEGRATION WITH APPS
 - FLEXIBLE MESSAGE DELIVERY OVER MULTIPLE TRANSPORT PROTOCOLS
 - INEXPENSIVE, PAY-AS-YOU-GO MODEL. NO UPFRONT COSTS
 - AWS CONSOLE OFFERS SIMPLE UI TO MANAGE
- SNS vs SQS
 - BOTH ARE MESSAGING SYSTEMS
 - SNS IS PUSH-BASED. SQS IS PULL-BASED

AWS – ELASTIC TRANSCODER

- IS A CLOUD-BASED, MEDIA TRANSCODER
- CONVERTS MEDIA FILES FROM THEIR ORIGINAL SOURCE FORMAT INTO DIFFERENT FORMATS THAT WILL PLAY ON DIFFERENT DEVICES (TABLETS, SMART PHONES, PCs, ETC)

AWS – API GATEWAY

- API GW – IS A “DOOR” FROM THE INTERNET TO YOUR AWS ENVIRONMENT (DIFFERENT THAN INTERNET GW TO EC2 INSTANCES)
- IT HAS CACHING CAPABILITIES TO INCREASE PERFORMANCE (NEED TO SET A TTL IN SECONDS)
- IT HAS A LOW COST
- SCALES AUTOMATICALLY
- IT CAN BE THROTTLED TO PREVENT ATTACKS
- RESULTS CAN BE LOGGED TO CLOUDWATCH
- IF YOUR APPLICATION USES MULTIPLE DOMAINS WITH API GW, ENSURE THAT CORS IS ENABLED
- REMEMBER THAT CORS IS ENFORCED BY YOUR BROWSER TO REDUCE VULNERABILITIES AND INJECTION ATTACKS OF MALICIOUS CODE

AWS – KINESIS

- IT IS AN AWS SERVICE/PLATFORM FOR PROCESSING DATA STREAMS
- COMES IN 3 TYPES: STREAMS, FIREHOSE, ANALYTICS
- K-STREAMS
 - DATA PRODUCERS CAN SEND DATA TO KINESIS
 - DATA HAS PERSISTENCE AND IS STORED 24 HOURS (DEFAULT) UP TO 7 DAYS
 - DATA IS STORED IN “SHARDS” WHICH ARE LOGICAL GROUPINGS THAT DATA CONSUMERS (EC2) INSTANCES CAN ACCESS, PROCESS AND STORE THE RESULTS OF THE WORK
- PERFORMANCE OF SHARDS
 - 5 TPS FOR READ RATE UP TO 2 MB MAXIMUM READ RATE PER SECOND
 - 1000 TPS FOR WRITES UP TO 1 MB MAXIMUM WRITE RATE PER SECOND (INCLUDING PARTITION KEYS)
 - DATA CAPACITY FOR A STREAM DEPENDS ON THE NUMBER OF SHARDS THAT ARE PROVISIONED FOR IT
- K-FIREHOSE
 - DATA IS NOT PERSISTED. IT HAS TO BE PROCESSED AS IT IS RECEIVED
 - LAMBDA FUNCTIONS ARE GOOD PROCESSORS OF IN-COMING DATA
- K-ANALYTICS
 - WORKS WITH K-STREAMS OR K-FIREHOSE AND ANALYZES THE DATA ON THE FLY
 - AFTER ANALYSIS, IT STORES THE DATA IN S3, REDSHIFT OR ELASTICSEARCH CLUSTER

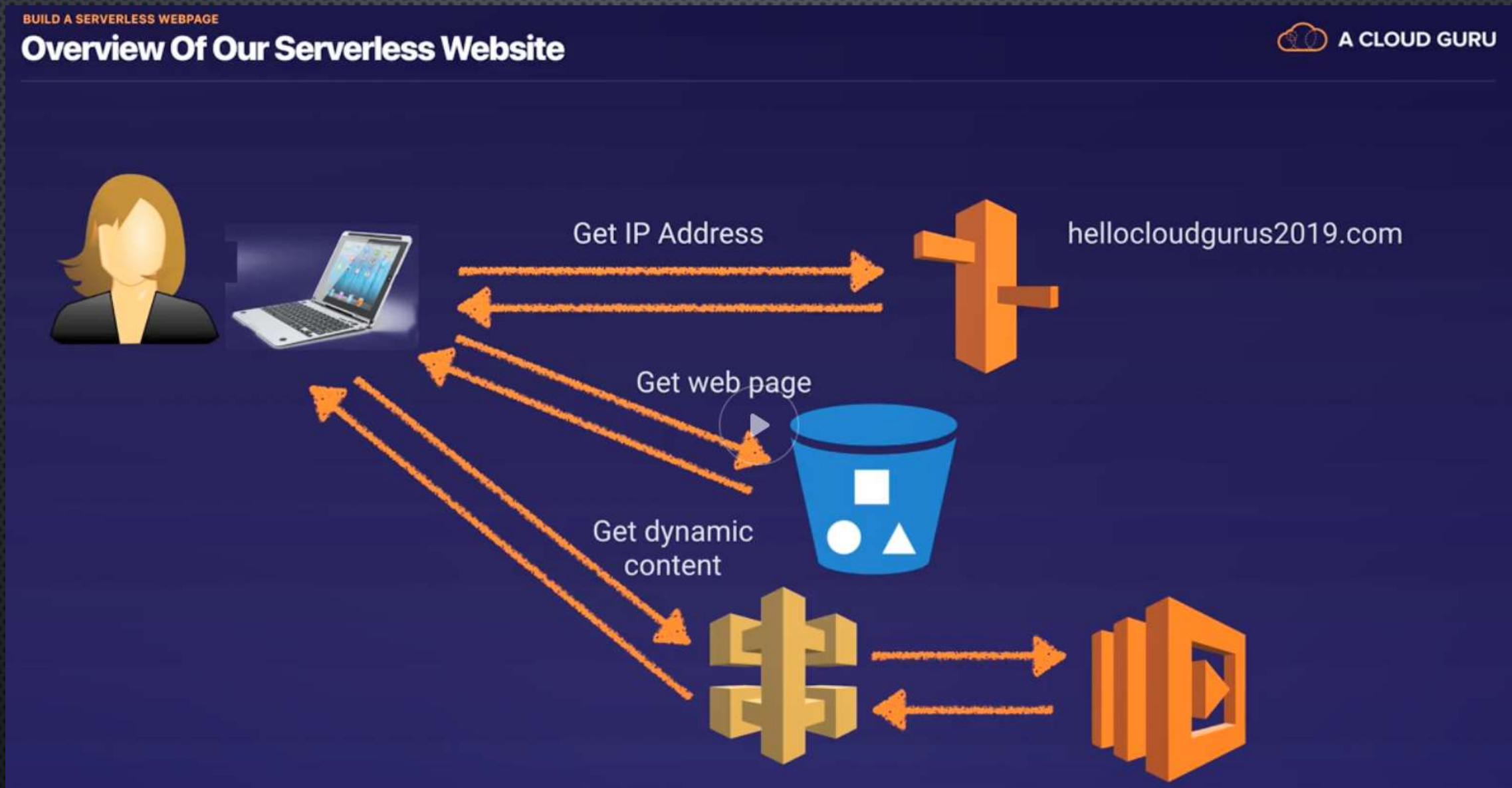
IAM – COGNITO

- AWS ALLOWS YOU TO LEVERAGE WEB IDENTITY FEDERATION SERVICES (I.E. GOOGLE ID, FACEBOOK ID, AMAZON ID, ETC) TO PROVIDE TEMPORARY ACCESS TO AWS SERVICES FOR AUTHENTICATED USERS VIA AMAZON COGNITO
- ALLOWS SIGN-UP AND SIGN-IN TO APPS
 - IF NOT USING FEDERATED IDs, THEN THE APP CAN USE COGNITO USER POOL
- ALLOWS ACCESS FOR “GUEST” USERS
- ACTS AS AN ID BROKER BETWEEN WEB-ID PROVIDERS AND AWS
- SUCCESSFUL AUTHENTICATION RESULTS IN THE CREATION OF JWT TO GRANT ACCESS TO SERVICES
- SYNCHRONIZES USER DATA FOR MULTIPLE DEVICES
- RECOMMENDED FOR ALL MOBILE APPLICATIONS THAT ACCESS AWS SERVICES
- IDENTITY POOLS
 - THESE DEAL WITH AUTHORIZATION OF USERS
 - PROVIDE TEMPORARY ACCESS TO AWS RESOURCES

AWS – LAMBDA FUNCTIONS

- LAMBDA SCALES OUT (NOT UP) AUTOMATICALLY
- LAMBDA FUNCTIONS ARE INDEPENDENT, 1 EVENT = 1 FUNCTION
- LAMBDA IS SERVER-LESS. UNDERSTAND WHAT AWS SERVICES ARE SERVER-LESS:
 - AURORA DB
 - DYNAMO DB
 - API GATEWAY
 - S3
 - LAMBDA
- LAMBDA FUNCTIONS CAN TRIGGER OTHER LAMBDA FUNCTIONS
- DRAW-BACKS/TRADE OFFS
 - SOLUTIONS CAN GET REALLY COMPLEX WITH LAMBDA
 - NETWORKS OF LAMBDA FUNCTIONS CAN BE DIFFICULT TO DEBUG
 - USE AWS X-RAY SERVICE TO DEBUG SERVER-LESS APPLICATIONS
- LAMBDA FUNCTIONS CAN HELP YOU TRIGGER ACTIONS GLOBALLY ACROSS AWS' GLOBAL INFRASTRUCTURE (I.E. BACK-UP AN S3 BUCKET TO ANOTHER ONE)
- UNDERSTAND WHAT ARE THE DIFFERENT TRIGGERS FOR LAMBDA FUNCTIONS

LAMBDA SERVERLESS WEB PAGE EXAMPLE



SERVER-LESS SUMMARY

- REMEMBER THE DIFFERENCE BETWEEN
 - TRADITIONAL ARCHITECTURE
 - MANUAL CONSTRUCTION OR AUTOMATIC (VIA AWS CLOUDFORMATION OR AWS ELASTIC BEANSTALK)
 - HAS LIMITATIONS FOR SCALABILITY
 - COULD BE COSTLY AS YOU SCALE OUT
 - SERVER-LESS ARCHITECTURE
 - HIGHLY SCALABLE
 - HIGHLY AVAILABLE
 - LOWEST COST
- LAMBDA FUNCTION BENEFITS
 - SCALE OUT NOT UP
 - ARE INDEPENDENT... 1 EVENT = 1 FUNCTION
 - FUNCTIONS ARE SERVERLESS
 - KNOW WHAT SERVICES TRIGGER (OR NOT) LAMBDA FUNCTIONS
 - LAMBDA FUNCTIONS CAN TRIGGER OTHER FUNCTIONS
 - CAN DO THINGS GLOBALLY WITHIN THE AWS INFRASTRUCTURE AND SERVICES (I.E. USE IT TO BACK UP S3 BUCKETS)
- LAMBDA FUNCTION DRAWBACKS/TRADEOFFS
 - ARCHITECTURES CAN GET REALLY COMPLEX AND HARD TO DEBUG (MITIGATION: AWS X-RAY SERVICE)