# Most Important Charts & Diagrams

## A Solutions Architect Associate Study Guide

## AWS Support Plans Chart

Via https://aws.amazon.com/premiumsupport/plans/

Basic Support is included for all AWS customers and includes:

**Customer Service & Communities** - 24x7 access to customer service, documentation, whitepapers, and support forums.

**AWS Trusted Advisor** - Access to the 7 core Trusted Advisor checks and guidance to provision your resources following best practices to increase performance and improve security.

**AWS Personal Health Dashboard** - A personalized view of the health of AWS services, and alerts when your resources are impacted.

| | **Developer** | **Business** | **Enterprise** |
|---|---|---|---|
| | *Recommended if you are experimenting or testing in AWS.* | *Recommended if you have production workloads in AWS.* | *Recommended if you have business and/or mission critical workloads in AWS.* |
| **AWS Trusted Advisor Best Practice Checks** | 7 Core checks | Full set of checks | Full set of checks |
| **Enhanced Technical Support** | Business hours** email access to Cloud Support Associates | 24x7 phone, email, and chat access to Cloud Support Engineers | 24x7 phone, email, and chat access to Cloud Support Engineers |
| | Unlimited cases / 1 primary contact | Unlimited cases / unlimited contacts (IAM supported) | Unlimited cases / unlimited contacts (IAM supported) |
| **Case Severity / Response Times*** | General guidance: < 24 business hours**  System impaired: < 12 business hours** | General guidance: < 24 hours  System impaired: < 12 hours  Production system impaired: < 4 hours  Production system down: < 1 hour | General guidance: < 24 hours  System impaired: < 12 hours  Production system impaired: < 4 hours  Production system down: < 1 hour  Business-critical system down: < 15 minutes |
| **Architectural Guidance** | General | Contextual to your use-cases | Consultative review and guidance based on your applications |
| **Programmatic Case Management** | | AWS Support API | AWS Support API |

| | | | |
|---|---|---|---|
| **Third-Party Software Support** | | Interoperability & configuration guidance and troubleshooting | Interoperability & configuration guidance and troubleshooting Infrastructure Event Management |
| **Proactive Programs** | | Access to Infrastructure Event Management for additional fee. | Well-Architected Reviews<br><br>Operations Reviews |
| **Technical Account Management** | | | Technical Account Manager (TAM) coordinates access to programs and other AWS experts as needed. Designated Technical Account Manager (TAM) to proactively monitor your environment and assist with optimization. |
| **Training** | | | Access to online self-paced labs |
| **Account Assistance** | | | Concierge Support Team |
| **Pricing** | Greater of $29 / month***<br><br>- or -<br><br>3% of monthly AWS usage<br><br>See pricing detail and example. | Greater of $100 / month***<br><br>- or -<br><br>10% of monthly AWS usage for the first $0–$10K<br><br>7% of monthly AWS usage from $10K–$80K<br><br>5% of monthly AWS usage from $80K–$250K<br><br>3% of monthly AWS usage over $250K<br><br>See pricing detail and example. | Greater of $15,000<br><br>- or -<br><br>10% of monthly AWS usage for the first $0–$150K<br><br>7% of monthly AWS usage from $150K–$500K<br><br>5% of monthly AWS usage from $500K–$1M<br><br>3% of monthly AWS usage over $1M<br><br>See pricing detail and example. |

S3 Storage Classes Chart

Via https://aws.amazon.com/s3/storage-classes/

# **Performance across the S3 Storage Classes**

| | S3 Standard | S3 Intelligent-Tiering* | S3 Standard-IA | S3 One Zone-IA† | S3 Glacier | S3 Glacier Deep Archive |
|---|---|---|---|---|---|---|
| Designed for durability | 99.999999999% (11 9's) | 99.999999999% (11 9's) | 99.999999999% (11 9's) | 99.999999999% (11 9's) | 99.999999999% (11 9's) | 99.999999999% (11 9's) |
| Designed for availability | 99.99% | 99.9% | 99.9% | 99.5% | 99.99% | 99.99% |
| Availability SLA | 99.9% | 99% | 99% | 99% | 99.9% | 99.9% |
| Availability Zones | ≥3 | ≥3 | ≥3 | 1 | ≥3 | ≥3 |
| Minimum capacity charge per object | N/A | N/A | 128KB | 128KB | 40KB | 40KB |
| Minimum storage duration charge | N/A | 30 days | 30 days | 30 days | 90 days | 180 days |
| Retrieval fee | N/A | N/A | per GB retrieved | per GB retrieved | per GB retrieved | per GB retrieved |
| First byte latency | milliseconds | milliseconds | milliseconds | milliseconds | select minutes or hours | select hours |
| Storage type | Object | Object | Object | Object | Object | Object |
| Lifecycle transitions | Yes | Yes | Yes | Yes | Yes | Yes |

† Because S3 One Zone-IA stores data in a single AWS Availability Zone, data stored in this storage class will be lost in the event of Availability Zone destruction.

* S3 Intelligent-Tiering charges a small tiering fee and has a minimum eligible object size of 128KB for auto-tiering. Smaller objects may be stored but will always be charged at the Frequent Access tier rates. See the Amazon S3 Pricing for more information.

# S3 Storage Classes (continued)



The Above Diagram is from ACG

# S3 Storage Classes (continued)



The Above Information is from ACG

# S3 Storage Classes (continued)



The Above Diagram is from ACG

# S3 Storage Classes (continued)



**S3 PRICING TIERS**
## S3 - What Drives the Price?

A CLOUD GURU

| | |
|---|---|
| **S3 Standard - Infrequent Access *** - For long lived but infrequently accessed data that needs millisecond access | |
| All Storage / Month | $0.0125 per GB |
| **S3 One Zone - Infrequent Access *** - For re-createable infrequently accessed data that needs millisecond access | |
| All Storage / Month | $0.01 per GB |
| **S3 Glacier **** - For long-term backups and archives with retrieval option from 1 minute to 12 hours | |
| All Storage / Month | $0.004 per GB |
| **S3 Glacier Deep Archive **** - For long-term data archiving that is accessed once or twice in a year and can be restored within 12 hours | |
| All Storage / Month | $0.00099 per GB |

The Above Diagram is from ACG

# AWS Organizations



The Above Diagram is from ACG
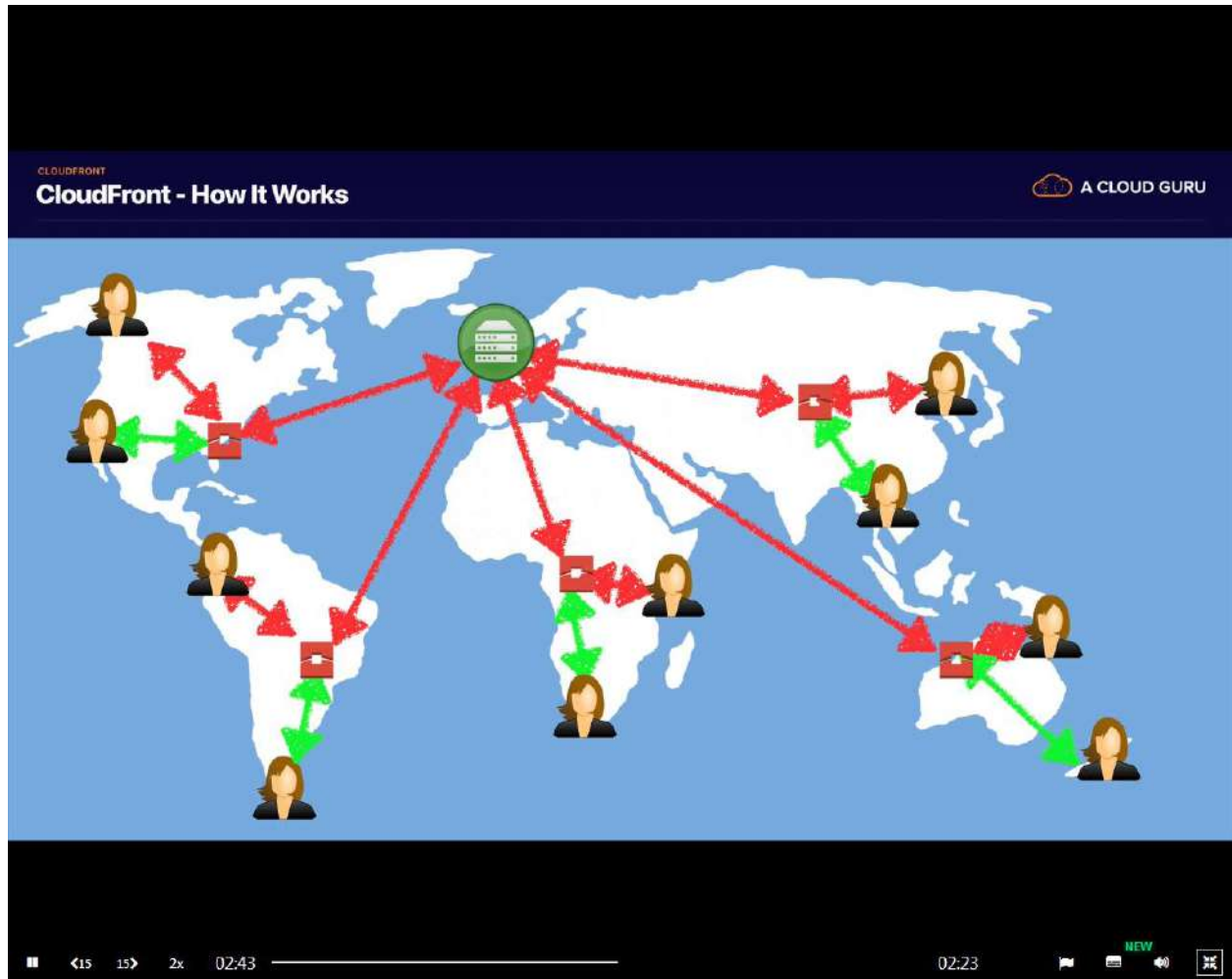
# AWS Organizations & Consolidated Billing



The Above Diagram is from ACG

# CloudFront Diagram

*The squares represent Edge Locations

*The green rack icon represents the origin of the file to be distributed



The above diagram is from ACG

# S3 Transfer Acceleration Diagram

*The squares represent Edge Locations

*The red bucket icon represents S3



The above diagram is from ACG

EBS Volume Types Chart

Via https://docs.aws.amazon.com/AWSEC2/latest/UserGuide/ebs-volume-types.html

# Volume characteristics

The following table describes the use cases and performance characteristics for each volume type. The default volume type is General Purpose SSD (gp2).

| | Solid-state drives (SSD) | | Hard disk drives (HDD) | |
|---|---|---|---|---|
| **Volume type** | General Purpose SSD (gp2) | Provisioned IOPS SSD (io1) | Throughput Optimized HDD (st1) | Cold HDD (sc1) |
| **Description** | General purpose SSD volume that balances price and performance for a wide variety of workloads | Highest-performance SSD volume for mission-critical low-latency or high-throughput workloads | Low-cost HDD volume designed for frequently accessed, throughput-intensive workloads | Lowest cost HDD volume designed for less frequently accessed workloads |
| **Use cases** | <ul><li>Recommended for most workloads</li><li>System boot volumes</li><li>Virtual desktops</li><li>Low-latency interactive apps</li><li>Development and test environments</li></ul> | <ul><li>Critical business applications that require sustained IOPS performance, or more than 16,000 IOPS or 250 MiB/s of throughput per volume</li><li>Large database workloads, such as:<ul><li>MongoDB</li><li>Cassandra</li><li>Microsoft SQL Server</li><li>MySQL</li><li>PostgreSQL</li><li>Oracle</li></ul></li></ul> | <ul><li>Streaming workloads requiring consistent, fast throughput at a low price</li><li>Big data</li><li>Data warehouses</li><li>Log processing</li><li>Cannot be a boot</li></ul> | <ul><li>Throughput-oriented storage for large volumes of data that is infrequently accessed</li><li>Scenarios where the lowest storage cost is important</li><li>Cannot be a boot volume</li></ul> |

volume

| API name | gp2 | io1 | st1 | sc1 |
|---|---|---|---|---|
| **Volume size** | 1 GiB - 16 TiB | 4 GiB - 16 TiB | 500 GiB - 16 TiB | 500 GiB - 16 TiB |
| **Max IOPS per volume** | 16,000 (16 KiB I/O) * | 64,000 (16 KiB I/O) † | 500 (1 MiB I/O) | 250 (1 MiB I/O) |
| **Max throughput per volume** | 250 MiB/s * | 1,000 MiB/s † | 500 MiB/s | 250 MiB/s |
| **Max IOPS per instance ††** | 80,000 | 80,000 | 80,000 | 80,000 |
| **Max throughput per instance ††** | 2,375 MB/s | 2,375 MB/s | 2,375 MB/s | 2,375 MB/s |
| **Dominant performance attribute** | IOPS | IOPS | MiB/s | MiB/s |

# EBS Volume Types (continued)



## Compare EBS Types

| | Solid-State Drives (SSD) | | Hard disk Drives (HDD) | | |
|---|---|---|---|---|---|
| Volume Type | General Purpose SSD | Provisioned IOPS SSD | Throughput Optimized HDD | Cold HDD | EBS Magnetic |
| Description | General purpose SSD volume that balances price and performance for a wide variety of transactional workloads | Highest-performance SSD volume designed for mission-critical applications | Low cost HDD volume designed for frequently accessed, throughput-intensive workloads | Lowest cost HDD volume designed for less frequently accessed workloads | Previous generation HDD |
| Use Cases | Most Work Loads | Databases | Big Data & Data Warehouses | File Servers | Workloads where data is infrequently accessed |
| API Name | gp2 | io1 | st1 | sc1 | Standard |
| Volume Size | 1 GiB - 16 TiB | 4 GiB - 16 TiB | 500 GiB - 16 TiB | 500 GiB - 16 TiB | 1 GiB-1 TiB |
| Max. IOPS**/Volume | 16,000 | 64,000 | 500 | 250 | 40-200 |

The above chart is from ACG

# Placement Group Diagrams

via https://docs.aws.amazon.com/AWSEC2/latest/UserGuide/placement-groups.html#placement-groups-cluster

Cluster Placement Group



Partition Placement Group

Spread Placement Group

# RDS – Restoring Backups



The above diagram is from ACG

## Multi-AZ

Multi-AZ keeps an exact copy of your production database in another AZ (for disaster recovery only). When you write to your primary AWS automatically synchronizes the changes to the standby database. In the event of failure, RDS will failover to the standby instance automatically. You keep the same DNS Endpoint (Amazon updates the IP address to point from one RDS instance to the other and Amazon does that for you so you don't have to do it)



The above diagram is from ACG

# Multi-AZ (continued)



The above chart is from ACG

# Read Replicas

Allow you to have a read only copy of your production database, for increased performance.



The above diagram is from ACG

# Read Replicas (continued)

Read Replicas are replicated **asynchronously**. It is possible to architect EC2 instances to read from different read replicas but only write to a single database.



The above diagram is from ACG

# Elasticache – Memcached vs Redis Chart



| 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 |

The above chart is from ACG

# DNS/Route53 Diagrams

## Simple Routing Policy



The above diagram is from ACG

# Weighted Routing Policy



The above diagram is from ACG

## Latency-Based Routing



The above diagram is from ACG

# Failover Routing Policy



The above diagram is from ACG

# Geolocation Routing Policy



The above diagram is from ACG

# Geoproximity Routing (Traffic Flow only mode)



The above diagram is from ACG

# Geoproximity Routing (continued)



The above information is from ACG

# Multivalue Answer Policy



The above diagram is from ACG

# VPC Diagram

When you create a custom VPC a Route Table, Network ACL, and Security Group are created by default. Then it is up to you to create subnets, create an internet Gateway to be attached to your VPC, create instances, etc.



The above diagram is from ACG

## VPC Peering

No Transitive Peering. You cannot peer through one VPC to another. You have to establish a new peering relationship. So if VPC B wants to talk to VPC C you have to establish a connection between the two.



The above diagram is from ACG

# NAT Gateway Diagram

A Nat Gateway (as well as a NAT instance) is used to provide internet traffic to EC2 instances in a private subnet. A NAT Gateway is redundant. A Nat instance is an individual EC2 instance located behind a security group.



The above diagram is from ACG

## Bastion Host Diagram

A Bastion is used to SSH or RDP into an instance in your private subnet.

Whereas, A Nat Gateway (as well as a NAT instance) is used to provide internet traffic to EC2 instances in a private subnet. A NAT Gateway is redundant. A Nat instance is an individual EC2 instance located **behind** a security group.



The above diagram is from ACG

# Direct Connect

-Directly connects your data center to AWS

-Useful for high throughput workloads (i.e. lots of network traffic)

-Stable Reliable Secure Connection



The above diagram is from ACG

# Direct Connect Setup Steps


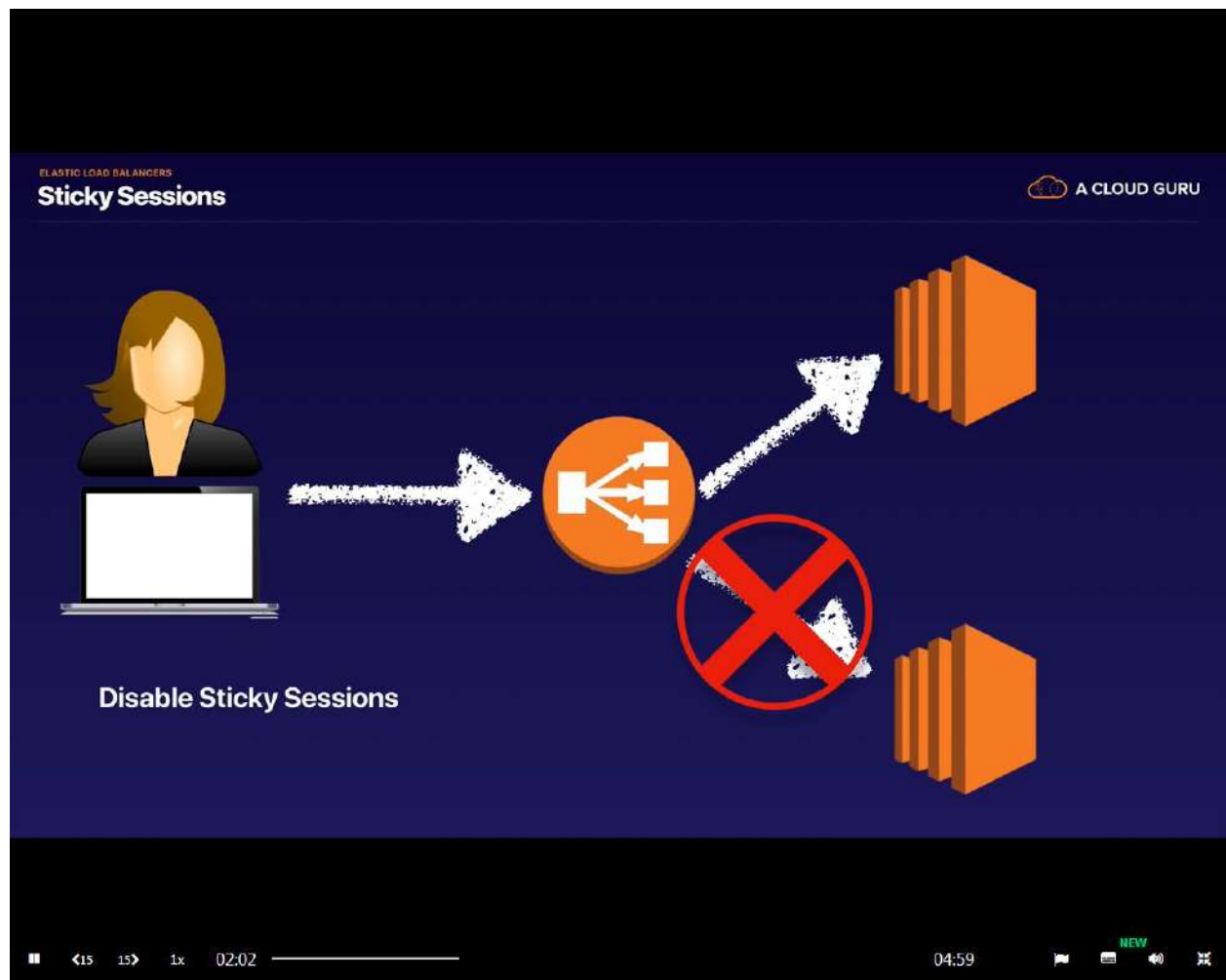
The above information is from ACG
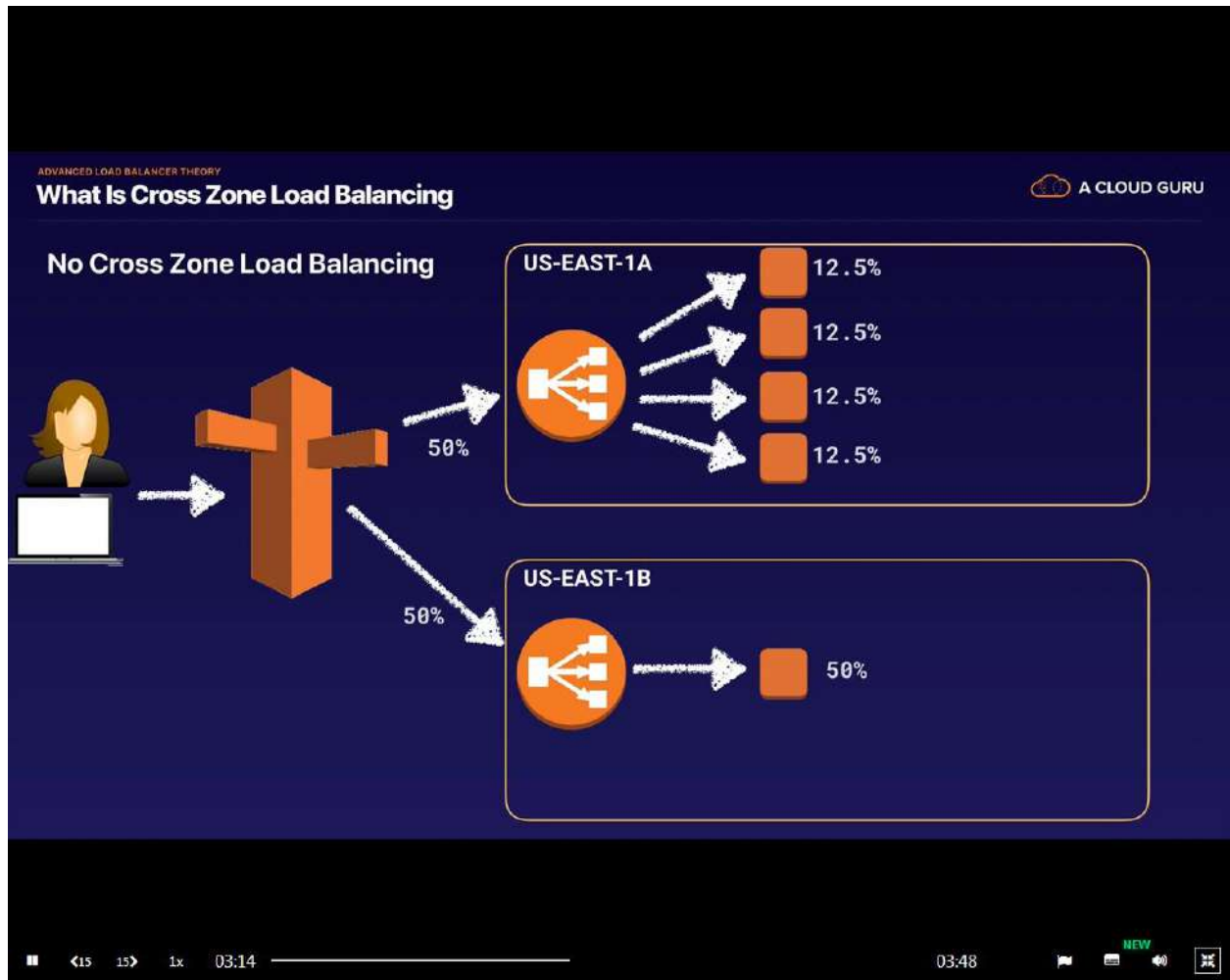
# Global Accelerator



The above diagram is from ACG

# Sticky Sessions

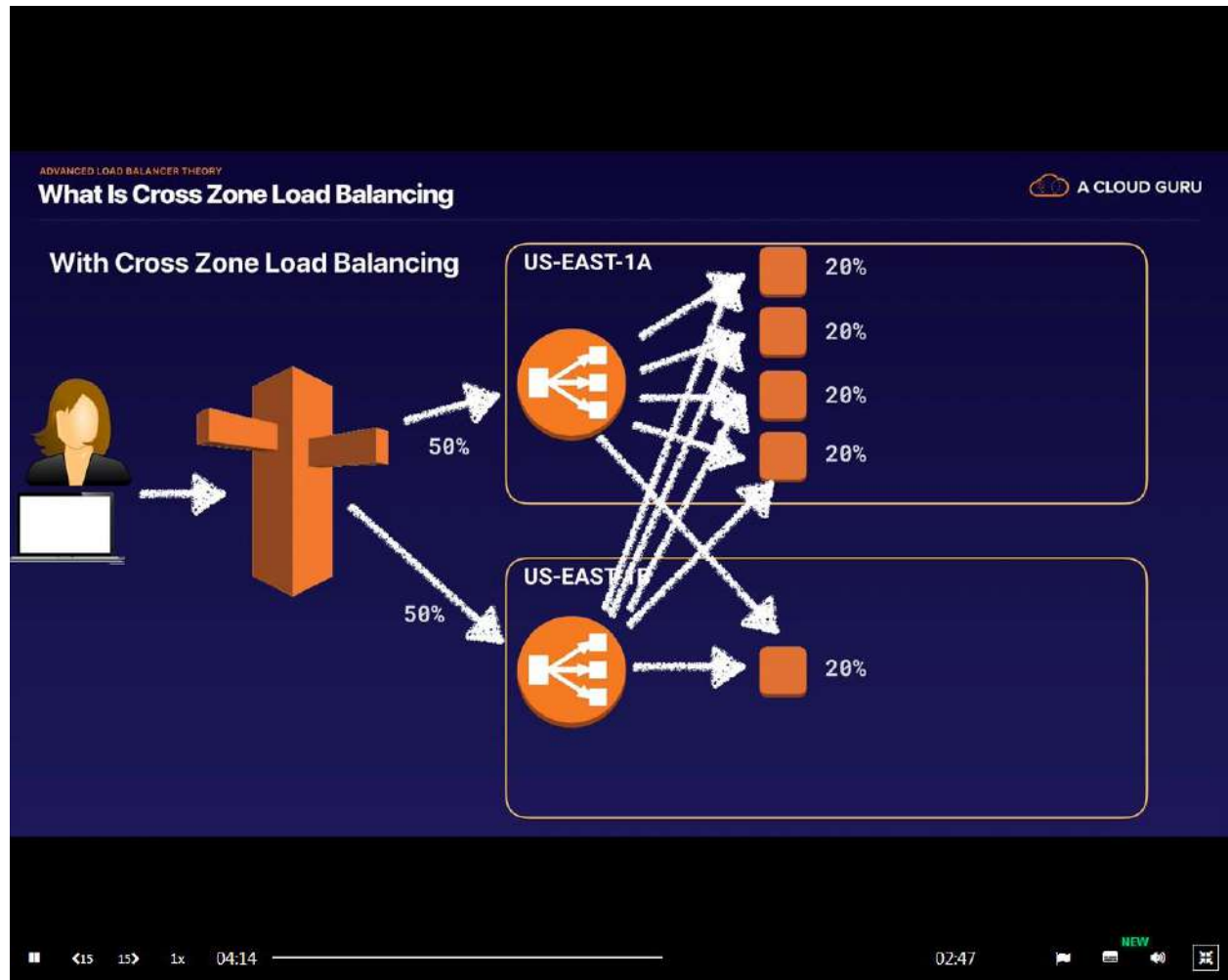Enabling and Disabling sticky sessions given the scenario



The above diagram is from ACG

# Cross Zone Load Balancing **NOT ENABLED**



The above diagram is from ACG
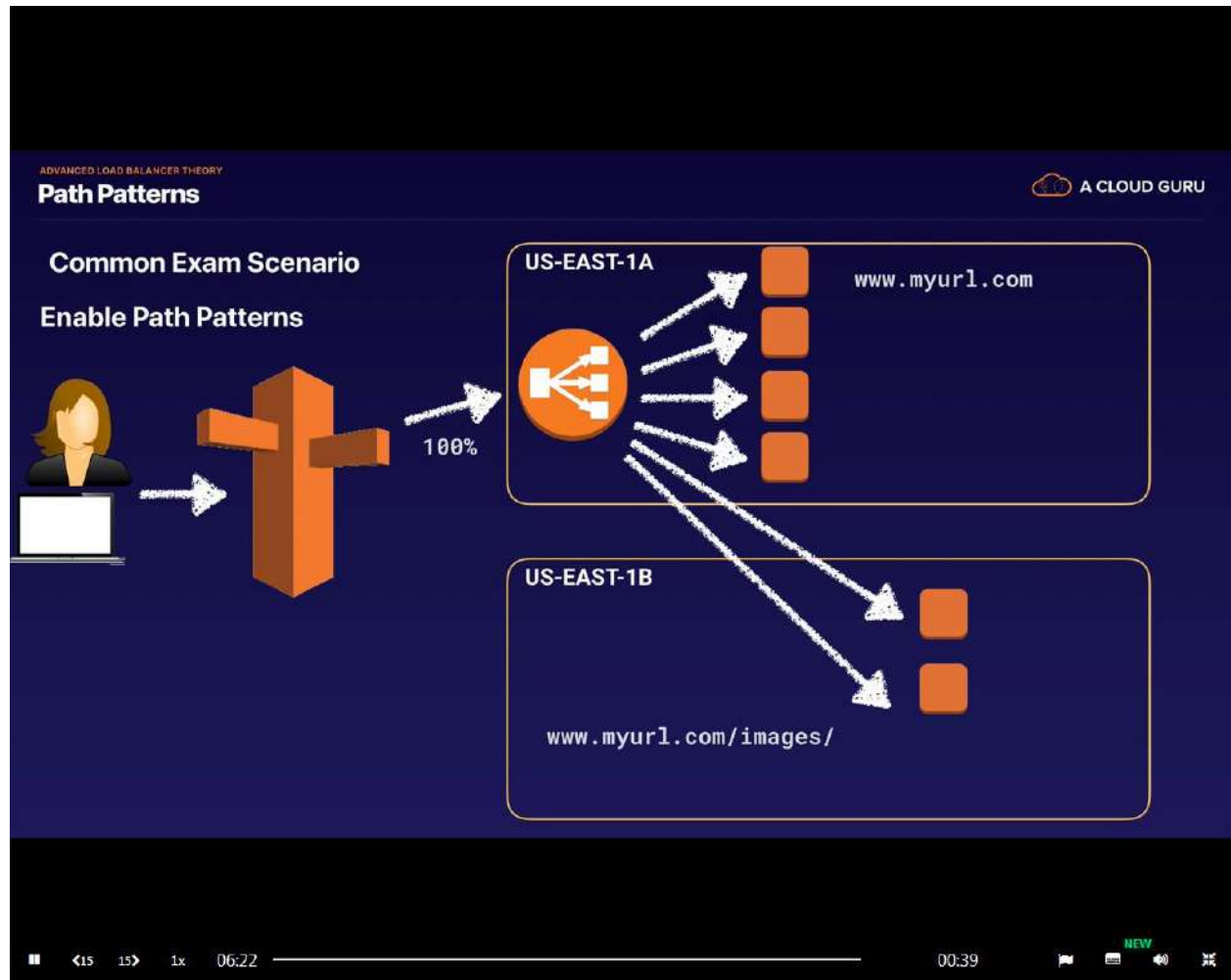
# Cross Zone Load Balancing **ENABLED**



The above diagram is from ACG

# Path Patterns

Example: Sending the Normal Path myurl.com to the 4 web servers in US-EAST-A1

And sending the myurl.com/images traffic to your media instances located in a separate availability zone
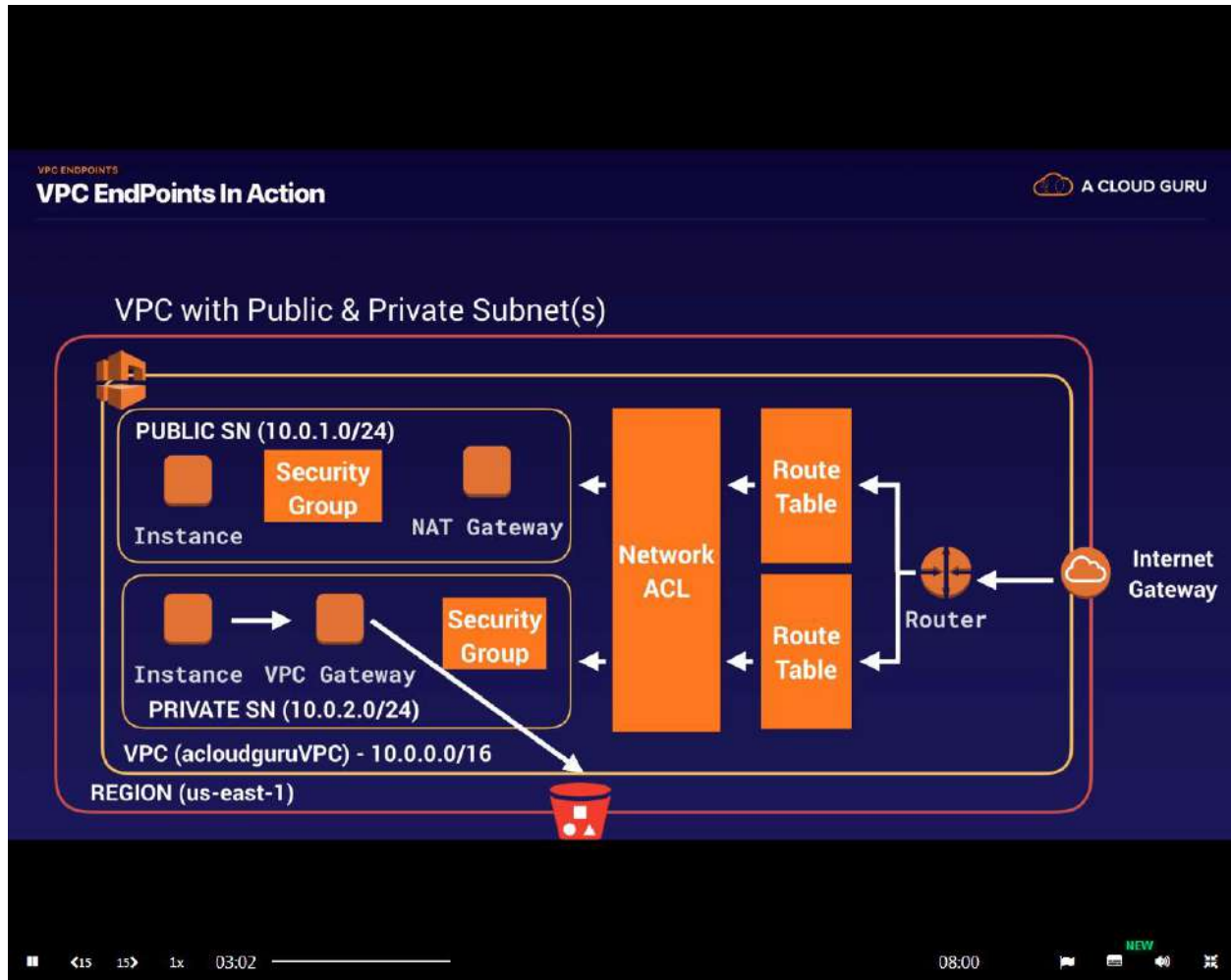


The above diagram is from ACG

# VPC Gateway Endpoint Diagram
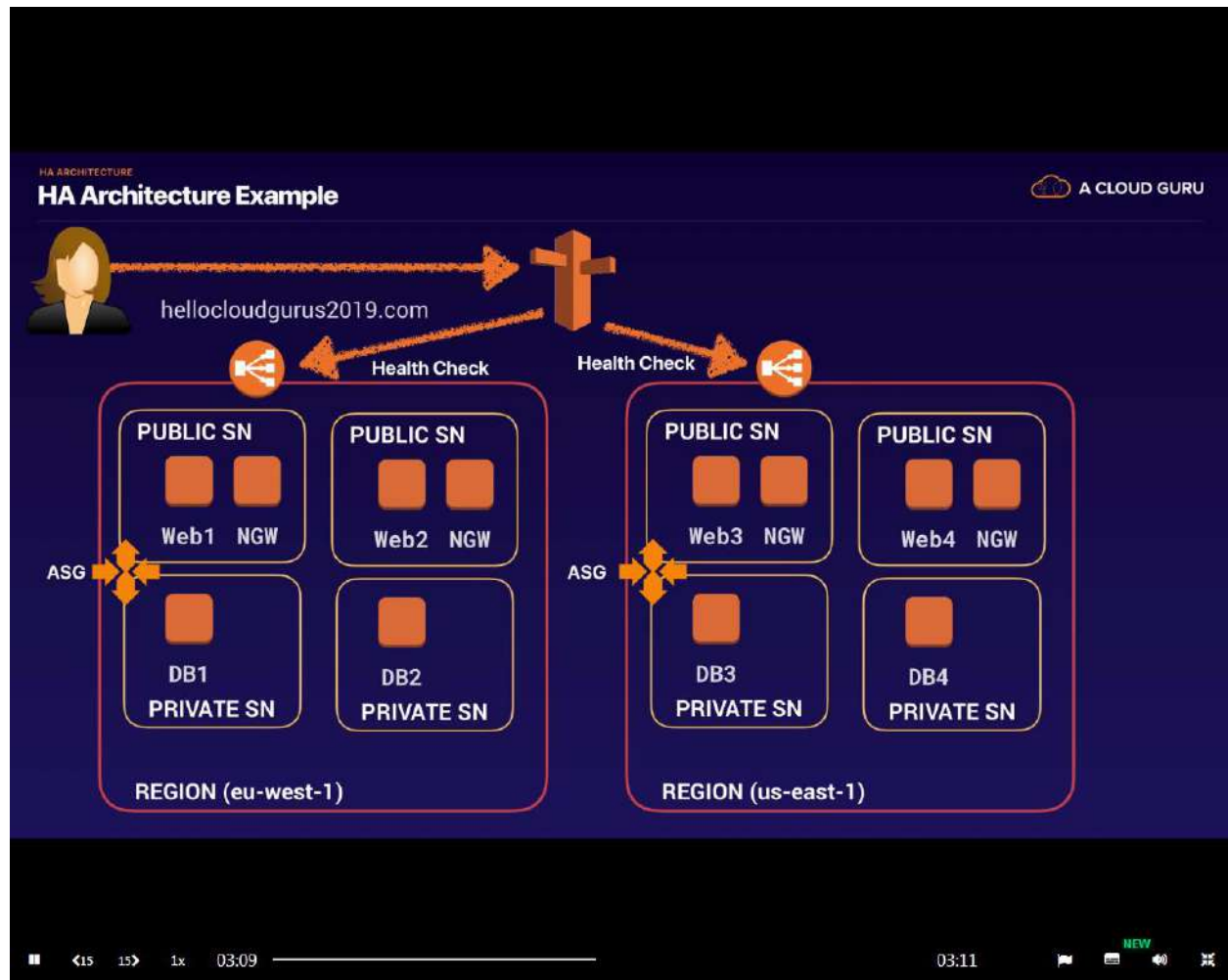
*The bucket represents S3

Our instance sends files to the VPC Gateway and that Gateway is going to send the file to our S3 bucket and it will not leave the Amazon network.



The above diagram is from ACG

# HA Architecture Example

If one of the regions goes down or an AZ goes down then you've got failover. You can failover from one region to another or one AZ to another.



The above diagram is from ACG
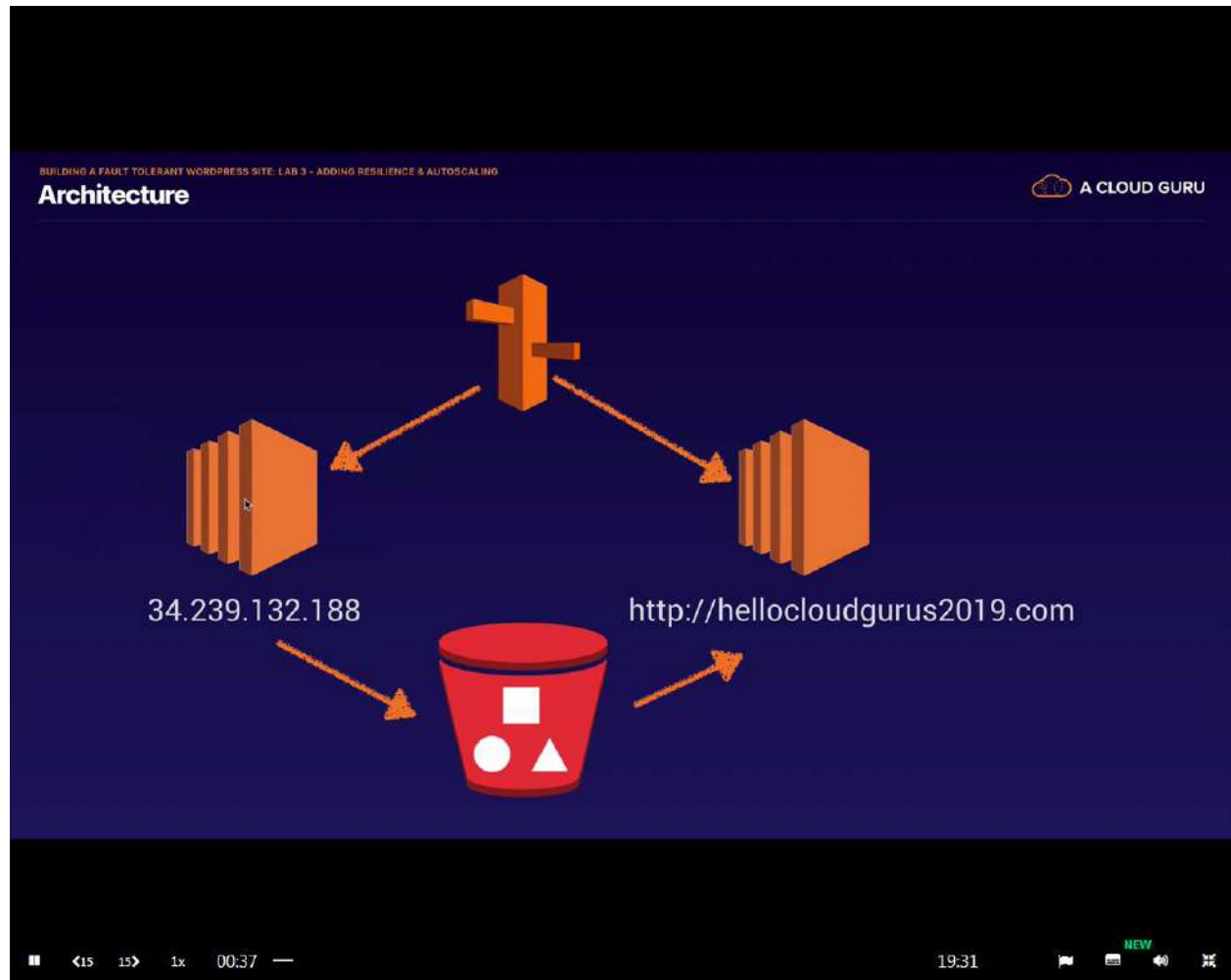
# Building a fault tolerant WordPress Site

Example: Here we have a user browsing the internet to our Route53 domain name which will connect up to an elastic load balancer. We have some EC2 instances behind an autoscaling group which are going to be in separate AZ's. We have RDS instances that are multi-AZ. We have two S3 buckets, one for our media and one for our code. And we serve our pictures from our wordpress site through cloudfront.



The above diagram is from ACG
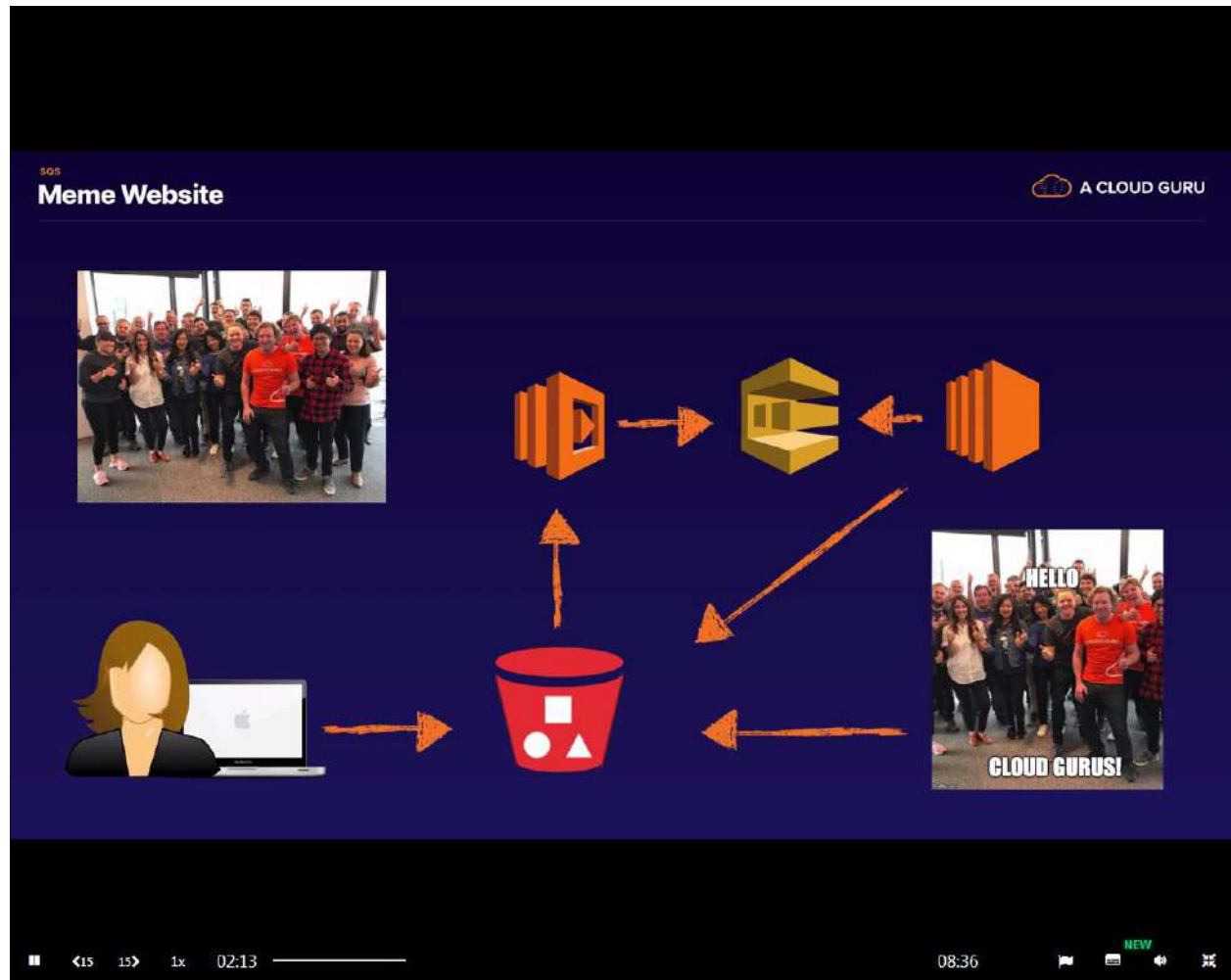
# Fault Tolerant Wordpress Site Network Diagram

On the left is the IP of our writer node. Our blog writers will navigate directly to the IP address on the left side of the diagram. That EC2 instance will push any changes to S3. The Fleet of EC2 instances on the right will be pulling the S3 bucket every minute looking for changes. Our site visitors will visit the domain and Route53 will send them to the fleet of EC2 instances on the right side of the diagram, so it is just sending them to our read nodes only. (We will also have the writer node instances on the right be an autoscaling group situated behind an ELB)
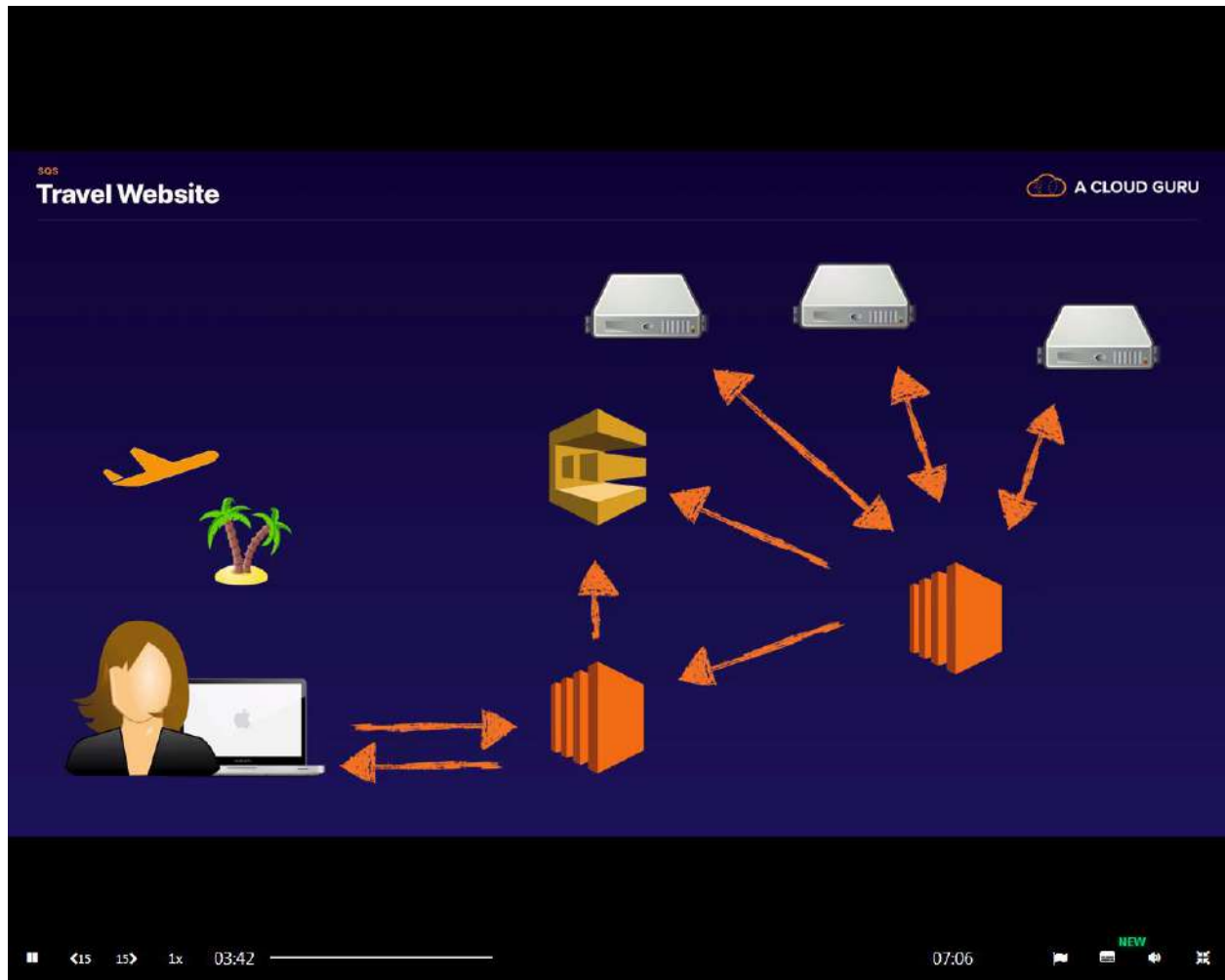


The above diagram is from ACG

# SQS

User uploads a photo to s3 which triggers a lambda function which will take the image and write text over it. It will store that text in SQS. Then a fleet of EC2 instances will pull the message queue for work and it will then create the meme and store it in S3. If an EC2 instance fails while creating the meme the message will become available again in the queue and another instance will create the meme. So SQS is storing the message independently so that even if an EC2 instance can't process the message, another EC2 instance will come along and take the message.



The above diagram is from ACG

# SQS (Continued)

User  goes to ec2 server and says they want to go to Rome on these particular dates. The web server passes that information to an SQS queue. A fleet of application servers are configured to pull that message and look for different airlines. Once the information is retrieved it is passed back to the web server and then back to the end user. If we lose an individual ec2 instance we won't lose the information and another one will come along and poll the queue to do the work and ultimately return the result back to the end user.



The above diagram is from ACG

# SQS Standard Queues



The above diagram is from ACG

# FIFO Queues



The above diagram is from ACG

# FIFO Queues



The above diagram is from ACG

## SNS Availability



The above diagram is from ACG

## Elastic Transcoder

Upload a video into S3 Bucket…then a lambda function will take the metadata of the video and send it to Elastic Transcoder which will transcode the video so that it looks good on various devices and is in high resolution and then stores the transcoded video in another S3 Bucket.



The above diagram is from ACG

## API Gateway

Users do a call to our API Gateway. API Gateway is the 'front-end/front door' to our AWS environment. That API call could be passed to Lambda, or an EC2 instance, or writing to DynamoDB, etc.



The above diagram is from ACG

# API Gateway Caching

For example, User 1 makes a get request that is forwarded on to lambda and lambda returns a response. Then User 2 makes the same get request and since API Gateway has cached it, it does not go to the Lambda function.



The above diagram is from ACG

## Kinesis Streams

Data Producers (such as an EC2 monitoring stock prices, or perhaps IoT monitoring farm data) stream the data to kinesis and Kinesis Streams is a place to store that data coming in (by default stores for 24 hrs but can be stored for 7 days). The data is contained in 'shards'. EC2 instances (aka the data consumers) can then analyze that data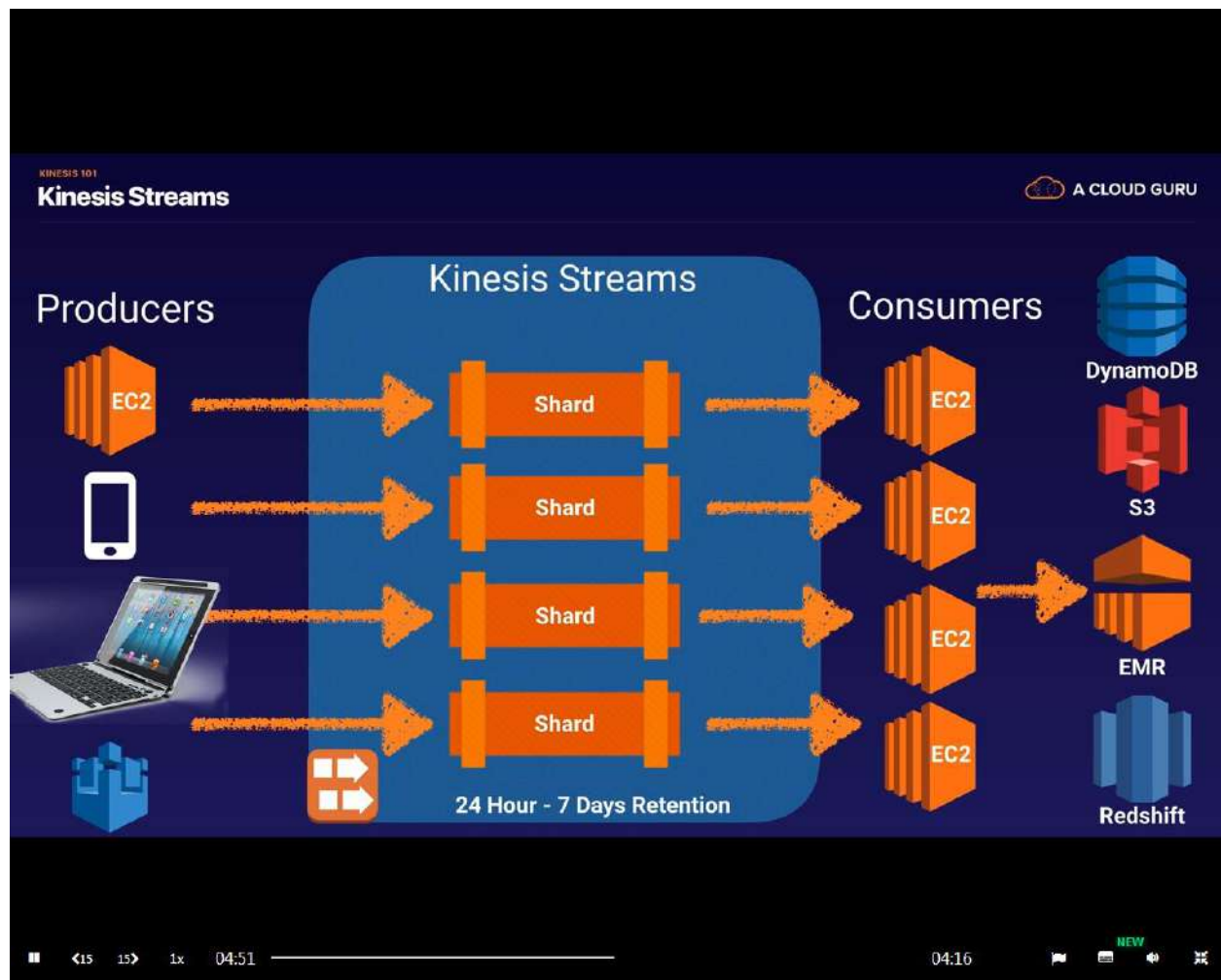 located within the shards and then can store it in various places such as DynamoDB, S3, EMR, Redshift, RDS etc. Kinesis allows you to persistently store your data for 24 hrs to 7 days, while your data consumers do something with that data.



The above diagram is from ACG

## Kinesis Firehose

Data Producers send the data to Kinesis Firehose (of which does NOT have persistent storage and does NOT have 'shards'). The data has to be analyzed as it comes in. So perhaps you have lambda function within your Kinesis Firehose running a set of code for the data as it comes in and outputs it somewhere. It can output it to S3. Or it could output it to S3 and then you can import it to Redshift. Or it can output to Elasticsearch Cluster, etc.



The above diagram is from ACG

# Kinesis Analytics

Works with Kinesis Streams, as well as with Kinesis Firehose. It can analyze the data on the fly inside either service and then it stores this data either on S3, Redshift, or Elasticsearch Cluster.
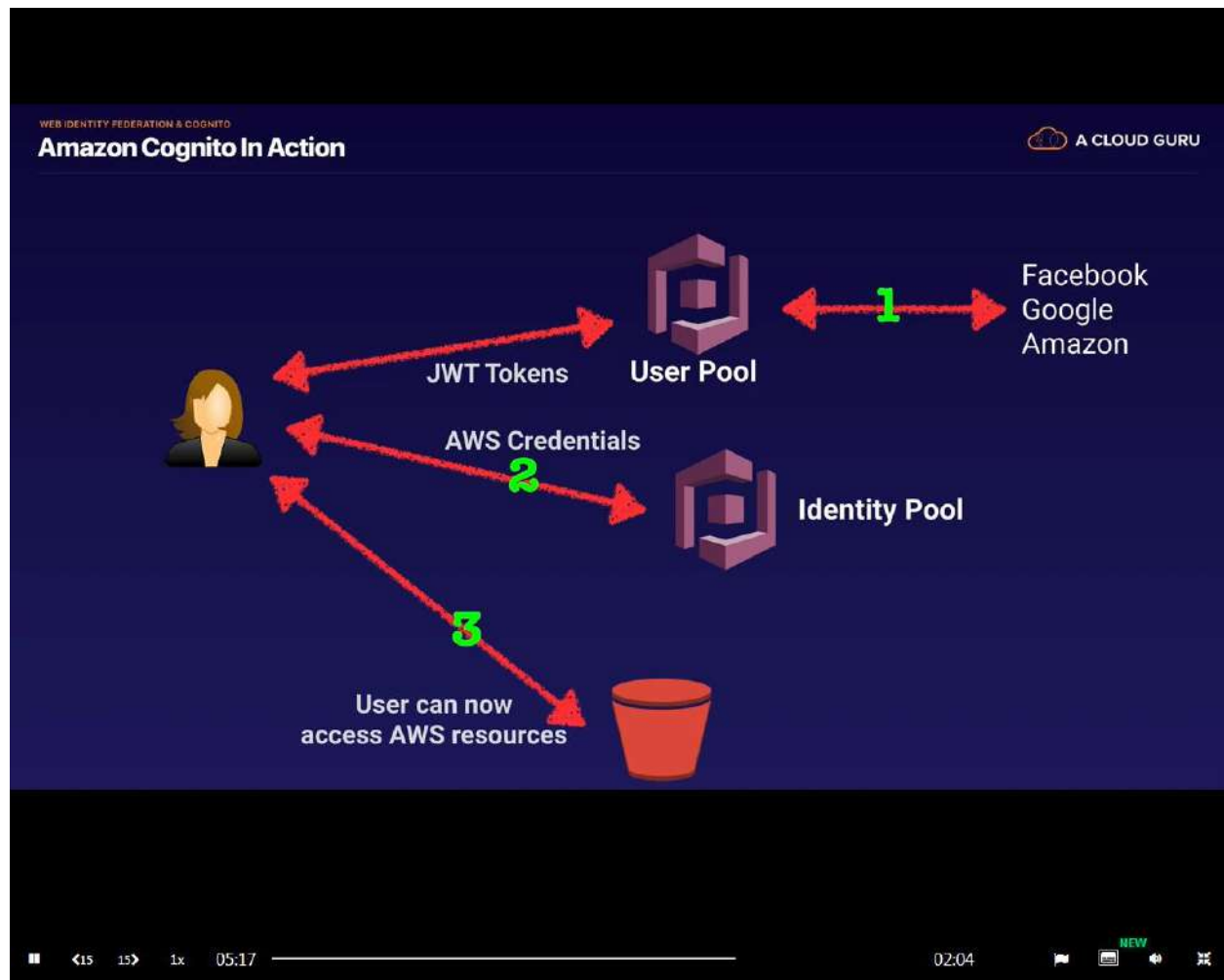


The above diagram is from ACG

## Amazon Cognito

User logs in with Facebook account. Facebook authenticates her account and passes back an authentication token to Cognito, which converts it to a JWT token. The user sends the token to an identity pool and that identity pool will grant her AWS credentials in the form of an IAM role and then she will have access to AWS resources. User pools pertains to registration, authentication, usernames, and passwords. Identity pools pertains to **the actual granting of access** to use AWS resources.
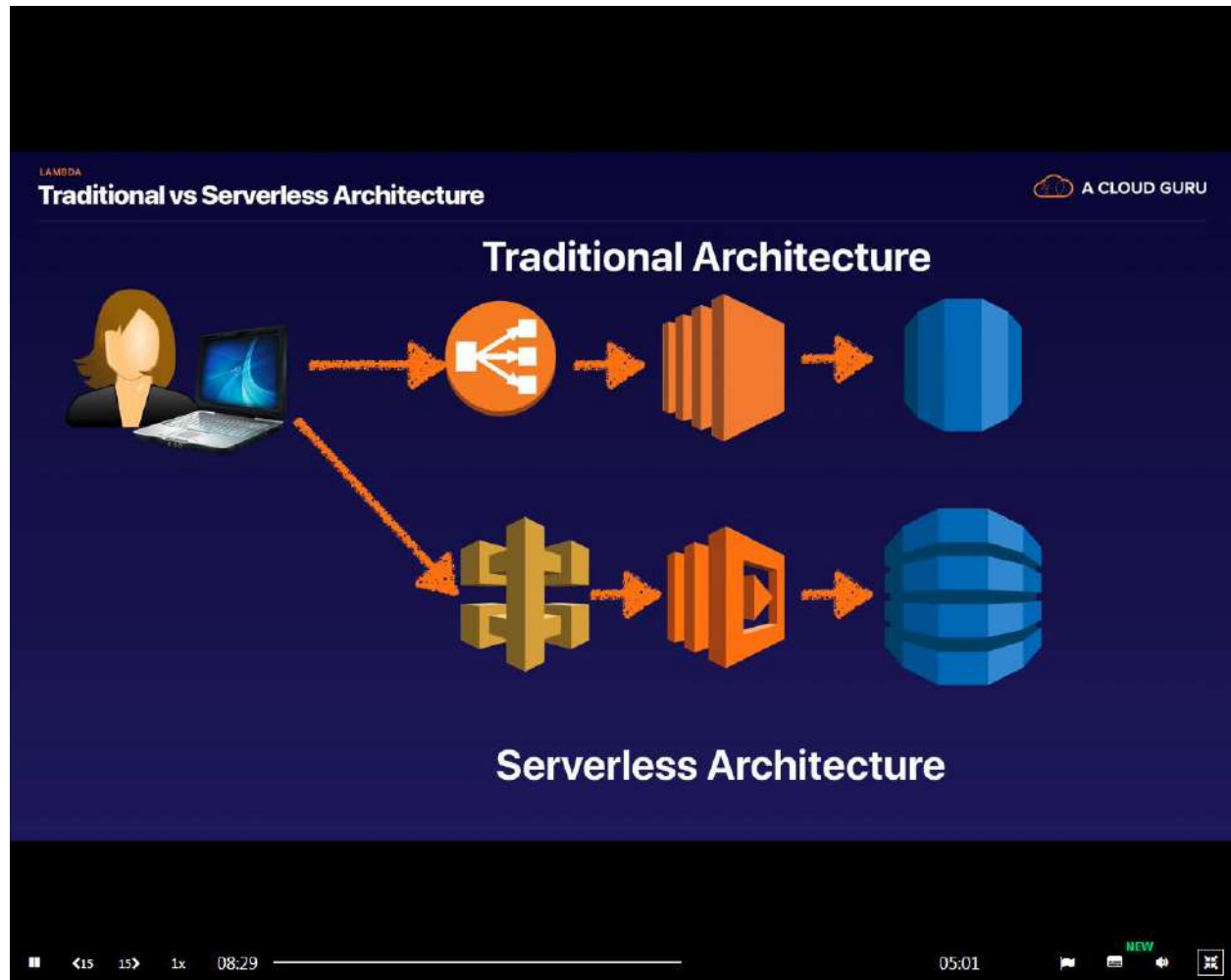


The above diagram is from ACG

# Traditional vs Serverless Architecture

Traditional

ELB ->EC2->Database Storage (RDS)

Serverless
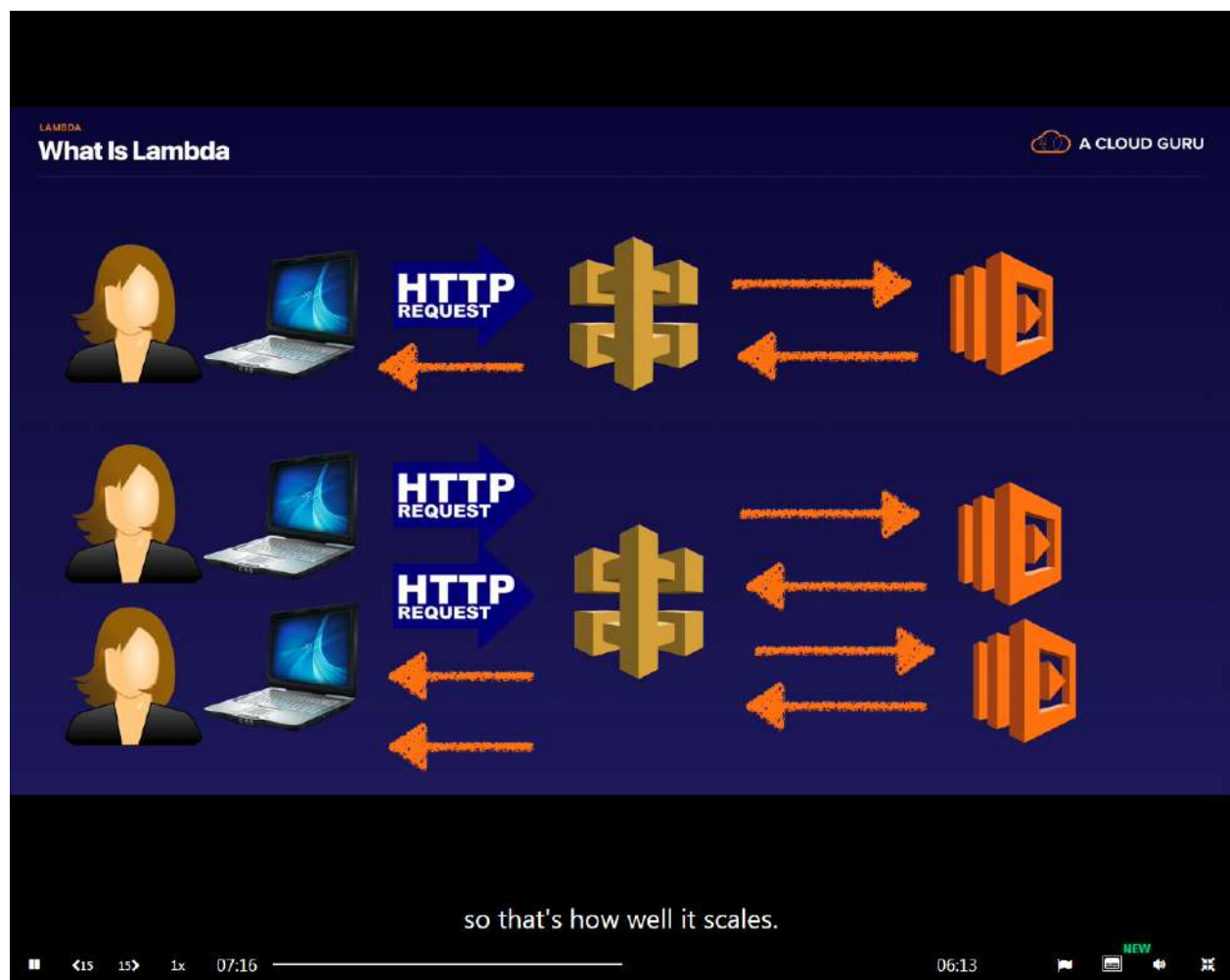
API Gateway ->Lambda->DynamoDB



The above diagram is from ACG

## Lambda

User sends http request to API Gateway. API Gateway proxies that to Lambda. Lambda will run the code in response to that HTTP request, and then send it back to API Gateway, that will send it back to the user.

If two users have two separate HTTP requests, the same lambda function will be triggered but it will be separately run, so it will be two separate lambda functions. If you have a million users hitting your API gateway at once, it will trigger a million lambda functions, so that's how well it scales.

# Serverless Website with API Gateway & Lambda

User wants to go to helloucloudgurus2019.com. User is sending a query across to Route53 which will respond with the bucket address for our website. So our user goes to our S3 bucket and they are going to go to our index.html, which will show up as a static page. But it will have a button. When the button is pushed they are going to get dynamic content. Because a request is send through to API gateway which will proxy a request to a lambda function which will take the data and return a result to API gateway, which will then return a result, to our user.