## What is DNS?

DNS, or Domain Name System, is a technology used to translate domain names into IP addresses and vice versa.

Here's the main components and functions of DNS:

1. **Domain Names:** A domain name is a user-friendly, human-readable address that is used to identify resources on the internet. Examples include "www.example.com" or "google.com."
2. **IP Addresses:** IP addresses (Internet Protocol addresses) are numeric addresses used to uniquely identify devices on the internet or within a network. There are two types of IP addresses: IPv4 (e.g., 192.168.1.1) and IPv6 (e.g., 2001:0db8:85a3:0000:0000:8a2e:0370:7334).
3. **DNS Resolution:** When you enter a domain name into a web browser, it needs to determine the corresponding IP address to establish a connection. This process is called DNS resolution.
4. **DNS Servers:** DNS servers are responsible for storing and managing domain name-to-IP address mappings. They maintain databases called DNS zones that contain records associating domain names with IP addresses.
5. **DNS Hierarchy:** DNS operates in a hierarchical manner, with a root domain at the top, followed by top-level domains (TLDs), second-level domains, and subdomains. For example:
    a. - Root domain: "."
    b. - TLDs: ".com," ".org," ".net," etc.
    c. - Second-level domains: "example.com," "google.com," "wikipedia.org," etc.
    d. - Subdomains: "www.example.com," "mail.google.com," "en.wikipedia.org," etc.
6. **DNS Records:** DNS zones contain various types of DNS records, including:
    a. A (Address) records: Map domain names to IPv4 addresses.
    b. AAAA (IPv6 Address) records: Map domain names to IPv6 addresses.
    c. CNAME (Canonical Name) records: Create aliases for domain names.
    d. MX (Mail Exchange) records: Specify mail server destinations for a domain.
    e. TXT (Text) records: Store textual information, often used for verification and authentication purposes.
    f. NS (Name Server) records: Identify authoritative DNS servers for a domain.

7. **DNS Resolvers:** DNS resolvers are software components or DNS servers that initiate DNS queries to resolve domain names. They send DNS queries to authoritative DNS servers and cache the results for future use.

## AWS Route53

Amazon Route 53 is a scalable and highly available Domain Name System (DNS) web service provided by Amazon Web Services (AWS). It's designed to route end-user requests to endpoints such as EC2 instances/Web Servers, S3 buckets, load balancers etc. It performs three functions:

**1. Domain Name Registration:** Your website needs a name, such as example.com. Route 53 allows you to register domain name for your website or transfer existing ones. You can manage the DNS records for these domains, including A records, CNAME records, MX records, and more.

**2. Traffic Routing:** Route 53 can route traffic based on various routing policies, including simple routing, weighted routing, latency-based routing, geolocation-based routing, and more.

**3. Health Checks:** You can set up health checks for your resources, and Route 53 can automatically reroute traffic away from unhealthy endpoints to healthy ones.

## Route53 – Domain Registration

Here's an overview of how you register a domain name with Amazon Route 53:

1. You choose a domain name and confirm that it's available, meaning that no one else has registered the domain name that you want.
2. When you register a domain, you provide names and contact information for the domain owner and other contacts. When you register a domain with Route 53, the service automatically makes itself the DNS service for the domain by doing the following:
   a. Creates a hosted zone that has the same name as your domain.
   b. Assigns a set of four name servers to the hosted zone. When someone uses a browser to access your website, such as www.example.com, these name servers tell the browser where to find your resources, such as a web server or an Amazon S3 bucket.
   c. Gets the name servers from the hosted zone and adds them to the domain.
3. At the end of the registration process, we send your information to the registrar for the domain.
4. The registrar sends your information to the *registry* for the domain. A registry is a company that sells domain registrations for one or more top-level domains, such as .com.
5. The registry stores the information about your domain in their own database and also stores some of the information in the public WHOIS database.

To register a domain using Amazon Route 53, you can follow these steps:

1. Sign in to the AWS Management Console and open the Route 53 service console.
2. In the Route 53 dashboard, select "Registered domains" from the left-hand navigation pane and Click the "Register Domain" button.
3. Enter Domain Name you want to register (e.g., example.com) in the provided field.
4. Click the "Check" button to check if the domain name is available for registration. If it's available, you'll be able to proceed.
5. If domain name is not available, you can try other names.
6. If the domain name is available, you'll be presented with details about the domain and its pricing.
7. Review the information, select the desired registration duration (usually in years), and click the "Add to Cart" button.
8. Review the domain registration details, and make sure everything is correct. Click the "Continue" button to proceed.
9. Configure Domain Settings like contact information, DNS routing, and privacy protection (WHOIS privacy) if desired.
10. Follow the on-screen instructions to provide the necessary information.
11. Review and Checkout
12. Enter your billing and payment information to complete the domain registration process.
13. Review the final details, including the total cost and Click the "Confirm and Buy" button to complete the domain registration.
14. Once the registration is successful, you'll receive a confirmation email, and the domain will appear in your Route 53 registered domains list.
15. After domain registration, you can configure DNS records in the Route 53 hosted zone associated with your domain to point to AWS resources or services.

## Route 53 - Record Sets

Here are the most common records in AWS. When you create records, you specify the following values for each record: Name, Type, Alias, TTL (Time to Live), Value, Routing Policy.

- **A Record (URL to IPv4)**
  The "A" record stands for Address record. The A record is used by computer to translate the name of the domain to an IP address.

- **CNAME (Canonical Records- URL to URL)**
  CNAME Points a URL to any other URL. (gaurav.gupta.com => gkg.example.com), We use it only for Non-Root Domain (aka. something.mydomain.com)

- **Alias Record:**
  Alias record points a URL to an AWS Resource, Alias record are used to map resource record sets in your hosted zone to Elastic Load Balancer, CloudFront or S3 Buckets websites.
- **AAAA: (URL to IPv6)**
  An **AAAA record** maps a domain name to the IP address (Version 6) of the computer hosting the domain. An **AAAA record** is used to find the IP address of a computer connected to the internet from a name.
- **MX Record (Main Exchange Record)**
  A mail Exchanger **record** (**MX record**) specifies the mail server responsible for accepting email messages on behalf of a domain name. It is a resource **record** in the Domain Name System (**DNS**). It is possible to configure several **MX records**, typically pointing to an array of mail servers for load balancing and redundancy.

## Route53 – Create our First Records

To create your first DNS records in Amazon Route 53, follow these steps:

1. Sign in to the AWS Management Console and navigate to the Route 53 service.
2. Create a Hosted Zone: If you haven't already created a hosted zone for your domain, you'll need to do create. Click on "Hosted zones" in the left navigation pane and then click the "Create hosted zone" button.
3. Enter your domain name (e.g., example.com) and choose a comment (optional). Click the "Create hosted zone" button.
4. Create DNS Records: Inside your hosted zone, you can create DNS records. To create records, follow these steps:
   a. Click on your hosted zone in the list of hosted zones.
   b. Under the "Record sets" tab, click the "Create Record Set" button.
   c. Configure the record set:
      i. Name: This is typically the subdomain or a domain apex (leave empty for the domain itself or enter "www" for the www subdomain).
      ii. Type: Select the type of DNS record you want to create (e.g., A for IPv4 address, CNAME for canonical name).
      iii. Alias: If you're pointing the record to an AWS resource (e.g., S3 bucket, CloudFront distribution, ELB), you can enable "Alias" and select the target resource from the dropdown list.
      iv. Value: Enter the value for the DNS record. This will vary depending on the record type. For example, for an A record, it's an IPv4 address; for a CNAME, it's the canonical name.
   d. Click the "Create" button to add the DNS record.

5. Repeat for Additional Records: To create additional records, repeat the process by clicking the "Create Record Set" button, specifying the name, type, alias, and value for each record.
6. Review and Save Changes: After creating your DNS records, review them to ensure they are accurate.
7. Click the "Save changes" button to save your DNS records.
8. DNS Propagation: It may take some time (usually minutes to a few hours) for DNS changes to propagate across the internet.
9. Test Your DNS Records: Once DNS propagation is complete, you can test your DNS records by entering your domain name or subdomain in a web browser and confirming that it resolves to the expected resource.

## Route 53 - TTL

- In Amazon Route 53, TTL (Time-to-Live) is a setting associated with DNS (Domain Name System) records that specifies how long a DNS resolver or cache should cache the information about a DNS record before making a new request to the authoritative DNS server.
- TTL is measured in seconds.

Here's what you need to know about TTL in Route 53:

- Default TTL: When you create a new DNS record in Route 53, it typically comes with a default TTL value. The default TTL value varies depending on the record type but is commonly set to 300 seconds (5 minutes).
- Custom TTL: You can customize the TTL value for each DNS record you create in Route 53. This allows you to control how long resolvers and caches should store the record's information. Common TTL values are 60 seconds, 300 seconds, 3600 seconds (1 hour), and 86400 seconds (1 day), but you can set it to any value you prefer.
- Impact of TTL:
  - Short TTL: Setting a short TTL means that DNS resolvers and caches will only store the record's information for a short period. This is useful when you anticipate frequent changes to the record and want updates to propagate quickly. However, it can increase the load on DNS servers due to more frequent lookups.
  - Long TTL: A longer TTL reduces the frequency of DNS queries, reducing the load on DNS servers and improving performance for end-users. However, it means that changes to the record may take longer to propagate when updates are made.
- Changing TTL: You can change the TTL value for a DNS record at any time within the Route 53 console. Simply edit the record set, modify the TTL value, and save the changes. Keep in mind that changing the TTL won't affect cached records until the existing cached records expire.

## Route 53 CNAME Vs. ALIAS Record

- In Amazon Route 53, both CNAME (Canonical Name) records and ALIAS records are used to map one domain name to another, allowing for flexible DNS configurations. However, they serve slightly different purposes and have distinct characteristics:

## CNAME (Canonical Name) Records

- Usage: CNAME records are typically used to create an alias for one domain name to another domain name. For example, you can use a CNAME to map a subdomain like "www.example.com" to another domain like "myapp.example.net."
- Targets: CNAME records can only point to other domain names, not to IP addresses or AWS resources directly.
- TTL: CNAME records have their own Time-to-Live (TTL) value, which determines how long DNS resolvers and caches should cache the CNAME record.

## ALIAS Records

- Usage: ALIAS records are specific to Route 53 and are used to map a domain name to AWS resources or other Route 53 DNS records. They are used when you want to map a domain name to an AWS resource like an S3 bucket, CloudFront distribution, ELB (Elastic Load Balancer), or an AWS Global Accelerator.
- Targets: ALIAS records can point to AWS resources or other Route 53 DNS records that are configured as aliases.
- No TTL: ALIAS records do not have their own TTL values. They inherit the TTL of the target resource, which is managed automatically by AWS.

## CNAME Vs. ALIAS Records - Comparison and Use Cases

- Use CNAME records when you need to create an alias for a domain name to point to arbitrary resources, including non-AWS resources.
- Use ALIAS records when you want to map a domain name to an AWS resource or another Route 53 DNS record that is configured as an alias. ALIAS records are especially useful when you want to map to AWS resources that may change IP addresses or endpoints, as they automatically update.
- In summary, CNAME records are more general-purpose and can point to various resources, while ALIAS records are specific to Route 53 and are used to map domain names to AWS resources with automatic updates.

# Route 53 – Routing Policies

There are total 6 types of routing policy in Route53.

- Simple Routing Policy
- Weighted Routing Policy
- Latency Based Routing Policy
- Failover Routing Policy
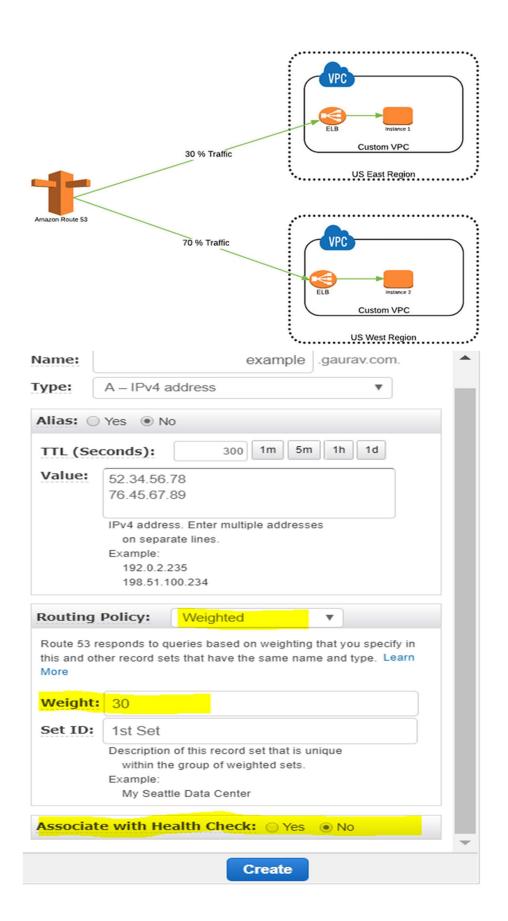- Geo-location Routing Policy
- Multivalue Answer Routing Policy

## Simple Routing Policy

- With simple routing, you typically route traffic to a single resource, for example, to a web server for your website.
- If you choose the simple routing policy in the Route 53 console, you can't create multiple records that have the same name and type, but you can specify multiple values in the same record, such as multiple IP addresses.
- If you specify multiple values in a record, Route 53 returns all values to the recursive resolver in random order, and the resolver returns the values to the client (such as a web browser) that submitted the DNS query.
- The client then chooses a value and resubmits the query.
- With simple routing policy, although you can specify multiple IP addresses, these IP addresses are not health checked.
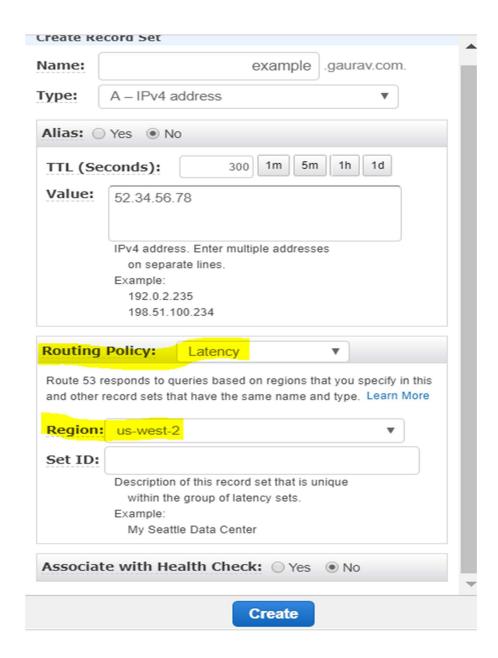
## Weighted Routing Policy

- Weighted routing lets you associate multiple resources with a single domain name (example.com) or subdomain name (acme.example.com) and choose how much traffic is routed to each resource.

- This can be useful for a variety of purposes, including load balancing and testing new versions of software.

- To configure weighted routing, you create records that have the same name and type for each of your resources.

- You assign each record a relative weight that corresponds with how much traffic you want to send to each resource.

- Weighted Routing Policy controls the what percentage % of the requests that go to specific resource. It's helpful to test 1% of traffic on new app version.

- It is also helpful to split traffic between two regions.

- We can associate Health checks with it.

- Example: Suppose you have two regions hosting your application, one in North America and one in Europe. You want to direct 70% of your traffic to the North American servers and 30% to the European servers. You can create two A records with weights 70 and 30, respectively.
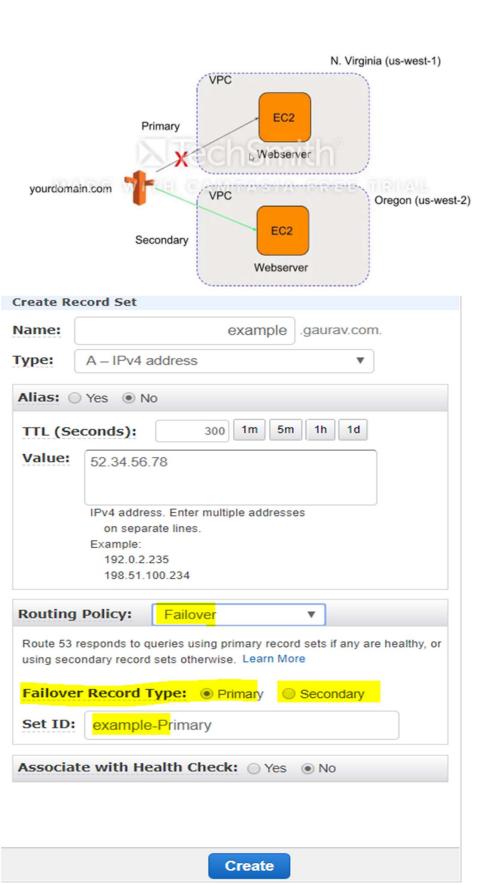
**Name:** example .gaurav.com.

**Type:** A – IPv4 address ▼

**Alias:** ○ Yes  ● No

**TTL (Seconds):** 300  | 1m | 5m | 1h | 1d |

**Value:**
52.34.56.78
76.45.67.89

IPv4 address. Enter multiple addresses
on separate lines.
Example:
192.0.2.235
198.51.100.234

**Routing Policy:** Weighted ▼

Route 53 responds to queries based on weighting that you specify in this and other record sets that have the same name and type. Learn More

**Weight:** 30

**Set ID:** 1st Set

Description of this record set that is unique
within the group of weighted sets.
Example:
My Seattle Data Center

**Associate with Health Check:** ○ Yes  ● No

**Create**

## Latency Routing Policy

- If your application is hosted in multiple AWS Regions, you can improve performance for your users by serving their requests from the AWS Region that provides the lowest latency.

- It allows you to route your traffic based on lowest network latency for your end user. It redirects to the server that has the least latency close to us.

- This policy is particularly useful when you want to direct users to the resource that provides the best performance or lowest response time for their location.

- When you have resources or services distributed across multiple AWS regions, you can use latency routing to direct users to the region with the lowest latency, optimizing the user experience.

- To use latency-based routing, you create latency records for your resources in multiple AWS Regions.

- When Route 53 receives a DNS query for your domain or subdomain (example.com or acme.example.com), it determines which AWS Regions you've created latency records for, determines which Region gives the user the lowest latency, and then selects a latency record for that Region.

- Route 53 responds with the value from the selected record, such as the IP address for a web server.

- For example, suppose you have Elastic Load Balancing load balancers in the US West (Oregon) Region and in the Asia Pacific (Singapore) Region. You create a latency record for each load balancer. Here's what happens when a user in London enters the name of your domain in a browser:
    - DNS routes the query to a Route 53 name server.
    - Route 53 refers to its data on latency between London and the Singapore Region and between London and the Oregon Region.
    - If latency is lower between the London and Oregon Regions, Route 53 responds to the query with the IP address for the Oregon load balancer. If latency is lower between London and the Singapore Region, Route 53 responds with the IP address for the Singapore load balancer.
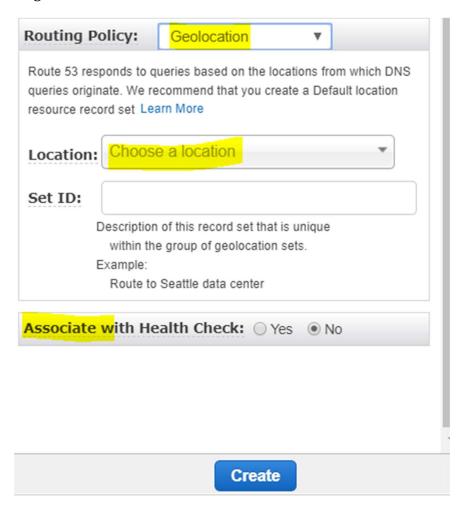
## Failover Routing Policy

- Failover routing lets you route traffic to a resource when the resource is healthy or to a different resource when the first resource is unhealthy.
- The primary and secondary records can route traffic to anything from an Amazon S3 bucket that is configured as a website to a complex tree of records.
- You can associate health check with this type of policy.

## Create Record Set

**Name:** `example` .gaurav.com.

**Type:** A – IPv4 address ▼

**Alias:** ○ Yes  ◉ No

**TTL (Seconds):** `300`  [1m] [5m] [1h] [1d]

**Value:**
```
52.34.56.78
```
IPv4 address. Enter multiple addresses
    on separate lines.
Example:
    192.0.2.235
    198.51.100.234

**Routing Policy:** Failover ▼

Route 53 responds to queries using primary record sets if any are healthy, or using secondary record sets otherwise. Learn More

**Failover Record Type:** ◉ Primary   ○ Secondary

**Set ID:** example-Primary

**Associate with Health Check:** ○ Yes  ◉ No

[ Create ]

# Geo Location Routing Policy

- Geolocation routing lets you choose the resources that serve your traffic based on the geographic location of your users, meaning the location that DNS queries originate from.
- For example, you might want all queries from Europe to be routed to an Elastic Load Balancing load balancer in the Frankfurt Region.
- When you use geolocation routing, you can localize your content and present some or all of your website in the language of your users.
- Geolocation routing lets you choose the resources that serve your traffic based on the geographic location of your users, meaning the location that DNS queries originate from.
- For example, you might want all queries from Europe to be routed to an ELB load balancer in the Frankfurt region.
- This is routing based on user location. Health check associated.

## Multi Value Routing Policy

- It helps distribute DNS responses across multiple resources.
- For example, use multivalue answer routing when you want to associate your routing records with a Route 53 health check.
- Use multivalue answer routing when you need to return multiple values for a DNS query and route traffic to multiple IP addresses.
- Up to 8 healthy records are returned for each Multi Value query.
- Multi Value is not a substitute for having an ELB.



## Simple Routing Policy

Use Case: The Simple routing policy is the most basic and straightforward routing policy. It is used when you want to route DNS traffic evenly to multiple resources.

- How it Works: With the Simple routing policy, you associate one or more DNS records with a domain name. When a DNS query is made for that domain, Route 53 responds with all applicable records in a random order, effectively distributing the traffic evenly among those resources.
- Example: Suppose you have multiple web servers, and you want to distribute incoming web traffic evenly among them. You can create multiple A records for your domain, each pointing to the IP address of one of your web servers. When users access your domain, they will be directed to one of the web servers in a round-robin fashion.

## Weighted Routing Policy

- Use Case: The Weighted routing policy is used when you want to route traffic to different resources based on assigned weights. This policy allows you to distribute traffic unevenly among resources, giving you fine-grained control over the traffic distribution.
- How it Works: With the Weighted routing policy, you assign a weight to each DNS record associated with your domain. The weight represents the proportion of traffic that should be sent to that resource. For example, if you have two records with weights 70 and 30, approximately 70% of the traffic will be directed to the resource with the higher weight.
- Example: Suppose you have two regions hosting your application, one in North America and one in Europe. You want to direct 70% of your traffic to the North American servers and 30% to the European servers. You can create two A records with weights 70 and 30, respectively.
- Flexibility: The Weighted routing policy gives you flexibility to perform various traffic routing scenarios, such as A/B testing, gradual resource rollout, or resource failover testing.
- Health Checks: You can also associate health checks with weighted records, allowing you to route traffic based on the health of the resources. If a resource fails a health check, Route 53 will reduce the traffic sent to it.
- In summary, the Simple routing policy evenly distributes traffic among resources, while the Weighted routing policy allows you to distribute traffic unevenly based on weights.

## Latency Routing Policy

- In Amazon Route 53, the Latency routing policy is a DNS routing policy that allows you to route traffic based on the lowest network.
- This policy is particularly useful when you want to direct users to the resource that provides the best performance or lowest response time for their location.

Here are the key points about the Latency routing policy:

- Geographic-Based Routing: Latency routing takes into account the geographic location of the DNS resolver (the system querying DNS for a domain). It uses the resolver's geographic location as a proxy for the location of the end user.
- Latency Measurements: Amazon Route 53 constantly monitors the latency between DNS resolvers around the world and the AWS resources or IP addresses associated with your DNS records.
- Routing Decision: When a DNS query is made for your domain, Route 53 responds by selecting the record associated with the AWS region or resource that offers the lowest latency for the querying resolver's location.

- Failover Support: You can configure health checks for the resources associated with latency records. If a resource becomes unhealthy, Route 53 will automatically route traffic away from it to healthy resources.
- Latency routing is most effective when you have resources or services located in multiple geographic regions or when latency differences between regions are significant.

## AWS Route53 Health Checks

- Amazon Route 53 Health Checks are a feature that allows you to monitor the health and availability of your AWS resources or other endpoints.
- Health checks help ensure that Route 53 routes traffic to only healthy and operational resources.

### Key Features and Concepts

- Resource Monitoring: Health checks can be associated with a specific resource, such as an AWS Elastic Load Balancer (ELB), Amazon EC2 instance, or an IP address of a resource.
- Health Check Types: Route 53 supports several health check types, including HTTP/HTTPS, TCP, and ICMP.
- You can choose the type that matches the protocol of your resource. For web-based resources, you can set up HTTP or HTTPS checks, which verify the resource's HTTP status codes.
- Advanced Configuration: Health checks offer advanced configuration options. You can specify the request interval, the number of consecutive failures required to mark a resource as unhealthy, and the threshold for marking a resource as healthy again.
- Regions: Health checks can be associated with specific AWS regions. This allows you to monitor resources from multiple geographic locations, ensuring accurate health status information.
- DNS Failover: Health checks are commonly used in DNS failover scenarios. If a resource fails its health check, Route 53 can automatically route traffic away from the unhealthy resource to a healthy one, helping ensure high availability.

## Steps to Set Up Health Checks**:

Here's a simplified overview of how to set up Health Checks in Route 53:

- Create a Health Check
    - In the Route 53 console, navigate to "Health checks."
    - Click "Create health check."
    - Choose the health check type, specify the target (e.g., endpoint or IP address), and configure advanced settings.

- Associate Health Check with Resources: Associate the health check with the resources you want to monitor. This can be done when creating or editing DNS records or when configuring a load balancer.
- Configure Failover Policies (Optional): If you want to set up DNS failover, configure routing policies based on health check status. For example, you can configure a primary and secondary resource, and if the primary resource fails its health check, Route 53 will automatically route traffic to the secondary resource.
- Monitor Health Checks: Regularly monitor the status of health checks in the Route 53 console. You can see the results of health checks, including the time of the last check and the status (healthy, unhealthy, etc.).

Summary: Route 53 Health Checks are a crucial tool for ensuring the reliability and availability of your infrastructure. They can help you detect and respond to issues with your resources quickly, improving the overall performance of your applications and services.

## Failover Routing Policy

- Use Case: Failover routing is used to create a high-availability setup where traffic is directed to a primary resource by default but can fail over to a secondary resource if the primary becomes unavailable.
- Configuration: You create two sets of resource records: one for the primary resource and one for the secondary resource. Each set has the same name but different values.
- You can associate health checks with each resource to monitor their availability.
- You configure a primary and secondary routing policy and set a failover condition. If the primary resource fails its health check, Route 53 automatically fails over to the secondary resource.
- Example: You have a website hosted on an Elastic Load Balancer (ELB) as the primary resource. You also have a backup static website hosted on an S3 bucket as the secondary resource. If the ELB becomes unhealthy (e.g., due to ELB issues or issues with the application servers behind it), Route 53 can automatically fail over to the S3 bucket, ensuring continued availability.

## GeoLocation Routing Policy

- Use Case: GeoLocation routing is used to route traffic based on the geographic location of the DNS resolver (the client making the DNS query).
- It enables you to provide region-specific content or direct users to the nearest data center.
- Configuration: You create resource record sets for different geographic regions or countries. For each record set, you specify the location and associated resources.
- You can use "Default" for unspecified regions or countries.

- Example: You have a global e-commerce application with servers in North America, Europe, and Asia. You use GeoLocation routing to direct users to the nearest server based on their location. North American users are directed to the North American server, European users to the European server, and so on.

## Geoproximity Routing Policy**:

- Use Case: Geoproximity routing is similar to GeoLocation routing but provides more granular control and allows for more sophisticated routing decisions.
- Configuration**: Geoproximity routing lets you define geographic regions using geographic coordinates (latitude and longitude) and associate resources with those regions. You can also configure "bias" to influence routing decisions.
- You can use "Traffic Bias" to prioritize certain resources when they are in proximity to the client's location.
- Example: You have a global application and want to ensure that users are directed to the nearest AWS region but also have a secondary region as a backup. You configure Geoproximity routing with bias to favor the nearest region but still allow failover to a secondary region if necessary.

## Multi-Value Answer routing policy

- Amazon Route 53's Multi-Value Answer routing policy is designed to provide simple load balancing and failover capabilities for multiple resources associated with a single DNS record.
- This routing policy is particularly useful when you have multiple resources (e.g., servers, endpoints) that should receive traffic in a random or round-robin fashion for load distribution or failover purposes. Here are the key points about the Multi-Value Answer routing policy:
- Multiple Resource Records**: With the Multi-Value Answer routing policy, you create multiple resource records (A, AAAA, or CNAME records) for the same DNS name. Each resource record is associated with a different resource or endpoint.
- DNS Responses**: When a DNS query is made for the domain, Route 53 responds with all the applicable resource records in a random order or in a round-robin fashion. This allows for basic load balancing across multiple resources.
- Health Checks**: You can associate health checks with each resource record to monitor the health and availability of the resources. If a resource becomes unhealthy, Route 53 will stop including it in the DNS responses.
- Use Case: You can use Multi-Value Answer routing to distribute traffic across multiple servers or endpoints to balance the load.

# AWS Route 53 Routing Policies

Amazon Route 53 offers several routing policies that allow you to control how DNS requests for your domain are routed to different endpoints. Here are the most common routing policies in Route 53:

1.  **Simple Routing Policy**
    - ✓ This is the most basic routing policy.
    - ✓ You can associate one or more DNS records (e.g., A or CNAME) with this policy.
    - ✓ All records are treated equally, and Route 53 responds with all applicable records in a random order when DNS queries are made.
    - ✓ Typically used when you want to distribute traffic evenly across multiple endpoints.

2.  **Weighted Routing Policy**
    - ✓ Allows you to assign weights to different DNS records.
    - ✓ The DNS resolver receives records based on their assigned weights.
    - ✓ Useful for load balancing or distributing traffic unevenly among resources.
    - ✓ For example, you can send 70% of traffic to one resource and 30% to another.

3.  **Latency-Based Routing Policy**
    - ✓ Routes traffic based on the lowest latency to the target endpoints.
    - ✓ Route 53 uses latency measurements to determine which endpoint is closest to the user, minimizing response time.
    - ✓ Ideal for applications with geographically distributed resources.

4.  **Failover Routing Policy**
    - ✓ Used in conjunction with health checks to route traffic to a standby or failover resource when the primary resource becomes unhealthy.
    - ✓ Commonly used for creating high availability and fault tolerance in applications.

5.  **Geolocation Routing Policy**
    - ✓ Routes traffic based on the geographic location of the DNS resolver making the request.
    - ✓ You can specify different resources for different regions or countries.
    - ✓ Useful for directing users to the nearest or most appropriate resource based on their location.

6.  **Geoproximity Routing Policy (Traffic Flow Only)**
    - ✓ Similar to geolocation routing but more advanced and available through Traffic Flow.
    - ✓ Allows you to route traffic based on geographic regions, but also consider the health of the resources in those regions.

7. **Multivalue Answer Routing Policy**
   - ✓ Similar to simple routing but allows you to return multiple IP addresses in response to DNS queries.
   - ✓ Each response includes a random selection of records, providing a basic form of load balancing.
   - ✓ Useful for distributing traffic across multiple endpoints without weighted routing.

# Content Delivery Network (CDN)

- ✓ A Content Delivery Network (CDN) is a distributed network of servers strategically located across the globe to speed up the delivery of web content to end-users.
- ✓ CDNs are designed to improve the speed, reliability, and performance of delivering web content, such as web pages, images, videos, and other static or dynamic assets, to users from various geographic locations.

Here are the key features and functions of CDNs:

1. **Content Caching:** CDNs cache (store) copies of web content in multiple data centers or edge server locations. This means that content is closer to end-users, reducing the latency and load times associated with retrieving data from a distant origin server.
2. **Geographic Distribution:** CDNs have a network of servers in various regions and cities worldwide. These servers, often referred to as "edge servers" or "PoPs" (Points of Presence), are strategically placed to reduce the physical distance between users and content, improving performance.
3. **Load Balancing:** CDNs distribute incoming traffic across multiple servers or locations to ensure that no single server becomes overloaded. This load balancing helps maintain optimal performance during traffic spikes or high-demand periods.
4. **Distributed Caching:** Cached content is distributed across multiple servers within the CDN, allowing for redundancy and fault tolerance. If one server becomes unavailable, another server can serve the cached content.
5. **Security:** Many CDNs offer security features such as DDoS (Distributed Denial of Service) protection, web application firewalls, and SSL/TLS (Secure Sockets Layer/Transport Layer Security) encryption to protect content and applications from online threats.
6. **Content Compression:** CDNs often use techniques like content compression and image optimization to reduce the size of assets, leading to faster load times and lower bandwidth consumption.

# Content Delivery Networks (CDNs) offer several advantages

- Improved Website Performance
  - Reduced Latency: CDNs cache content closer to end-users, reducing the physical distance data needs to travel. This results in faster loading times for web pages, reducing latency.
  - Faster Page Load: Caching static assets (images, scripts, stylesheets) on CDN edge servers means users can access these resources more quickly, leading to a better user experience.
- Global Reach
  - Geographic Distribution: CDNs have a network of servers in various regions and countries, allowing businesses to serve content to users worldwide with minimal latency.
  - Scalability: CDNs can handle traffic spikes and sudden increases in demand, ensuring that websites remain available even during high-traffic events like product launches or viral content.
- High Availability and Redundancy
  - Load Balancing: CDNs distribute traffic across multiple servers, balancing the load to prevent any single server from becoming overwhelmed.
  - Redundancy: CDNs often have multiple copies of content stored across various edge servers. If one server becomes unavailable, another can serve the content.
- Content Compression and Optimization
  - Content Compression: CDNs often compress text-based content (HTML, CSS, JavaScript) and optimize images to reduce the size of assets. This lowers bandwidth usage and accelerates content delivery.
  - Minimized Data Transfer Costs: Optimized content results in lower data transfer costs for both content providers and users.
- Security
  - DDoS Mitigation: Many CDNs offer Distributed Denial of Service (DDoS) protection, mitigating attacks by filtering malicious traffic and ensuring the availability of content.
  - Web Application Firewall (WAF): CDNs often include WAF capabilities to protect websites and applications from common web threats, such as SQL injection and cross-site scripting (XSS).

# AWS Cloudfront

- Amazon CloudFront is a highly scalable Content Delivery Network (CDN) service offered by Amazon Web Services (AWS).
- CloudFront is designed to accelerate the delivery of web content, including web pages, images, videos, etc. to users worldwide with low latency and high transfer speeds.

Here are some key aspects and features of AWS CloudFront:

- Global Network of Edge Servers: CloudFront operates a global network of edge servers, which are distributed in major cities and regions worldwide. These edge servers cache content and serve it to users from locations geographically closer to them, reducing latency.
- Low Latency and High Performance: CloudFront automatically routes user requests to the nearest edge location, minimizing the distance content needs to travel. The result is faster content delivery and improved website and application performance.
- Content Caching and Distribution: CloudFront caches and distributes both static and dynamic content, including HTML, CSS, JavaScript, images, videos, and APIs. Content can be cached for a specified duration or until it's invalidated or updated.
- Origin Fetching: CloudFront can fetch content from multiple origins, including Amazon S3 buckets, Amazon EC2 instances, on-premises servers etc.
- Security: CloudFront integrates with AWS Web Application Firewall (WAF) to protect against web application attacks and threats. It supports SSL/TLS encryption to provide secure content delivery to end-users.
- High Availability and Redundancy: CloudFront is designed for high availability, with multiple redundant edge locations. Content is automatically distributed across these locations, ensuring reliability and fault tolerance.

## AWS Cloudfront with S3

- Amazon CloudFront can be integrated with Amazon S3 (Simple Storage Service) to enhance the delivery, performance, and security of content stored in S3 buckets.
- This integration allows you to distribute your static and dynamic assets globally with low latency, high transfer speeds, and other CDN benefits.

Here's how you can use CloudFront with S3:

- Step 1: Set Up an S3 Bucket: If you don't already have one, create an Amazon S3 bucket to store your content (e.g., images, videos, HTML files). Ensure that the S3 bucket and the objects within it are configured to be publicly accessible or have appropriate permissions.

- Step 2: Create a CloudFront Distribution:
  - o 1. Log in to the AWS Management Console.
  - o Navigate to the Amazon CloudFront service.
  - o Click "Create Distribution."
  - o Select the "Web" distribution type, which is suitable for most content delivery scenarios.
- Step 3: Configure the CloudFront Distribution: In the configuration wizard, you'll need to specify several settings:
  - o Origin Settings:
    - Choose "S3 Bucket" as the origin type.
    - Select the S3 bucket that you want to use as the origin for your distribution.
  - o Cache Behavior Settings: Define how CloudFront should cache and serve content from your S3 bucket. You can customize cache behaviors based on URL patterns, query string parameters, and more.
  - o Distribution Settings: Configure additional settings such as the default root object (e.g., "index.html"), logging, and custom SSL certificates if necessary.
  - o Security and Restrictions: Set up security options such as SSL/TLS settings, access control, and geo-restriction if needed.
- Step 4: Review and Create the Distribution: Review your settings and then click "Create Distribution."
- Step 5: Wait for Distribution to Deploy: It may take some time (usually 15-30 minutes) for the CloudFront distribution to deploy globally.
- Step 6: Update DNS Records: Once the distribution is deployed, CloudFront provides a unique domain name (e.g., d123456789abcdef.cloudfront.net) for your distribution. Update your DNS records (e.g., CNAME) to point to this CloudFront domain name. This enables your users to access your content through CloudFront.
- Step 7: Test Your CloudFront Distribution: Access your content using the CloudFront domain name or custom domain (if configured). You should experience improved performance and content delivery.

## AWS Cloudfront Signed URL and Signed Cookies

- Amazon CloudFront provides two mechanisms for controlling access to your content: signed URLs and signed cookies.
- These mechanisms allow you to restrict who can access specific content through your CloudFront distribution.
- Both options are useful for secure content delivery, but they serve slightly different use cases:

- The choice between them depends on your specific use case and whether you require temporary, time-limited access (signed URLs) or session-based access to multiple resources (signed cookies).

## AWS Cloudfront Signed URLs

- Use Case: Signed URLs are typically used when you want to provide temporary access to specific resources. For example, you might generate a signed URL for a private video that allows a user to access the video for a limited time.
- How It Works: To generate a signed URL, you use your AWS credentials (Access Key and Secret Key) to create a time-limited URL that includes authentication information. This URL grants access to a specific resource for a limited duration.
- Configuration: You configure your CloudFront distribution to require signed URLs for access. You can specify the expiration time, allowing fine-grained control over access duration.
- Use Cases: temporary access to protected content.

## AWS Cloudfront Signed Cookies

- Use Case: Signed cookies are used when you want to provide authenticated access to multiple resources with a single authentication token.
- They are often used for scenarios where users need access to multiple protected resources during a single session.
- How It Works: Instead of embedding authentication information in the URL, signed cookies store authentication tokens in browser cookies. The cookies grant access to multiple resources during a session.
- Configuration: You configure your CloudFront distribution to require signed cookies for access. You can set up custom cookie behaviors for different paths and resources.
- Use Cases: Secure web applications, membership sites, scenarios where users need authenticated access to multiple resources.

# AWS Global Accelerator

- AWS Global Accelerator is a service provided by Amazon Web Services (AWS) that helps improve the availability and performance of applications by utilizing a global network infrastructure.
- It offers static IP addresses (Anycast) that act as entry points to your application hosted in one or more AWS regions.
- Global Accelerator directs traffic from users to the optimal AWS endpoint based on several factors, including health, routing policies, and geographic proximity.

- AWS Global Accelerator is suitable for applications that require high availability, low-latency global access, and fault tolerance.
- It's commonly used for global websites, multi-region applications, and scenarios where improving application performance and availability across regions is essential.

Here are some key features and benefits of AWS Global Accelerator:

- Global Anycast IP Addresses
  - AWS Global Accelerator provides a set of Anycast IP addresses that are globally distributed and highly available.
  - These IP addresses serve as entry points to your application.
- Load Balancing
  - Global Accelerator balances incoming traffic across multiple AWS endpoints, such as Elastic Load Balancers (ELBs), Application Load Balancers (ALBs), or AWS Elastic IP addresses.
  - It distributes traffic based on the health of endpoints to ensure high availability.
- Health Checking
  - Global Accelerator continuously monitors the health of your endpoints by sending health checks.
  - Unhealthy endpoints are automatically excluded from receiving traffic, ensuring that users are routed to healthy resources.
- Accelerated Traffic Routing
  - Traffic is routed to the nearest AWS endpoint based on geographic proximity to optimize latency.
  - AWS Global Accelerator uses the AWS global network to efficiently route traffic to the closest regional endpoint.
- Fault Tolerance
  - Global Accelerator offers fault tolerance by automatically rerouting traffic in case of an AWS endpoint failure.
  - This helps ensure that your application remains available even during AWS regional outages.

## **Serverless Computing**

- Serverless computing, often referred to as Function as a Service (FaaS), is a cloud computing model where cloud providers manage the underlying infrastructure, allowing developers to focus solely on writing and deploying code in the form of small, event-driven functions.
- In a serverless architecture, developers don't need to worry about provisioning or managing servers, scaling resources, or dealing with infrastructure maintenance.

- o Instead, they can concentrate on writing code to execute specific tasks or respond to events.

Here are some key characteristics and benefits of serverless computing:

- Event-Driven: Serverless functions are triggered by events. These events can include HTTP requests (API Gateway), changes to data in databases (e.g., Amazon DynamoDB, Azure Cosmos DB), file uploads (e.g., Amazon S3), or custom events generated by applications.
- Automatic Scaling: Cloud providers automatically scale the serverless environment to accommodate changes in traffic and the number of incoming events. Functions are instantiated as needed and can scale from zero to handle spikes in demand.
- Pay-as-You-Go Billing: Serverless computing follows a pay-as-you-go pricing model. You are billed based on the number of function executions and the time it takes for each execution. There are no charges for idle resources.
- No Server Management: Developers are relieved from managing server infrastructure, operating systems, or server provisioning. This allows them to focus on writing code and delivering functionality.
- Rapid Development: Serverless platforms provide a simplified development process. Developers can write functions in their preferred programming languages (e.g., Node.js, Python, Java, Go) and deploy them quickly.
- Automatic Scaling: Serverless platforms automatically manage the scaling of functions in response to incoming events. This eliminates the need for manual resource provisioning and scaling.
- High Availability: Serverless architectures are designed to be highly available and fault-tolerant. Functions are distributed across multiple data centers or availability zones to ensure reliability.
- Event-Driven Architecture: Serverless computing adapts event-driven architectures, making it suitable for applications that respond to various types of events, such as IoT data, user actions, or data updates.

## Serverless computing benefits

Here are some of the key advantages of serverless computing:

- Simplified Development: Developers can focus exclusively on writing code for specific functions or tasks without having to manage server infrastructure, operating systems, or provisioning.
- Auto Scaling: Serverless platforms automatically handle the scaling of functions based on incoming events or requests.

- Cost-Efficiency: Serverless computing follows a pay-as-you-go pricing model, where you are billed only for the actual compute resources used during function execution. There are no charges for idle resources.
- No Server Management: Developers are relieved from the burden of server management tasks, such as server provisioning, patching, and maintenance. This reduces operational overhead and frees up time for more productive tasks.
- High Availability and Fault Tolerance: Serverless architectures are designed to be highly available and fault-tolerant. Functions are distributed across multiple data centers or availability zones.
- Event-Driven: Serverless computing is well-suited for event-driven architectures, where functions respond to specific events or triggers. This makes it suitable for applications that process real-time data, respond to user actions, or react to changes in data.
- Rapid Development and Deployment: Developers can write and deploy serverless functions quickly.

## AWS Lambda

- Serverless computing in AWS is provided through a service called AWS Lambda, which allows you to run code without provisioning or managing servers.
- AWS Lambda is at the core of AWS's serverless offerings and enables developers to build and deploy applications with minimal operational overhead.

Here's an overview of AWS Lambda:

- AWS Lambda is a compute service that executes code in response to various triggers or events. You upload your code (written in languages like Node.js, Python, Java, Go, and more), define the trigger, and AWS Lambda takes care of scaling and executing the code.
- Triggers and Events: AWS Lambda functions are triggered by various events, such as HTTP requests, changes in data (e.g., Amazon S3 object uploads, DynamoDB updates), message queue events (e.g., AWS SQS), custom events (e.g., AWS CloudWatch Events), and more.
- Auto Scaling: AWS Lambda automatically scales your functions in response to incoming events. Functions can scale from zero to handle high concurrency and then scale back down when idle.
- Pay-as-You-Go Pricing: AWS Lambda follows a pay-as-you-go pricing model. You are billed based on the number of function executions and the compute time required for each execution. There are no charges for idle time.
- Integration with AWS Services: AWS Lambda integrates with various AWS services, allowing you to build serverless applications.

## AWS Lambda Limits – Hard & Soft

- Concurrent Executions: Default limits can range from a few hundred to thousands of concurrent executions per account, depending on your AWS service limits.
- Memory and CPU Allocation
  - Each Lambda function can be allocated a certain amount of memory (from 128 MB to 3008 MB in 64 MB increments)
  - CPU power is allocated in proportion to the selected memory size.
- Execution Timeout: The maximum execution timeout for a Lambda function is 900 seconds (15 minutes).
- Deployment Package Size: The deployment package (the compressed archive containing your function code and dependencies) has a size limit, typically 50 MB for deployment directly from the AWS Management Console or 250 MB when using the AWS CLI or SDK.
- Function Package Storage: AWS Lambda provides a temporary storage location for your function code and dependencies. The available storage is 512 MB.
- Number of Function Versions: You can create multiple versions of a Lambda function, but there is a soft limit on the total number of versions per function.
- Account-Level Limits: There are account-level limits on the total number of Lambda functions, function versions, and triggers across all regions within your AWS account.

## AWS Lambda@Edge

- AWS Lambda@Edge is a service provided by Amazon Web Services (AWS) that allows you to run AWS Lambda functions at the edge locations of the AWS CloudFront content delivery network (CDN).
- This enables you to execute code in response to viewer requests and customize the content delivery and behavior of your web applications globally, closer to the end-users.

Here are some key aspects and use cases for AWS Lambda@Edge:

- Edge Locations and CDN: AWS CloudFront is a globally distributed content delivery network with edge locations around the world. Lambda@Edge uses these edge locations to execute code closer to end-users, reducing latency.
- Event-Driven Execution: Lambda@Edge functions are triggered by various CloudFront events, including viewer requests, origin requests, origin responses, and viewer responses.
  - ✓ Viewer Request: Lambda@Edge functions triggered by viewer requests can inspect and modify incoming requests before they reach your origin server. Common use cases include authentication, authorization, and URL rewriting.

✓ Origin Request: Functions triggered by origin requests allow you to customize requests made to your origin server. You can modify headers, query parameters, or even change the origin server based on request characteristics.

✓ Origin Response: Lambda@Edge functions triggered by origin responses can modify responses received from the origin server. This is useful for modifying headers or cache control directives.

✓ Viewer Response: Functions triggered by viewer responses enable you to customize responses before they are sent to the end-user's browser. Use cases include adding security headers, optimizing content, or handling A/B testing.

## AWS API Gateway

▪ Amazon API Gateway is a fully managed service provided by Amazon Web Services (AWS) that allows you to create, publish, maintain, monitor, and secure APIs (Application Programming Interfaces) for your applications.

▪ It acts as a front-end service for your backend resources, enabling you to expose HTTP, WebSocket, and RESTful APIs to external clients, such as web and mobile applications.

Here are some key features and aspects of AWS API Gateway:

✓ API Creation and Management: AWS API Gateway allows you to create and define APIs. You can define endpoints, methods (GET, POST, PUT, DELETE, etc.), request/response transformations, and API models.

✓ API Deployment and Staging: You can deploy different versions or stages of your APIs (e.g., development, testing, production) and manage the lifecycle of your APIs.

✓ API Types: AWS API Gateway supports multiple API types, including HTTP APIs and REST APIs.

✓ Security and Authentication: API Gateway provides various security mechanisms, such as AWS Identity and Access Management (IAM) integration, API keys, Lambda authorizers, and integration with AWS Cognito for user authentication.

✓ Request and Response Transformations: You can transform requests and responses using mapping templates, allowing you to customize data formats and structures.

✓ API Testing: You can use the built-in testing tool in the AWS Management Console to test your API endpoints.

✓ Scaling and High Availability: AWS API Gateway automatically scales to handle high levels of traffic and provides high availability through multiple AWS regions.

## Creating and Configuring an API in AWS Gateway

Below is a simplified hands-on example of how to create a basic API using AWS API Gateway and a simple AWS Lambda function as a backend.

## Step 1: Create an AWS Lambda Function

1. Go to the AWS Lambda Console.
2. Click "Create Function."
3. Choose "Author from scratch."
4. Give your function a name, select a runtime (e.g., Node.js), and choose an existing or create a new execution role with basic Lambda permissions.
5. Click "Create function."
6. In the Function code section, replace the default code with a simple Node.js Lambda function, such as:

   def lambda_handler(event, context):

      return {

            'message': "Hello From Lambda"

      }

7. Click "Deploy" to save your Lambda function.

## Step 2: Create an API in AWS API Gateway

1. Go to the AWS API Gateway Console.
2. Click "Create API.
3. Choose "HTTP API" or "REST API," depending on your use case. HTTP API is a simpler option for most scenarios.
4. Click "Build" next to the "HTTP API" option.
5. Enter a name for your API, such as "MyAPI," and click "Create API."

## Step 3: Create a Resource and Method

1. In the API Gateway Console, select your API ("MyAPI")
2. In the left-hand menu, click "Routes."
3. Under "Routes," click "Create" to create a resource (e.g., "/hello").
4. With the new resource selected, click "Create" again to create a method (e.g., "GET").
5. Choose "Lambda Function" as the integration type.
6. Select the region where your Lambda function is deployed.
7. In the "Lambda Function" field, start typing the name of your Lambda function and select it from the dropdown.
8. Click "Create."

**Step 4: Deploy Your API**

1. In the API Gateway Console, select your API.
2. In the left-hand menu, click "Deployments."
3. Click "Create."
4. Enter a name for your deployment stage (e.g., "prod").
5. Click "Create."

**Step 5: Test Your API**

1. In the API Gateway Console, select your API.
2. In the left-hand menu, click "Stages."
3. Under "Stages," select your deployment stage (e.g., "prod").
4. Copy the "Invoke URL" provided at the top of the page.
5. Open a web browser or use a tool like `curl` or Postman to make an HTTP GET request to your API's endpoint (e.g., `https://your-api-id.execute-api.your-region.amazonaws.com/prod/hello`).

You should receive a response with the message "Hello from Lambda!" indicating that your API is correctly invoking your Lambda function.

## AWS Cognito

- Amazon Cognito is a managed authentication and identity service provided by Amazon Web Services (AWS).
- It simplifies the development of secure and scalable user authentication and user management features for web and mobile applications.
- Cognito offers a range of authentication methods and user directory options, making it suitable for a variety of application scenarios.

Here are some key features and components of AWS Cognito:

- User Pools:
  - ✓ User Pools are user directories that allow you to create and manage user identities (such as usernames and passwords) for your applications.
  - ✓ User Pools provide features for user registration, authentication, account recovery, and user profile management.
  - ✓ You can configure multi-factor authentication (MFA) and password policies to enhance security.
- Federated Identities
  - ✓ AWS Cognito Federated Identities, also known as Identity Pools, enable you to grant temporary, limited-privilege AWS credentials to users authenticated by

external identity providers (such as Amazon, Google, Facebook, or SAML-based identity providers).
- ✓ Identity Pools are often used to provide secure access to AWS resources on behalf of authenticated users.
- Multi-Factor Authentication (MFA): MFA can be enforced for added security. Users can use SMS, Time-based One-Time Password (TOTP), or other MFA methods.
- Cognito User Pools vs. Identity Pools: User Pools are for managing user identities, while Identity Pools are for granting access to AWS resources and services. Often, both are used together in an application.

## AWS DynamoDB

- Amazon DynamoDB is a fully managed NoSQL database service provided by Amazon Web Services (AWS).
- It is designed to provide high availability, scalability, and low-latency performance for applications that require a flexible and highly available database.

Here are some key features and concepts associated with AWS DynamoDB:

- NoSQL Database: DynamoDB is a NoSQL (Not Only SQL) database, which means it doesn't rely on a traditional relational database structure. It is well-suited for handling unstructured or semi-structured data.
- Managed Service: AWS takes care of the operational aspects of DynamoDB, such as hardware provisioning, setup, configuration, patching, and scaling.
- Scalability: DynamoDB can scale both vertically and horizontally to accommodate varying workloads.
- Data Model: DynamoDB supports various data models, including key-value and document data models. It is schema-less, allowing you to add or modify attributes without the need to define a fixed schema.
- Primary Key: Every item in DynamoDB must have a primary key, which consists of either a single attribute (Simple Primary Key) or a combination of two attributes (Composite Primary Key): a partition key and an optional sort key.
- Encryption and Security: DynamoDB supports encryption at rest and in transit. It integrates with AWS Identity and Access Management (IAM) for fine-grained access control.

Overall, DynamoDB is a powerful and flexible database service that can be used for a wide range of applications, from web and mobile apps to IoT and gaming platforms, where high availability and scalability are essential.

## Components of AWS DynamoDB

Here are the primary components of AWS DynamoDB:

1. **<u>Tables</u>**: Tables are the basic component in DynamoDB. They store your data and are schema-less, meaning each item in a table can have a different set of attributes. You define the primary key when creating a table, which consists of a partition key and an optional sort key.
2. **<u>Items</u>**: Items are the individual records within a table. Each item is uniquely identified by its primary key attributes. Items can have attributes, which are key-value pairs that hold data.
3. **<u>Attributes</u>**: Attributes are the data elements within items. An item's attributes can include simple data types like strings, numbers, or binary data, as well as more complex types like lists and maps.
4. **<u>Primary Key</u>**: The primary key is used to uniquely identify items within a table. It can be one of two types:
   a. **Partition Key**: A single attribute that is used to distribute data across multiple partitions for scalability.
   b. **Sort Key (Optional)**: When used in combination with the partition key, it creates a composite primary key. It allows you to query and sort data efficiently within a partition.
5. **<u>Encryption</u>**: DynamoDB supports encryption at rest and in transit. Data at rest can be encrypted using AWS Key Management Service (KMS) keys, and data in transit can be encrypted using SSL/TLS.
6. **<u>Security and Access Control:</u>** DynamoDB integrates with AWS Identity and Access Management (IAM) for fine-grained access control, allowing you to define who can perform operations on your tables and items.
7. **<u>Backup and Restore</u>**: DynamoDB offers automated backups and on-demand backups, allowing you to restore your table data to any point in time within the backup retention period.