

Aws full document

Aws Networking

Aws Ec2

Aws S3

Aws IAM

Aws RDS

Aws Cloudwatch

Aws Certificate manager

Aws Rout53

Aws Lambda overview

VEERA

AWS Networking

Overview of AWS Networking

Amazon Web Services (AWS) Networking encompasses a broad set of services and capabilities that allow you to build highly scalable, reliable, and secure networking environments in the cloud. These services are designed to provide connectivity between AWS resources, between AWS and your on-premises data centers, as well as to and from the Internet.

Key Components of AWS Networking

1. Amazon VPC (Virtual Private Cloud)
 - o Amazon VPC is the foundational networking service that allows you to create a logically isolated network within the AWS Cloud. VPC enables you to define your own IP address ranges, create subnets, route traffic between instances, and configure security settings like firewalls.
 - o With VPC, you can:
 - Create Subnets: Segregate your network into subnets (e.g., public and private subnets) to organize resources in a way that suits your application.
 - Private and Public IP Addresses: Assign private IPs to your EC2 instances and optionally, public IPs for internet-facing instances.
 - Route Tables: Control how traffic is routed within your VPC and to external networks.

- Security Groups and Network ACLs: Control inbound and outbound traffic to/from resources using security groups (stateful) and network ACLs (stateless).
- VPC Peering: Connect multiple VPCs to share resources and services between them.

2. Elastic IP (EIP)

- Elastic IP is a static, public IP address designed for dynamic cloud computing. You can associate an EIP with an EC2 instance or a network interface, providing a fixed address that doesn't change, even if you stop and start your instance.
- Useful for high availability and failover configurations, as you can remap the EIP to another instance in case of failure.

Key Networking Concepts in AWS

1. VPC Subnets:

- Subnets are segments of your VPC's IP address range, where you place your AWS resources like EC2 instances. You can create public subnets (which can route traffic to/from the internet) and private subnets (isolated from the internet).

2. Route Tables:

- A route table defines how traffic is routed within a VPC and to/from the internet. Each subnet in a VPC must be associated with a route table.

3. Security Groups and Network ACLs:

- Security Groups are virtual firewalls that control inbound and outbound traffic at the instance level. They are stateful, meaning that if you allow inbound traffic, the response is automatically allowed.

- Network ACLs (Access Control Lists) are stateless firewalls that apply at the subnet level. Unlike security groups, network ACLs require explicit rules for both inbound and outbound traffic.

4. NAT Gateway and NAT Instance:

- NAT Gateway allows instances in a private subnet to connect to the internet for software updates and other purposes, while preventing inbound internet traffic from directly accessing the private instances.
- NAT Instance is a similar concept, but with more flexibility and management responsibilities.

5. Elastic Network Interfaces (ENI):

- Elastic Network Interface is a virtual network interface that you can attach to an EC2 instance or Lambda function. It enables communication between AWS resources, allowing more advanced network configurations (e.g., multi-homed instances).

Security in AWS Networking

1. VPC Security:

- Security groups and network ACLs help protect your VPC and resources by controlling traffic at both the instance and subnet levels.
- VPC Flow Logs: Capture detailed logs of network traffic in and out of your VPC, which can be useful for troubleshooting, monitoring, and auditing.

2. Private Communication:

- PrivateLink and VPC Peering allow private communication between VPCs and AWS services without routing traffic over the public internet.

3. Encryption:

- Encryption is supported for all traffic that moves between VPC resources. You can enable encryption for services like Amazon EFS, S3, and CloudFront, as well as traffic between EC2 instances.
-

Benefits of AWS Networking

1. Scalability and Flexibility:

- AWS networking services are designed to scale with your applications. Whether you're deploying a small website or a large enterprise system, AWS can scale to meet the demand.

2. Global Reach:

- AWS offers global infrastructure, so you can deploy applications closer to your users with low-latency connections via services like CloudFront, Global Accelerator, and Direct Connect.

3. Security:

- AWS provides a range of security features, including encryption, firewalls, and DDoS protection (via AWS Shield), helping you secure your network and applications.

4. Reliability:

- AWS networking components like Elastic Load Balancing and AWS Direct Connect ensure high availability and fault tolerance for your applications.

5. Cost-Effectiveness:

- Pay-as-you-go pricing models allow you to optimize costs based on usage. AWS helps you minimize operational costs by abstracting much of the complexity involved in network management.

1. Create VPC Network

Multicloud Devops by veera nareshit

vpc

The screenshot shows the AWS search interface with the query 'vpc'. The results page displays various services under the 'Services' category. The 'VPC' service card is highlighted with a yellow arrow. The card includes the service icon, name ('VPC'), star rating, and description ('Isolated Cloud Resources'). Below the card, under 'Top features', are links to 'Your VPCs', 'Subnet', 'Route table', 'Internet gateway', and 'Egress-only internet gateways'. Other visible cards include 'AWS Firewall Manager', 'Detective', and 'Managed Services'.

The screenshot shows the VPC dashboard. At the top, there is a 'Create VPC' button highlighted with a purple arrow. The dashboard also displays 'Launch EC2 Instances' and a note: 'Note: Your Instances will launch in the US East region.' Below this, there are sections for 'Resources by Region' (listing 'VPCs' and 'NAT Gateways') and 'Service Health' (with a link to 'View complete service health details'). On the right side, there are 'Settings' and 'Additional Information' sections.

The screenshot shows the 'Create VPC' configuration page in the AWS Management Console. At the top, there's a navigation bar with the AWS logo, 'Services' (selected), a search bar, and a keyboard shortcut '[Alt+S]'. Below the navigation is a breadcrumb trail: 'VPC > Your VPCs > Create VPC'. The main title is 'Create VPC' with an 'Info' link. A descriptive text states: 'A VPC is an isolated portion of the AWS Cloud populated by AWS objects, such as Amazon EC2 instances.' The first section is 'VPC settings'.

Resources to create [Info](#)
Create only the VPC resource or the VPC and other networking resources.

VPC only VPC and more

Name tag - optional
Creates a tag with a key of 'Name' and a value that you specify.
my-vpc

IPv4 CIDR block [Info](#)
 IPv4 CIDR manual input IPAM-allocated IPv4 CIDR block
10.0.0.0/16
CIDR block size must be between /16 and /28.

IPv6 CIDR block [Info](#)
 No IPv6 CIDR block IPAM-allocated IPv6 CIDR block
 Amazon-provided IPv6 CIDR block
 IPv6 CIDR owned by me

Tenancy [Info](#)

[CloudShell](#) [Feedback](#)

Multicloud Devops by veera nareshit

The screenshot shows the AWS VPC creation process. The first part is a configuration wizard:

- IPv4 CIDR block:** IPv4 manual input (selected), CIDR: 10.0.0.0/16.
- IPv6 CIDR block:** No IPv6 CIDR block selected.
- Tenancy:** Default.

Tags: A tag named "Name" is added with value "my-vpc".

Create VPC button is highlighted with a blue checkmark.

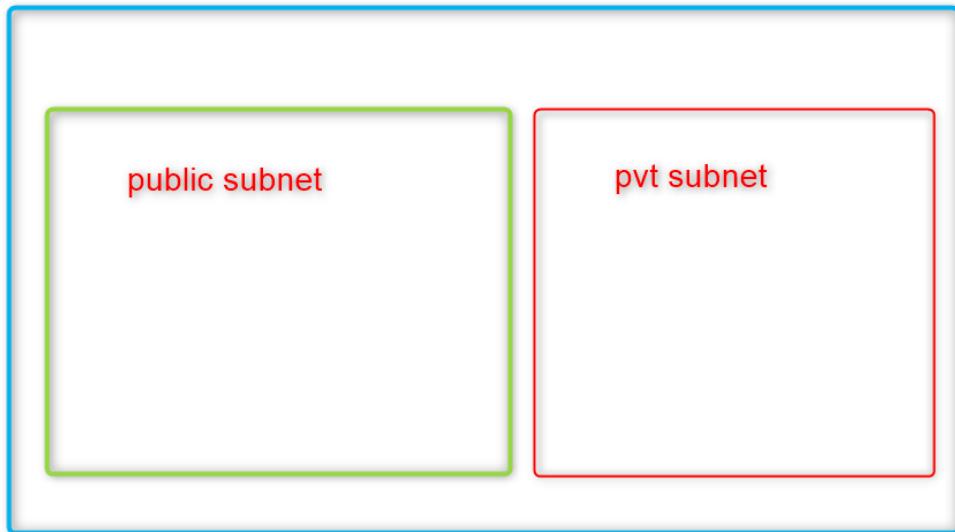
The second part is the VPC details page for the newly created VPC (ID: `vpc-02d126d121a842a1c`):

| Details | |
|---|---|
| VPC ID vpc-02d126d121a842a1c | State Available |
| Tenancy Default | DHCP option set dopt-0371dce67f9579a94 |
| Default VPC No | IPv4 CIDR 10.0.0.0/16 |
| Network Address Usage metrics Disabled | Route 53 Resolver DNS Firewall rule groups - |
| | DNS hostnames Disabled |
| | Main route table rtb-057b662b7075468ab |
| | IPv6 pool - |
| | DNS resolution Enabled |
| | Main network ACL acl-051f3b5a634d8dc80 |
| | IPv6 CIDR (Network border group) - |
| | Owner ID 944572535326 |

Resource map tab is selected.

Now vpc created.
Multicloud Devops by veera nareshit

2. Create Subnet (Public Subnet and Private Subnet)



Screenshot of the AWS VPC Subnets page:

Subnets (6) Info

| Name | Subnet ID | State | VPC |
|------|---------------------------|-----------|------------------------|
| - | subnet-0bc689d507bfa930b | Available | vpc-0f2e55bbcd8266ff60 |
| - | subnet-02003938f08dc7a51 | Available | vpc-0f2e55bbcd8266ff60 |
| - | subnet-035ec2aeb88662cd6 | Available | vpc-0f2e55bbcd8266ff60 |
| - | subnet-04530d4844be7ddff0 | Available | vpc-0f2e55bbcd8266ff60 |
| - | subnet-011da722ee126d6de | Available | vpc-0f2e55bbcd8266ff60 |
| - | subnet-004a0e8e42edb08a7 | Available | vpc-0f2e55bbcd8266ff60 |

Create subnet

Subnets (Selected)

Search:

Actions

N. Virginia

Whiz_User_183176.75957162 @ 9445-7255-5326

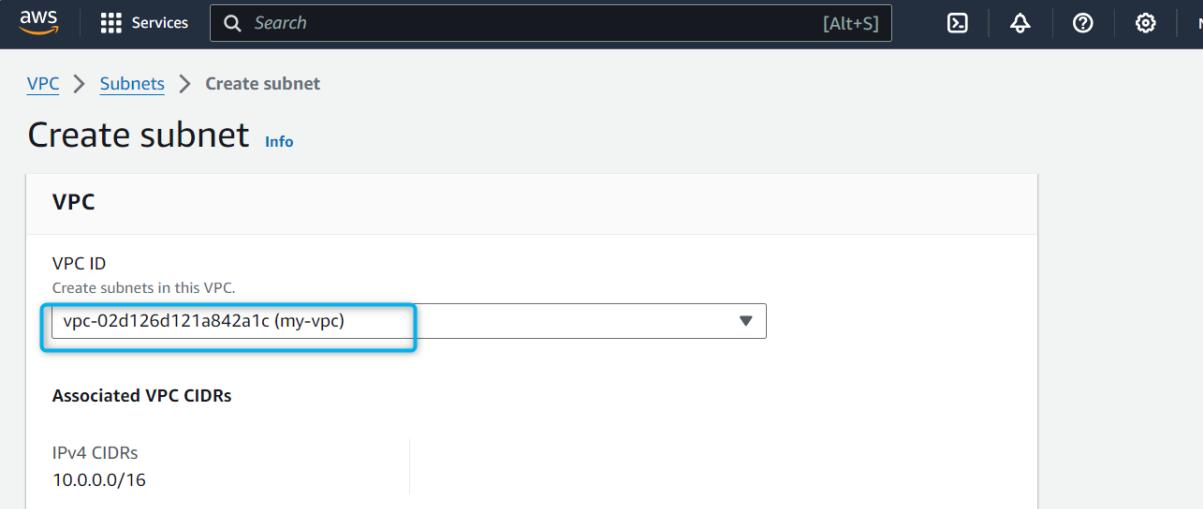
Subnets (6) Info

Find resources by attribute or tag

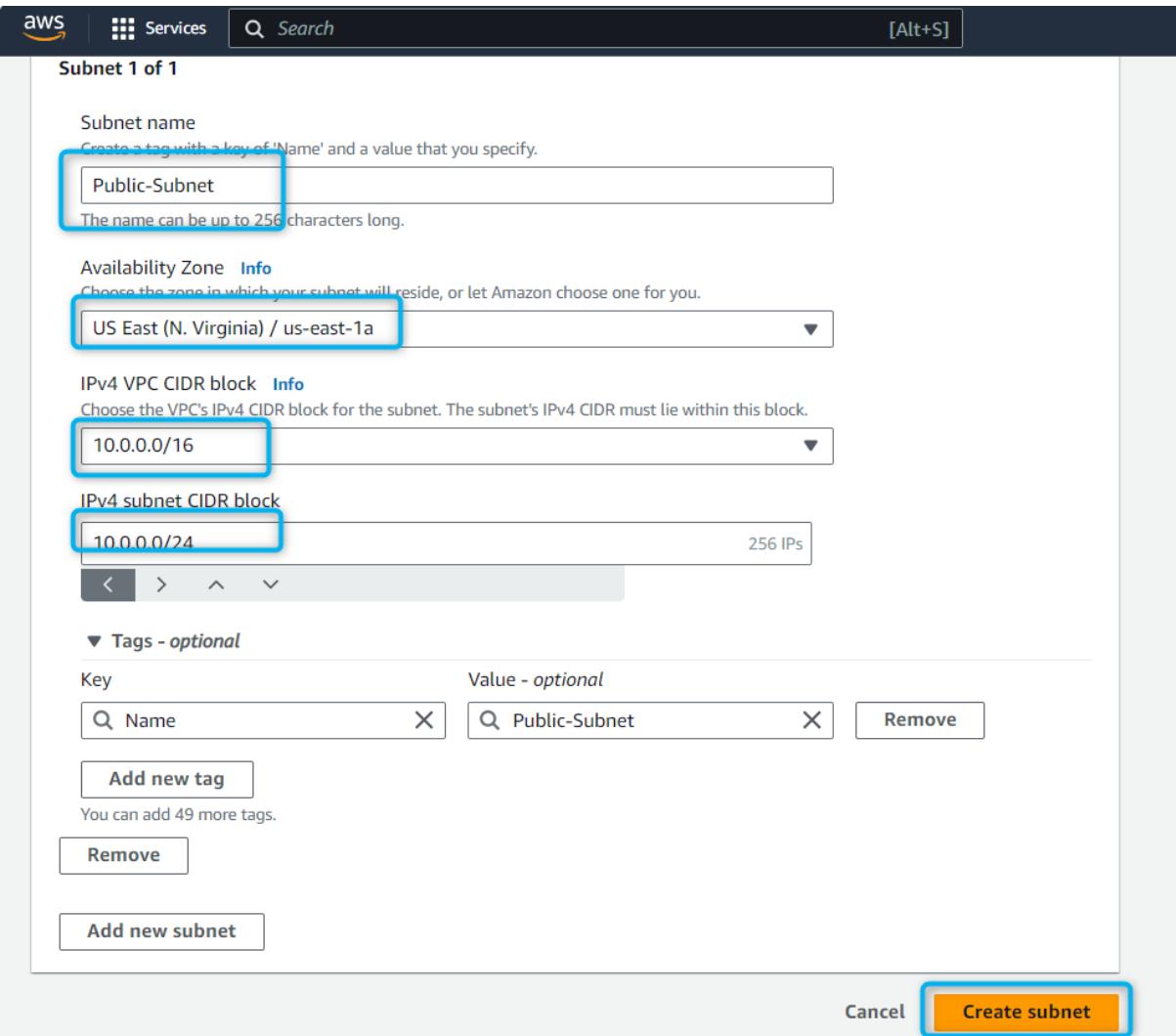
Select a subnet

VPC dashboard EC2 Global View Filter by VPC: Select a VPC Virtual private cloud Your VPCs Subnets Route tables Internet gateways Egress-only internet gateways Carrier gateways DHCP option sets Elastic IPs Managed prefix lists Endpoints Endpoint services NAT gateways Peering connections

Multicloud Devops by veera nareshit



The screenshot shows the 'Create subnet' wizard on the AWS VPC service. The first step, 'VPC', is selected. A dropdown menu under 'VPC ID' shows 'vpc-02d126d121a842a1c (my-vpc)'. Below it, 'Associated VPC CIDRs' lists 'IPv4 CIDRs' as '10.0.0.0/16'.



The screenshot shows the 'Subnet 1 of 1' configuration page. The 'Subnet name' field contains 'Public-Subnet'. The 'Availability Zone' dropdown is set to 'US East (N. Virginia) / us-east-1a'. The 'IPv4 VPC CIDR block' dropdown is set to '10.0.0.0/16'. The 'IPv4 subnet CIDR block' dropdown is set to '10.0.0.0/24' with '256 IPs'. Under 'Tags - optional', there is a key-value pair 'Name: Public-Subnet'. At the bottom, there are 'Cancel' and 'Create subnet' buttons, with 'Create subnet' highlighted with a blue border.

Multicloud Devops by veera nareshit

The screenshot shows two screenshots of the AWS VPC console. The top screenshot displays the 'Subnets' page with a success message: 'You have successfully created 1 subnet: subnet-095fdd16584af7df'. It lists one subnet named 'Public-Subnet' with details: Subnet ID: subnet-095fdd16584af7df, State: Available, VPC: my-vpc, IPv4 CIDR: 10.0.0.0/24. The bottom screenshot shows the 'Create subnet' wizard. The 'VPC' step is selected, showing the 'VPC ID' field with 'my-vpc' selected, highlighted with a blue border. A red annotation to the right of the field says 'for subnet 2 pvt same vpc'. Below the VPC section, the 'Associated VPC CIDRs' section shows 'IPv4 CIDRs: 10.0.0.0/16'.

Multicloud Devops by veera nareshit

Screenshot of the AWS Subnet creation interface:

Subnet name: Private-subnet (highlighted)

Availability Zone: US East (N. Virginia) / us-east-1b (highlighted)

IPv4 VPC CIDR block: 10.0.0.0/16 (highlighted)

IPv4 subnet CIDR block: 10.0.1.0/24 (highlighted)

Tags - optional:

| Key | Value - optional |
|------|------------------|
| Name | Private-subnet |

Create subnet button (highlighted)

Screenshot of the AWS Subnets list interface:

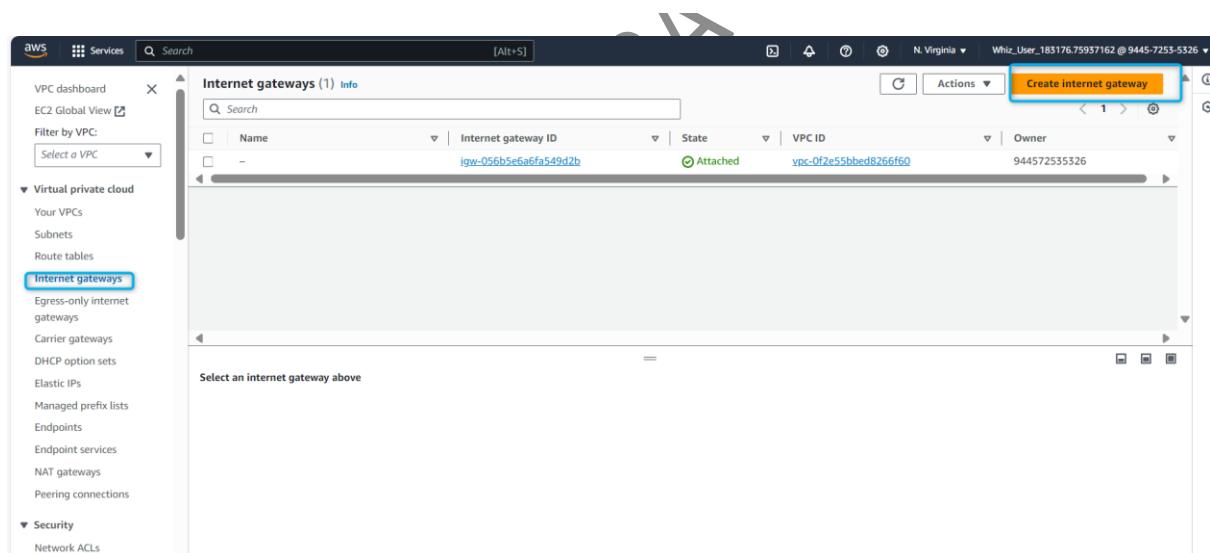
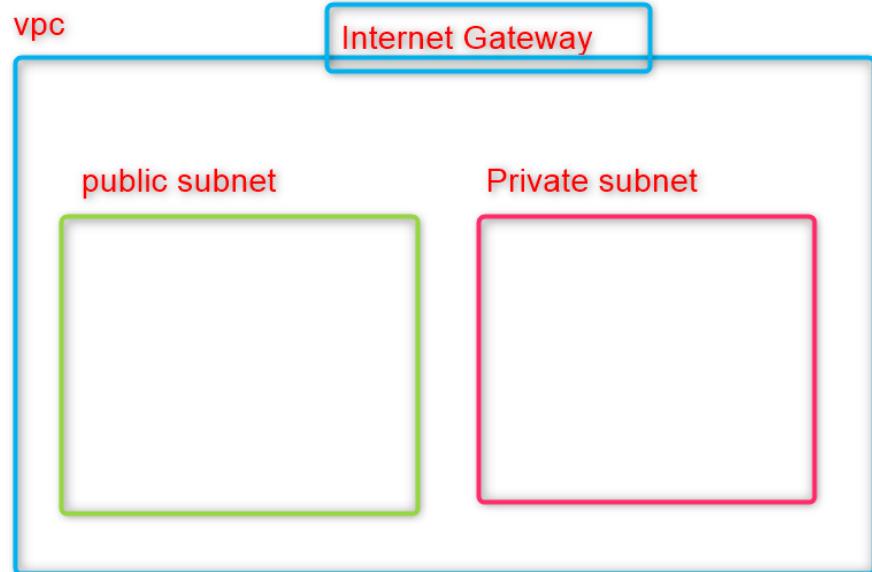
You have successfully created 1 subnet: subnet-005100391fb690199

| Name | Subnet ID | State | VPC | IPv4 CIDR |
|----------------|--------------------------|-----------|--------------------------------|----------------|
| - | subnet-02003938f08dc7a51 | Available | vpc-0f2e5bbbed8266f60 | 172.31.80.0/20 |
| - | subnet-035ec2ae88662cd6 | Available | vpc-0f2e5bbbed8266f60 | 172.31.48.0/20 |
| - | subnet-04530d4844be7df0 | Available | vpc-0f2e5bbbed8266f60 | 172.31.0.0/20 |
| - | subnet-011da72ee126d6de | Available | vpc-0f2e5bbbed8266f60 | 172.31.64.0/20 |
| - | subnet-004a0e8e42cd08a7 | Available | vpc-0f2e5bbbed8266f60 | 172.31.16.0/20 |
| Public-Subnet | subnet-095fdd16584af7df | Available | vpc-02d126d121a842a1c my-vpc | 10.0.0.0/24 |
| Private-subnet | subnet-005100391fb690199 | Available | vpc-02d126d121a842a1c my-vpc | 10.0.1.0/24 |

3. Create Internet Gateway

Multicloud Devops by veera nareshit

Multicloud Devops by veera nareshit



Multicloud Devops by veera nareshit

The screenshot shows the AWS VPC Internet Gateways page. On the left, there's a sidebar with options like 'Virtual private cloud', 'Internet gateways', and 'Egress-only internet gateways'. The main area displays an internet gateway with the ID 'igw-08acf94092139d318'. It shows the state as 'Detached' and the owner as '944572535326'. A 'Tags' section lists a single tag: 'Name' with the value 'Public Internet-Gateway'. There are 'Actions' buttons for managing the gateway.

Internet gateway created now attached this internet gateway to vpc.

This screenshot shows the same internet gateway page as above, but with a blue box highlighting the 'Actions' button in the top right corner. Below it, the 'Attach to VPC' option is also highlighted with a blue box. The rest of the interface is identical to the first screenshot.

This screenshot shows the 'Attach to VPC' dialog box. It has a search bar for selecting a VPC and a dropdown menu showing 'vpc-02d126d121a842a1c - my-vpc'. A purple checkmark is drawn next to this item. At the bottom right, there are 'Cancel' and 'Attach internet gateway' buttons, with the latter being highlighted with a yellow box and a purple checkmark drawn next to it.

Multicloud Devops by veera nareshit

VPC > Internet_gateways > igw-08acf94092139d318 / Public Internet-Gateway

igw-08acf94092139d318 / Public Internet-Gateway

Actions ▾

Details Info

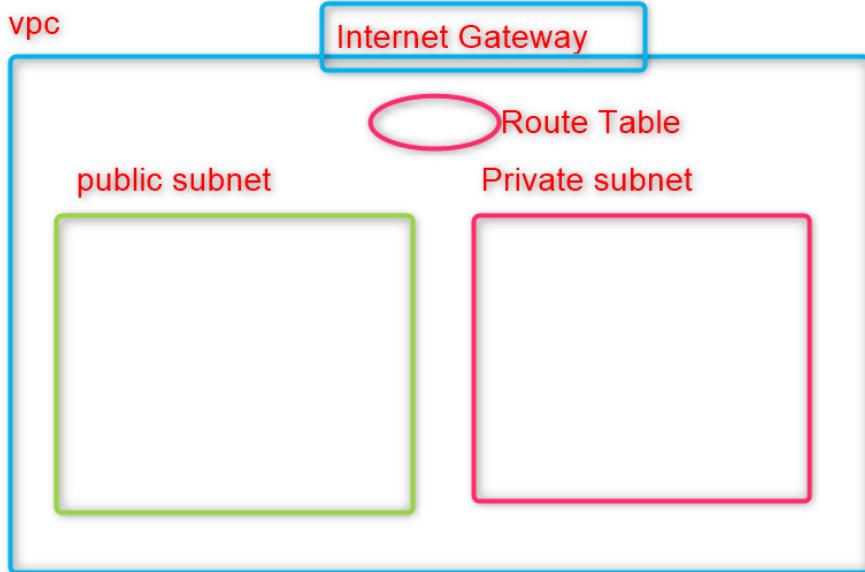
Internet gateway ID: igw-08acf94092139d318 State: Attached VPC ID: vpc-02d126d121a842a1c | my-vpc Owner: 944572535326

Tags

Search tags Manage tags

| Key | Value |
|------|-------------------------|
| Name | Public Internet-Gateway |

4.Create Route Table



Multicloud Devops by veera nareshit

The screenshot shows the AWS VPC dashboard. On the left, there's a sidebar with options like 'Your VPCs', 'Subnets', and 'Route tables'. The 'Route tables' option is selected and highlighted with a blue box. In the main pane, the title 'Route tables (1) Info' is displayed above a table. The table has columns for Name, Route table ID, Explicit subnet associ..., Edge associations, Main, and VPC. One row is listed with a name of '-' and a Route table ID of 'rtb-0a1c98f8e0e115256'. The 'Main' column shows 'Yes' and the 'VPC' column shows 'vpc-0f2e55bbbed8266f60'. At the top right of the main pane, there's a 'Create route table' button, which is also highlighted with a blue box.

The screenshot shows the 'Create route table' wizard. The first step, 'Route table settings', is visible. It includes fields for 'Name - optional' (containing 'Route-Table') and 'VPC' (containing 'vpc-02d126d121a842a1c (my-vpc)'). Both of these fields are highlighted with blue boxes. Below this, there's a 'Tags' section with a table for adding tags. A single tag is present: 'Name' (Key) and 'Route-Table' (Value). An 'Add new tag' button is also visible. At the bottom right of the step, there's a 'Create route table' button, which is also highlighted with a blue box.

Multicloud Devops by veera nareshit

Route table rtb-0c7bace0bd77d39de | Route-Table was created successfully.

VPC > Route tables > rtb-0c7bace0bd77d39de

rtb-0c7bace0bd77d39de / Route-Table

Actions ▾

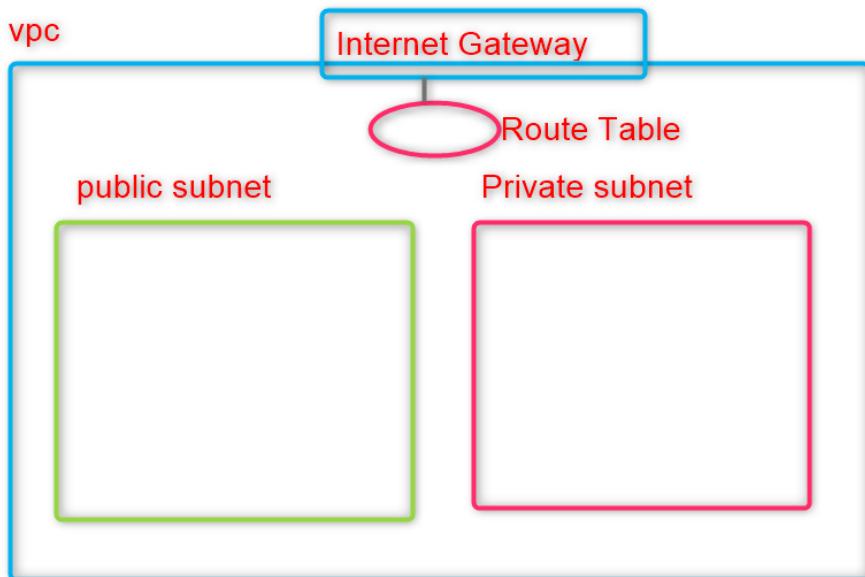
| Details | | Info | |
|------------------------------|--------------------------------|----------|--------------|
| Route table ID | rtb-0c7bace0bd77d39de | Main | No |
| VPC | vpc-02d126d121a842a1c my-vpc | Owner ID | 944572535326 |
| Explicit subnet associations | | | |
| Edge associations | | | |

Routes Subnet associations Edge associations Route propagation Tags

Routes (1)

| Filter routes | | | |
|---------------|--------|--------|------------|
| Destination | Target | Status | Propagated |
| 10.0.0.0/16 | local | Active | No |

- Attach this route table to the internet Gateway.



Multicloud Devops by veera nareshit

The image consists of three vertically stacked screenshots of the AWS VPC Route Tables interface.

Screenshot 1: Route Table Details

This screenshot shows the details of a route table named "rtb-0c7bace0bd77d39de". The "Routes" tab is selected, displaying one route entry:

| Destination | Target | Status | Propagated |
|-------------|--------|--------|------------|
| 10.0.0.0/16 | local | Active | No |

A blue box highlights the "Edit routes" button in the top right corner of the routes table.

Screenshot 2: Edit Routes - Local Target

This screenshot shows the "Edit routes" page for the same route table. A new route entry is being added:

| Destination | Target | Status | Propagated |
|--------------|------------------|--------|------------|
| 10.0.0.0/16 | local | Active | No |
| Q. 0.0.0.0/0 | Internet Gateway | - | No |

A blue box highlights the "Add route" button in the top left. The "Save changes" button is highlighted with a yellow box in the bottom right.

Screenshot 3: Edit Routes - Internet Gateway Target

This screenshot shows the "Edit routes" page after adding the Internet Gateway target. The route table now contains two entries:

| Destination | Target | Status | Propagated |
|--------------------------|---|------------------------------|------------|
| 10.0.0.0/16 | local | Active | No |
| Q. 0.0.0.0/0 | Internet Gateway | - | No |
| Q. igw-08acf94092139d318 | igw-08acf94092139d318 (Public Internet-Gateway) | Use: "igw-08acf94092139d318" | No |

A blue box highlights the "Save changes" button in the bottom right.

Multicloud Devops by veera nareshit

⌚ Updated routes for rtb-0c7bace0bd77d39de / Route-Table successfully

► Details

VPC > Route tables > rtb-0c7bace0bd77d39de

rtb-0c7bace0bd77d39de / Route-Table

Actions ▾

Details Info

| | | | |
|---|--------------------------|--|------------------------|
| Route table ID rtb-0c7bace0bd77d39de | Main No | Explicit subnet associations - | Edge associations - |
| VPC vpc-02d126d121a842a1c my-vpc | Owner ID 944572535326 | My Internet G/w now connected to Route Table | |

Routes Subnet associations Edge associations Route propagation Tags

Routes (2)

| Destination | Target | Status | Propagated |
|-------------|-----------------------|--------|------------|
| 0.0.0.0/0 | igw-08acf94092139d318 | Active | No |
| 10.0.0.0/16 | local | Active | No |

■ Now let me rename my Route Table name to Public Route table.

Route tables (1/3) Info

Find resources by attribute or tag

| Name | Route table ID | Exp |
|---|-----------------------|-----|
| - | rtb-0a1c98f8e0e115256 | - |
| - | rtb-057b662b7075468ab | - |
| <input checked="" type="checkbox"/> Route-Table | rtb-0c7bace0bd77d39de | - |

Edit Name

Route-Table Public-Route-Table

Cancel Save

Route tables (1/3) Info

Find resources by attribute or tag

| Name | Route table ID |
|--|-----------------------|
| - | rtb-0a1c98f8e0e115256 |
| - | rtb-057b662b7075468ab |
| <input checked="" type="checkbox"/> Public-Route-Table | rtb-0c7bace0bd77d39de |

Multicloud Devops by veera nareshit

- Now as my Public Route Table is attached with internet gateway and getting internet into it. I can attach my Public Route Table with my Public Subnet so that my Public Subnet will also able to get Internet Access.

Route tables (1/3) Info

| Name | Route table ID | Explicit subnet associ... | Edge associations | Main | VPC | Ow... |
|---------------------------|------------------------------|---------------------------|-------------------|------|-------------------------------|-------|
| - | rtb-0a1c98f8e0e115256 | - | - | Yes | vpc-0f2e5bbbed8266f60 | 944 |
| - | rtb-057b662b7075468ab | - | - | Yes | vpc-02d126d121a842a1c my... | 944 |
| Public-Route-Table | rtb-0c7bace0bd77d39de | - | - | No | vpc-02d126d121a842a1c my... | 944 |

rtb-0c7bace0bd77d39de / Public-Route-Table

Details | Routes | **Subnet associations** | Edge associations | Route propagation | Tags

Explicit subnet associations (0)

No subnet associations

You do not have any subnet associations.

Subnets without explicit associations (2)

The following subnets have not been explicitly associated with any route tables and are therefore associated with the main route table:

| Name | Subnet ID | IPv4 CIDR | IPv6 CIDR |
|----------------|--------------------------|-------------|-----------|
| Public-Subnet | subnet-095fdd16584af7df | 10.0.0.0/24 | - |
| Private-subnet | subnet-005100391fb690199 | 10.0.1.0/24 | - |

Edit subnet associations

Change which subnets are associated with this route table.

Available subnets (1/2)

| Name | Subnet ID | IPv4 CIDR | IPv6 CIDR | Route table ID |
|---|--------------------------|-------------|-----------|------------------------------|
| <input checked="" type="checkbox"/> Public-Subnet | subnet-095fdd16584af7df | 10.0.0.0/24 | - | Main (rtb-057b662b7075468ab) |
| <input type="checkbox"/> Private-subnet | subnet-005100391fb690199 | 10.0.1.0/24 | - | Main (rtb-057b662b7075468ab) |

Selected subnets

subnet-095fdd16584af7df / Public-Subnet X

Save associations

Multicloud Devops by veera nareshit

The screenshot shows the AWS Route Tables page. At the top, a green success message box says "You have successfully updated subnet associations for rtb-0c7bace0bd77d39de / Public-Route-Table." Below this is a table titled "Route tables (1/3) Info". The table has columns: Name, Route table ID, Explicit subnet associa..., Edge associations, Main, VPC, and Owner. There are three rows: one row with "-" values, one row for "rtb-0a1c98f8e0e115256" associated with "vpc-Of2e55bbcd8266f60" (Owner 944), and one row for "rtb-057b662b7075468ab" associated with "vpc-02d126d121a842a1c | my..." (Owner 944). A checkbox next to "Public-Route-Table" is checked. The "rtb-0c7bace0bd77d39de / Public-Route-Table" details page is shown below, with tabs for Details, Routes, Subnet associations, Edge associations, Route propagation, and Tags. The Details tab is selected, showing sections for Route table ID, Main, Explicit subnet associations, and Edge associations.

So my public subnet can able to get internet from Route table.

Now we will use this same above custom network to connect our EC2 instance with internet so that we can able to ping google.com with help of internet into ec2 server after configuring this custom network.

Amazon Web Services - Security Group vs NACL

Security groups and NACL both act as virtual firewalls which control the traffic from Inbound and Outbound. In this article, we will discuss the difference between Security Groups and NACL on Amazon Web Services.

Security Group:

Security groups are virtual shields or protectors of EC2 instances. Unless specifically allowed By default all Inbound traffic is blocked whereas all Outbound traffic is allowed from the Instance.

We can edit inbound and outbound rules after creating the Security Group. Here is an example of default outbound rules which allow all traffic for all protocols.

It is virtual firewall for your EC2 instances to control incoming and outgoing traffic

Multicloud Devops by veera nareshit

Multicloud Devops by veera nareshit

Here, we are adding inbound rules for protocol SSH with the default port of 22 for our current IP address here.

The screenshot shows the 'Inbound rules' section of an AWS Security Group. It lists one rule:

| Security group rule ID | Type | Protocol | Port range | Source | Description - optional |
|------------------------|------|----------|------------|--------|------------------------|
| sgr-035c64c7dd9a4f66a | SSH | TCP | 22 | Custom | Abhilash |

A red box highlights the 'Delete' button next to the rule.

1. Above picture is for Inbound rule

In SG we can block traffic at inbound level only

The screenshot shows the 'Edit outbound rules' section of an AWS Security Group. It lists one rule:

| Security group rule ID | Type | Protocol | Port range | Destination | Description - optional |
|------------------------|-------------|----------|------------|-------------|------------------------|
| sgr-05cbdfad39ac263e1 | All traffic | All | All | Custom | 0.0.0.0/0 |

A red box highlights the 'Delete' button next to the rule.

2. Outbound

The above picture there is no provision to block outbound rules so in SG there by default outbound allowed

In the security groups, we cannot block a specific IP address because it doesn't have any DENY rule just like ALLOW rule

Limits of Security Groups :

For a specific Security Group, the maximum Inbound and Outbound rules is 60

Network Access Control List (NACL):

Network Access Control List is also a virtual firewall for subnets, which controls the Inbound and Outbound traffic of Subnets. After the creation of VPC, a Default NACL will be associated and allow all Inbound and Outbound Traffic.

In NACL just like Security Groups, it contains set of Inbound and Outbound Rules , that can either allow or deny Traffic into or out of subnets. Since we have option to allow or deny traffic the order of the rules becomes important so that AWS uses a concept of rule number.

The screenshot shows the AWS Network ACLs console. On the left, there's a navigation sidebar with various VPC-related options like Elastic IPs, Managed Prefix Lists, Endpoints, Endpoint Services, NAT Gateways, Peering Connections, and several collapsed sections under SECURITY, NETWORK ANALYSIS, DNS FIREWALL, and NETWORK FIREWALL. The main area is titled "Network ACLs (1/1) Info". It displays a table with one row: Name: acl-03eb49085424761..., Network ACL ID: acl-03eb49085424761..., Associated with: 3 Subnets, Default: Yes, and VPC ID: vpc-0d2617096b762eedc. Below this table is a message: "You can now check network connectivity with Reachability Analyzer" with buttons for "Run Reachability Analyzer" and "X". At the bottom, there's a section titled "Inbound rules (2)" with a table showing two rules:

| Rule number | Type | Protocol | Port range | Source | Allow/Deny |
|-------------|-------------|----------|------------|-----------|------------|
| 100 | All traffic | All | All | 0.0.0.0/0 | Allow |
| * | All traffic | All | All | 0.0.0.0/0 | Deny |

1. Inbound rule

The screenshot shows the AWS Network ACLs console. On the left, there's a navigation sidebar with options like New VPC Experience, Elastic IPs, Managed Prefix Lists, Endpoints, Endpoint Services, NAT Gateways, Peering Connections, SECURITY (Network ACLs selected), NETWORK ANALYSIS (Reachability Analyzer), and DNS FIREWALL (Rule Groups, Domain Lists). The main area displays 'Network ACLs (1/1)' with a table. The table has columns: Name, Network ACL ID, Associated with, Default, and VPC ID. One row is listed: Name is 'acl-03eb49085424761...', Network ACL ID is 'acl-03eb49085424761...', Associated with is '3 Subnets', Default is 'Yes', and VPC ID is 'vpc-0d2617096b762eedc'. Below the table, a message says 'You can now check network connectivity with Reachability Analyzer' with a 'Run Reachability Analyzer' button. Under 'Outbound rules (2)', there's another table with columns: Rule number, Type, Protocol, Port range, Destination, and Allow/Deny. The first rule (rule 100) has Type 'All traffic', Protocol 'All', Port range 'All', Destination '0.0.0.0/0', and Allow/Deny 'Allow'. The second rule '*' has Type 'All traffic', Protocol 'All', Port range 'All', Destination '0.0.0.0/0', and Allow/Deny 'Deny'.

2. Outbound Rule

Amazon EC2 (Elastic Compute Cloud)

- Introduction and Purpose of having EC2
- Instances
- Instance Types
- Regions and Availability Zones
- Purchasing Options
- AMI
- Key Pairs
- Launch Templates
- Placement Groups
- Elastic Block Store (EBS) and Utilizing
- Snapshots
- Elastic IPs
- Load Balancers
- Elastic Load Balancer (ELB)
- Target Groups

CloudWatch Monitoring

Auto Scaling Groups (ASG)

A Brief Introduction to EC2 Instance:

An EC2 instance, short for Elastic Compute Cloud instance, is a virtual server offered by Amazon Web Services (AWS) that provides resizable compute capacity in the cloud. Essentially, it allows users

to rent virtual servers to run their applications or host their websites without the need to invest in physical hardware.

- EC2 instances offer a variety of configurations for CPU, memory, storage, and networking, catering to diverse workload requirements.
- Users can select from optimized instance types, such as general-purpose, memory-intensive, or high-performance computing, to suit their specific needs.
- Deployment and management are simplified through AWS Management Console, CLI, or SDKs, facilitating easy provisioning, monitoring, and scaling of instances.
- EC2 serves a wide range of use cases, from basic web hosting to complex data analytics, meeting the computing requirements of businesses across industries and sizes.
- EC2 follows a pay-as-you-go pricing model, ensuring cost-effectiveness by billing users based on instance type, usage duration, and other factors, thus accommodating varying business needs and budget constraints.

Purpose of having an EC2 Instance

- **Flexibility:** With a wide range of instance types and configurations, EC2 instances cater to diverse workload requirements, from basic web hosting to high-performance computing and data analytics.
- **Scalability:** EC2 instances offer scalable compute capacity, allowing businesses to easily adjust resources based on demand fluctuations. This scalability ensures efficient resource utilization and cost-effectiveness.
- **Global Reach:** With AWS's global infrastructure, EC2 instances are available in multiple regions worldwide. This enables businesses to deploy applications closer to their target audience for improved performance and latency optimization.
- **Security and Reliability:** EC2 instances provide a secure and reliable environment for running applications, with features like

built-in encryption, firewalls, and monitoring tools to safeguard data and ensure high availability.

- Ease of Management: EC2 instances can be quickly deployed and managed through AWS Management Console, CLI, or SDKs, streamlining the process of provisioning, monitoring, and scaling instances as needed.

Amazon EC2 offers a variety of instance types tailored for diverse needs. These types feature different configurations of CPU, memory, storage, and networking capacity, providing flexibility in selecting resources for your applications. Additionally, each instance type includes multiple sizes, enabling easy scaling according to your workload's demands.

EC2 Instance Types Overview:

General Purpose:

- These instances offer a balanced mix of compute, memory, and networking resources.
- Ideal for applications like web servers and code repositories that require equal proportions of these resources.

Burstable Performance:

- Also known as the T instance family, these instances provide baseline CPU performance with the ability to burst above it when needed.
- Refer to the Amazon EC2 User Guide for Linux Instances for more details.

Compute Optimized:

- Designed for compute-intensive applications requiring high-performance processors.
- Suitable for tasks such as batch processing, media transcoding, and high-performance web servers.

Memory Optimized:

- Tailored for workloads processing large datasets in memory.
- Deliver fast performance for applications with substantial memory requirements.

Storage Optimized:

- Intended for workloads demanding high, sequential read and write access to large data sets on local storage.

- Optimized to provide high IOPS for applications requiring rapid random I/O operations.

Accelerated Computing:

- Utilizes hardware accelerators or co-processors for efficient processing of functions like floating-point calculations and graphics processing.
- Enhances performance compared to software-based processing on CPUs.

High-Performance Computing (HPC):

- Purpose-built for running HPC workloads at scale on AWS, offering optimal priceperformance.
- Ideal for complex simulations and deep learning tasks requiring high-performance processors.

Previous Generation:

- AWS continues to support previous generation instance types for users who have optimized their applications around them.
- While current generation instance types offer better performance, previous generation types are still available for legacy applications.

Understanding AWS Regions, Availability Zones, and Deployment Strategies

Regions:

- AWS divides its global infrastructure into geographical regions, such as North America, Europe, Asia Pacific, etc.
- Each region is completely isolated from others and consists of multiple Availability Zones.
- Users select a region when deploying resources like EC2 instances. The choice of region affects latency, compliance, and data residency.

Availability Zones (AZs):

- Availability Zones are distinct data centers within a region, designed for fault tolerance and high availability.
- They are physically separate from each other and are connected through highspeed, low-latency links.
- Deploying resources across multiple AZs enhances fault tolerance. If one AZ experiences an issue, applications can continue running in others without interruption.

Deployment Strategy:

- To ensure high availability and fault tolerance, it's recommended to distribute resources (like EC2 instances) across multiple AZs within a region.
- This strategy protects against failures or outages in a single AZ and ensures that applications remain accessible even if one AZ becomes unavailable.

Scalability and Performance:

- Deploying resources across multiple AZs allows for better scalability and performance.
- Load balancers can distribute incoming traffic evenly across multiple AZs, optimizing resource utilization and improving response times for users.

Data Residency and Compliance:

- Different regions may have varying data residency regulations and compliance requirements.
- Users need to select the appropriate region and AZs to ensure compliance with local regulations and keep data within designated geographic boundaries.

Cost Considerations:

- Costs for AWS services may vary across regions, so users should consider pricing differences when selecting a region for deployment.
- Additionally, data transfer costs may apply when transferring data between regions or AZs.

EC2 Major Purchasing options:

1. **On Demand Instances** – Pay for the compute capacity per second , no up-front payment. We can decide when to launch, terminate, start, stop, or reboot.

Use cases: short-term work-loads , irregular workloads that cannot be interrupted.

2. **Reserved Instances** – We can reserve instances in terms of 1 year and 3 years. Provide up to 70% discount compared to On-Demand. Payment method can be all upfront, partial upfront, no

upfront. If you purchase Reserved Instances of a specific type in an availability zone and the same type is available in same zone then automatically charge at lower rate.

Use cases — Long-term workloads , Applications with steady state or predictable usage.

3. **Spot Instances** — Most cost effective , We can request for unused instances up to 90% discount. But if someone is paying more than of yours, then instance will get loosed. So it is less reliable.

Use cases — Image processing , non-critical applications.

4. **Dedicated Hosts**— Physical server completely dedicated for your use. More expensive and can be purchase as on demand or reserved. Billing will be per-host.

Use cases — Companies that have strong regulatory or compliance needs.

5.**Dedicated Instances** — Instances running on hardware that's dedicated to a single user. It can share hardware with other instances in same account. Billing will be per-instance.

6.**Savings Plans** — This will assist the user to lower the EC2 costs if they agree to keep it for a specific period – ideally, one or three years. The number of active instances are unrestricted.

AMI (Amazon Machine Image):

- An AMI is a template used to create virtual machines (instances) within the Amazon EC2 service.
- It contains the operating system, application server, and applications necessary to launch an instance.
- You can create custom AMIs or use pre-built AMIs provided by AWS and other users.

AMI Images:

- AMI images are the actual files stored in Amazon S3 that represent the templates used to launch EC2 instances.
- These images include the root volume snapshot and metadata required to launch an instance.

- AMI images are versioned, allowing you to maintain multiple versions of an AMI and roll back to previous versions if needed.

Placement Groups

- Placement groups are the way of logically grouping interdependent instances together in a selected region.
- In other words, instances that are existing within a common availability zone can be grouped under placement group in order to suit workload requirements.
- By using the placement group, we will increase the performance or improve the availability.

Types

- Cluster placement group
- Spread placement group
- Partition placement group

Benefits

- Can be able to launch multiple instances within the same AZ.
- Low network latency.
- High network throughput.

Limitations

- Cannot merge placement groups.
- Cannot launch dedicated hosts in placement groups.
- Instances cannot span into multiple placement groups.

Cluster Placement Group

- All the instances are placed in a same hardware rack inside an availability zone.
- But if the rack fails, all instances fail at the same time.
- This type of group will enable to achieve low latency network.
- Good for High Performance applications.

Use Case – Application with low latency and high throughput.

Spread Placement Group

- All EC2 instances are placed in different hardware racks in a single availability zone.
- A rack failure will not affect more than one instance.

We can create up to 7 EC2 instances per availability zone in spread placement group.

- Good for high availability applications and not suitable for high performance applications.

- It can span multiple availability zones in the same region.
- **Use Case-** Critical applications that needs to be isolated from failure from each other.

Partition Placement Group

- Group of instances spread across racks, and the instances in one partition do not share the underlying hardware with instances in different partition.
- If a rack fails, it will affect on the instances within that particular partition only.
- Partitions can be in different availability zones in the same region.

It strikes a balance between high performance and high availability.

Use Case – Big data applications like HDFS, HBase, Cassandra, Kafka.

In AWS EC2, there are three main types of load balancers:

Classic Load Balancer (CLB):

- CLB is the original load balancer service provided by AWS.
- It operates at both the application and transport layers (Layer 4 and Layer 7) of the OSI model.
- CLB supports distributing incoming traffic across multiple EC2 instances in multiple Availability Zones.
- It is suitable for applications that require simple load balancing without advanced features.

Application Load Balancer (ALB):

- ALB operates at the application layer (Layer 7) of the OSI model.
 - It supports advanced routing features such as host-based routing, path-based routing, and routing based on HTTP headers.
 - ALB can handle HTTP and HTTPS traffic and provides support for WebSocket and HTTP/2 protocols.
- It is ideal for modern web applications with multiple microservices or containers.

Network Load Balancer (NLB):

- NLB operates at the transport layer (Layer 4) of the OSI model.
-

It provides ultra-low latency and high throughput, making it suitable for handling millions of requests per second.

- NLB supports routing traffic to targets using IP addresses and TCP ports, making it ideal for TCP and UDP-based applications.
- It is commonly used for applications that require high performance, scalability, and reliability, such as gaming, IoT, and real-time communication applications.

Target Groups are a fundamental component of Elastic Load Balancing (ELB) services in AWS, specifically associated with the Application Load Balancer (ALB) and Network Load Balancer (NLB). Here's a concise explanation:

Routing:

- When configuring a load balancer (ALB or NLB), you specify one or more target groups to route traffic to.
- Each target group defines the criteria for routing requests to its associated targets, such as based on a specific path, host, or HTTP header.

Health Checks:

- Target Groups perform health checks on their registered targets to ensure they are capable of handling requests.
- Unhealthy targets are automatically removed from the target group until they pass health checks again.

Load Balancing:

- Target Groups distribute incoming traffic evenly across their registered targets according to the configured load balancing algorithm (e.g., round-robin or least outstanding requests).

Port Configuration:

- Target Groups allow you to specify the port on which the targets receive traffic, allowing flexibility in routing requests to different ports on the same target.

Autoscaling:

- Autoscaling is a feature provided by AWS that automatically adjusts the number of EC2 instances in a fleet based on user-defined policies.

- It helps maintain the desired number of instances to handle varying levels of application traffic or workload demand.
- Autoscaling can scale both up (adding more instances) and down (removing instances) dynamically, ensuring optimal resource utilization and cost efficiency.

Autoscaling Groups:

- An Autoscaling Group (ASG) is a logical grouping of EC2 instances managed by the Auto Scaling service.
- ASGs define the scaling policies, launch configurations, and other parameters that govern the behavior of the instances within the group.
- ASGs automatically launch or terminate EC2 instances based on predefined scaling policies and health checks.
- ASGs can be configured to distribute instances across multiple Availability Zones to improve fault tolerance and high availability.

Ways to connect to EC2 instances via

- the command line: SSH connection
- EC2 Instance Connect:
- AWS CLI (Command Line Interface):
- Third-Party Tools (like PuTTY or MobaXterm) and more.....

The launch of an ec2 instance:

- Choose an AMI (Amazon Machine Image): This is a template for the virtual machine you want to create. It contains the operating system, software, configuration, and sometimes data.
- Select an Instance Type: This determines the computational resources (CPU, RAM, storage, etc.) of your EC2 instance. There are various types optimized for different use cases.
- Configure Instance Details: This includes settings like networking (VPC, subnet), IAM role, and user data (scripts or data to be run at instance launch).

- Add Storage: Configure the amount and type of storage for your instance. This is typically done using Amazon EBS (Elastic Block Store), which provides block-level storage volumes.
- Add Tags: Assign metadata to your instance in the form of key-value pairs. This helps with organization, management, and billing.
- Configure Security Group: Define firewall rules for controlling inbound and outbound traffic to your instance. Security groups act as virtual firewalls.
- Review and Launch: Review your instance configuration, make any necessary changes, select an existing key pair or create a new one for SSH access (for Linux instances), and then launch the instance.

After launching the instance, you can access it using SSH (for Linux instances) or RDP (for Windows instances) using the key pair you selected or created. Remember to properly secure your instance and follow best practices for managing AWS resources.

Let's walk through the process of launching an EC2 instance.

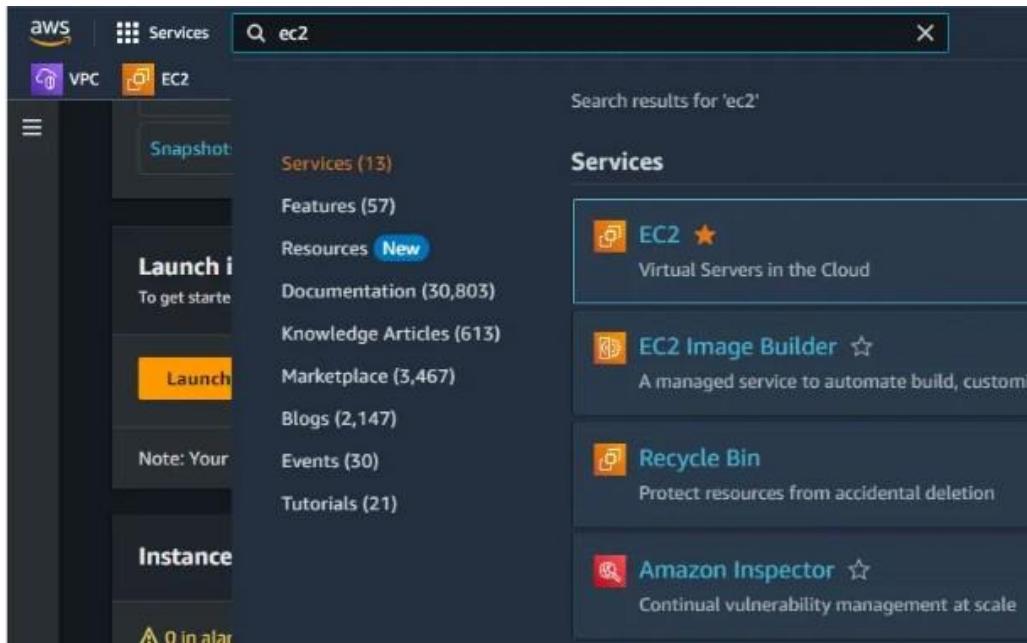
An EC2 instance essentially functions as a virtual machine.

Firstly, log in to your AWS console.

Next, choose your desired region. For instance, I'll be selecting the North Virginia region. It's advisable to select a region based on your current location, as some regions may incur higher costs compared to those in the US.

Each region is identified by a code, such as US East 1 for North Virginia.

To navigate to the EC2 service, click on “Services” and locate EC2 either by browsing through all services or by directly searching for EC2.



Once in the EC2 dashboard, you'll find an overview of your instances, volumes, and other relevant information. Keep in mind that the dashboard doesn't auto-refresh, so be sure to manually refresh after making any changes.

Scrolling down, you'll notice zones within the region, each identified by a code like 1A, 1B, etc. Typically, a region will have at least two zones.

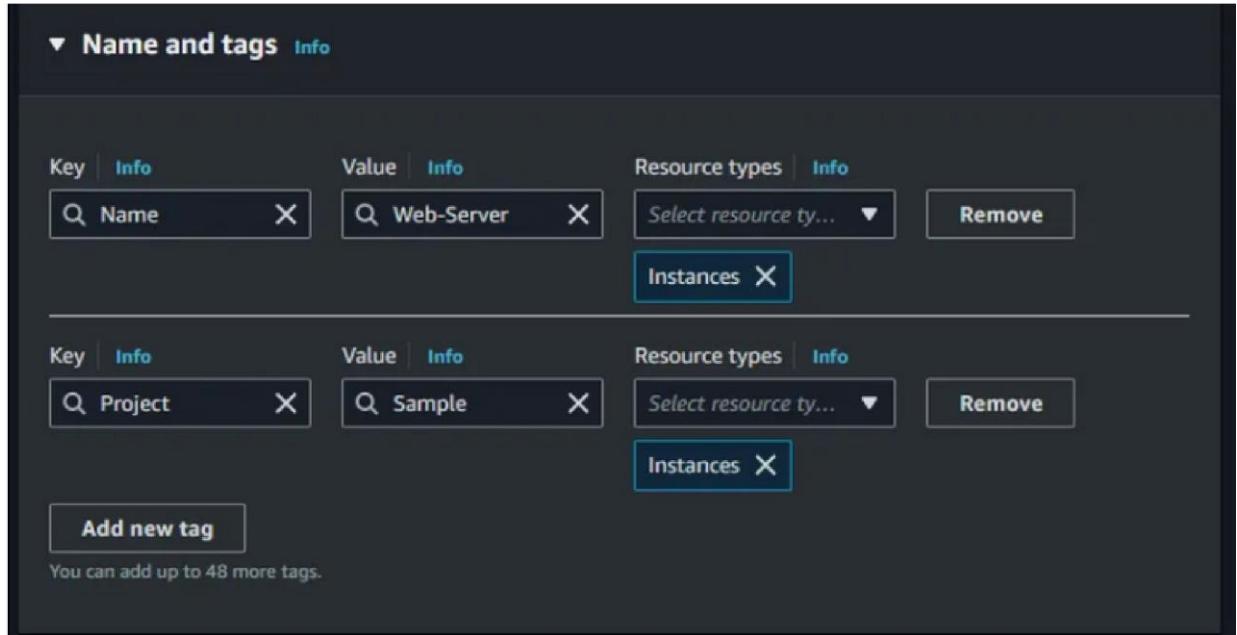
The screenshot shows the AWS Health Dashboard. At the top, it displays "Service health" and "AWS Health Dashboard" with a refresh icon. Below this, the "Region" is listed as "US East (N. Virginia)" and the "Status" is "This service is operating normally." A green checkmark icon is present next to the status message. The main section is titled "Zones" and contains a table with two columns: "Zone name" and "Zone ID". The zones listed are: us-east-1a (use1-az1), us-east-1b (use1-az2), us-east-1c (use1-az4), us-east-1d (use1-az6), us-east-1e (use1-az3), and us-east-1f (use1-az5). At the bottom of the table, there is a link labeled "Enable additional Zones".

Click on “Instances” and then “Launch Instance” to initiate the instance creation process.

Before proceeding, it’s worth noting that the AWS console interface may evolve over time, but the core options and functionalities remain consistent.

The screenshot shows the "Launch an instance" wizard. The first step, "Name and tags", is currently selected. It includes an informational message: "Amazon EC2 allows you to create virtual machines, or instances, that run on the AWS Cloud. Quickly get started by following the simple steps below." Below this, there is a "Name" field containing "Web-Server" and a "Add additional tags" button.

Now, let’s begin by assigning a name and optional tags to the instance. Tags can be useful for filtering and organization purposes.



Moving on to selecting the Amazon Machine Image (AMI), you'll find a variety of options including quick start AMIs like Amazon Linux, Ubuntu, and Windows, as well as community AMIs and those from the AWS Marketplace. Carefully review the details and charges associated with each AMI.

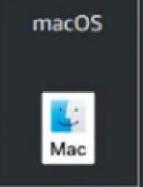
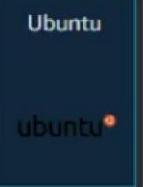
For this example, let's choose the Ubuntu AMI provided by Amazon Web Services, which is free.

▼ Application and OS Images (Amazon Machine Image) [Info](#)

An AMI is a template that contains the software configuration (operating system, application server, and applications) required to launch your instance. Search or Browse for AMIs if you don't see what you are looking for below.

Search our full catalog including 1000s of application and OS images

Quick Start

 [Browse more AMIs](#)
Including AMIs from AWS, Marketplace and the Community

Amazon Machine Image (AMI)

Ubuntu Server 22.04 LTS (HVM), SSD Volume Type Free tier eligible ▾

ami-080e1f13689e07408 (64-bit (x86)) / ami-0a55ba1c20b74fc30 (64-bit (Arm))
Virtualization: hvm ENA enabled: true Root device type: ebs

Description
Canonical, Ubuntu, 22.04 LTS, amd64 jammy image build on 2024-03-01

Architecture 64-bit (x86) ▾ **AMI ID** ami-080e1f13689e07408 Verified provider

Ensure that the instance type selected aligns with your requirements. For free-tier usage, T2 micro instances are suitable for most learning purposes.

▼ Instance type [Info](#) | [Get advice](#)

Instance type

t2.micro Free tier eligible

Family: t2 1 vCPU 1 GiB Memory Current generation: true
On-Demand Windows base pricing: 0.0162 USD per Hour
On-Demand SUSE base pricing: 0.0116 USD per Hour
On-Demand RHEL base pricing: 0.0716 USD per Hour
On-Demand Linux base pricing: 0.0116 USD per Hour

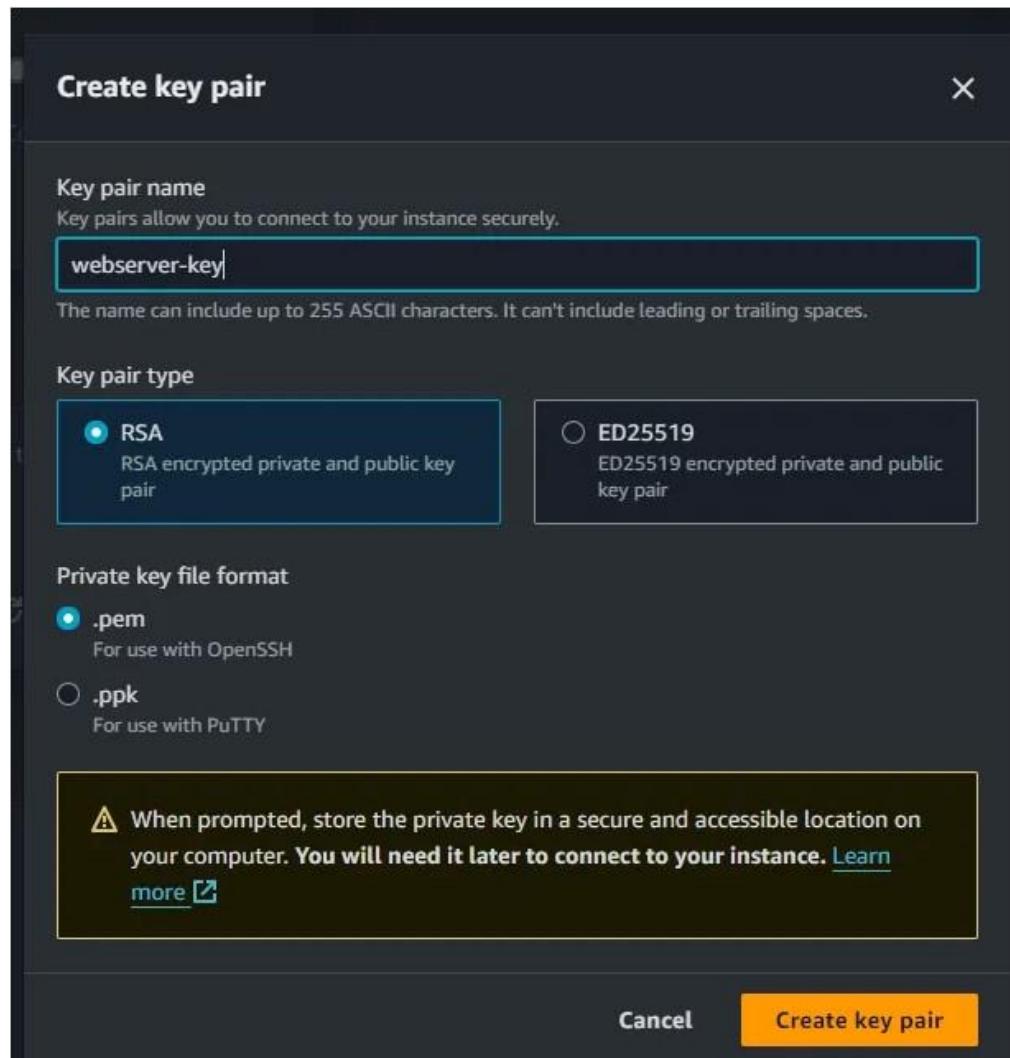
All generations [Compare instance types](#)

Additional costs apply for AMIs with pre-installed software

- Key pairs are used primarily for securely accessing EC2 instances. When

create an EC2 instance, you can specify a key pair to use for SSH access to the instance. The key pair consists of a public key that AWS stores, and a private key file that you download to your local machine.

- It's important to note that AWS doesn't store your private key, so if you lose it, there's no way to retrieve it from AWS. Make sure to securely store your private keys and avoid sharing them with unauthorized users.



Proceed to configure the network settings, including security groups. Security groups act as mandatory firewalls for EC2 instances.

▼ Network settings [Info](#)

[Edit](#)

Network | [Info](#)
vpc-03c70ece7f453f834

Subnet | [Info](#)
No preference (Default subnet in any availability zone)

Auto-assign public IP | [Info](#)
Enable
Additional charges apply when outside of [free tier allowance](#)

Firewall (security groups) | [Info](#)
A security group is a set of firewall rules that control the traffic for your instance. Add rules to allow specific traffic to reach your instance.

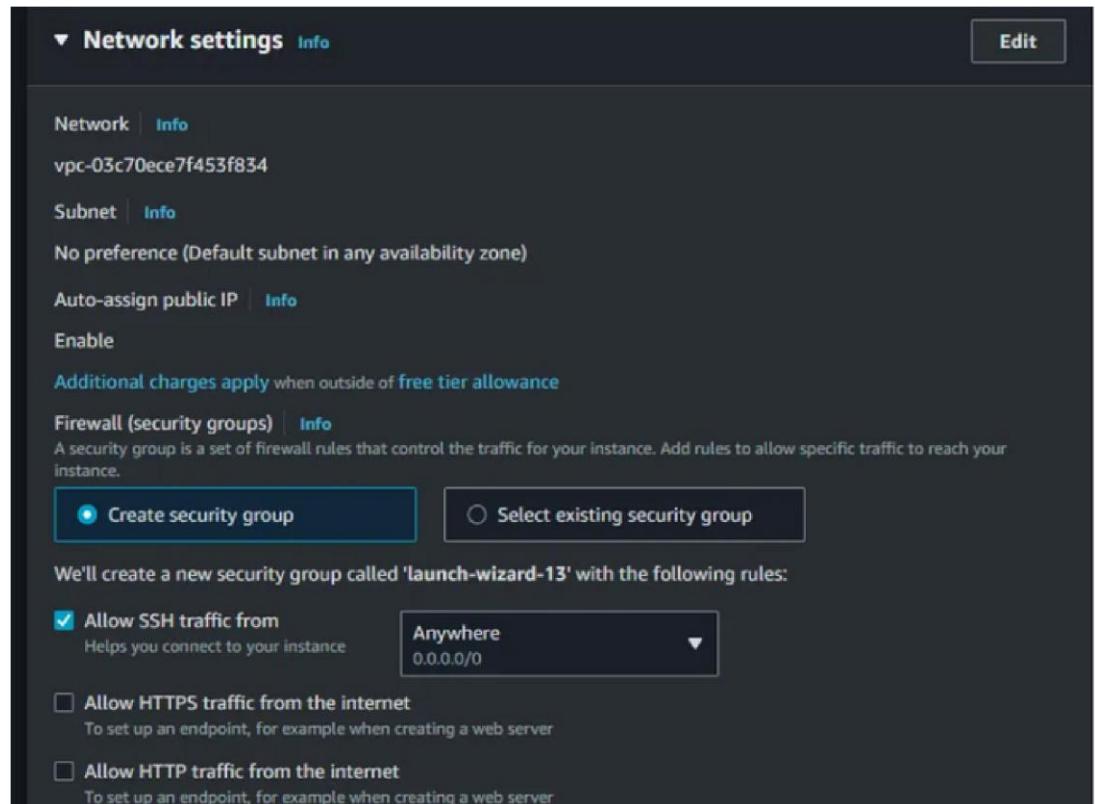
Create security group Select existing security group

We'll create a new security group called 'launch-wizard-13' with the following rules:

Allow SSH traffic from Anywhere
Helps you connect to your instance
0.0.0.0/0

Allow HTTPS traffic from the internet
To set up an endpoint, for example when creating a web server

Allow HTTP traffic from the internet
To set up an endpoint, for example when creating a web server



Customize the security group to allow necessary inbound and outbound traffic.

A security group is a set of firewall rules that control the traffic for your instance. Add rules to allow specific traffic to reach your instance.

Create security group Select existing security group

Security group name - *required*
Web-Server-SG

This security group will be added to all network interfaces. The name can't be edited after the security group is created. Max length is 255 characters. Valid characters: a-z, A-Z, 0-9, spaces, and _-:/()#,@[]+=&;!\$*

Description - *required* | [Info](#)
Web-Server-SG

Inbound Security Group Rules

▼ Security group rule 1 (TCP, 22, 0.0.0.0/0) [Remove](#)

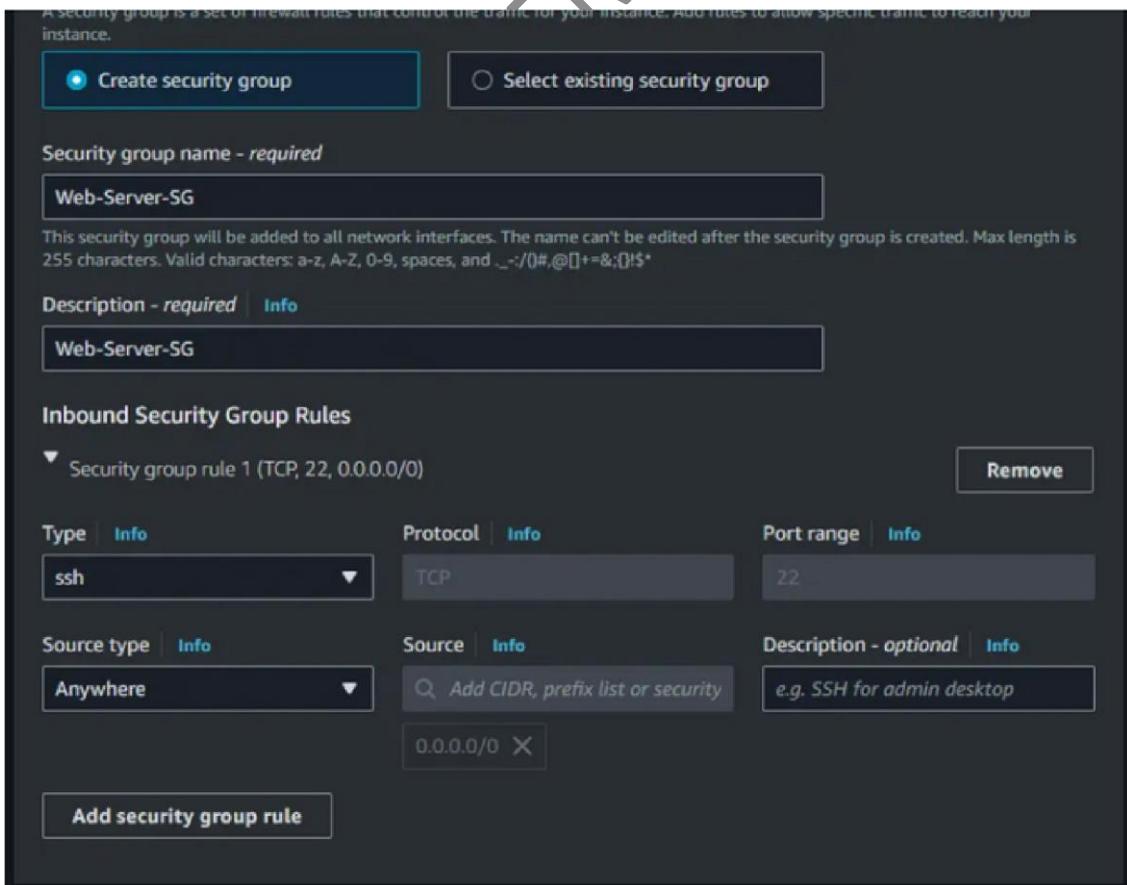
| Type Info | Protocol Info | Port range Info |
|-----------------------------|---------------------------------|-----------------------------------|
| ssh | TCP | 22 |

Source type | [Info](#)
Anywhere

Source | [Info](#)
 Add CIDR, prefix list or security
0.0.0.0/0 [X](#)

Description - *optional* | [Info](#)
e.g. SSH for admin desktop

[Add security group rule](#)



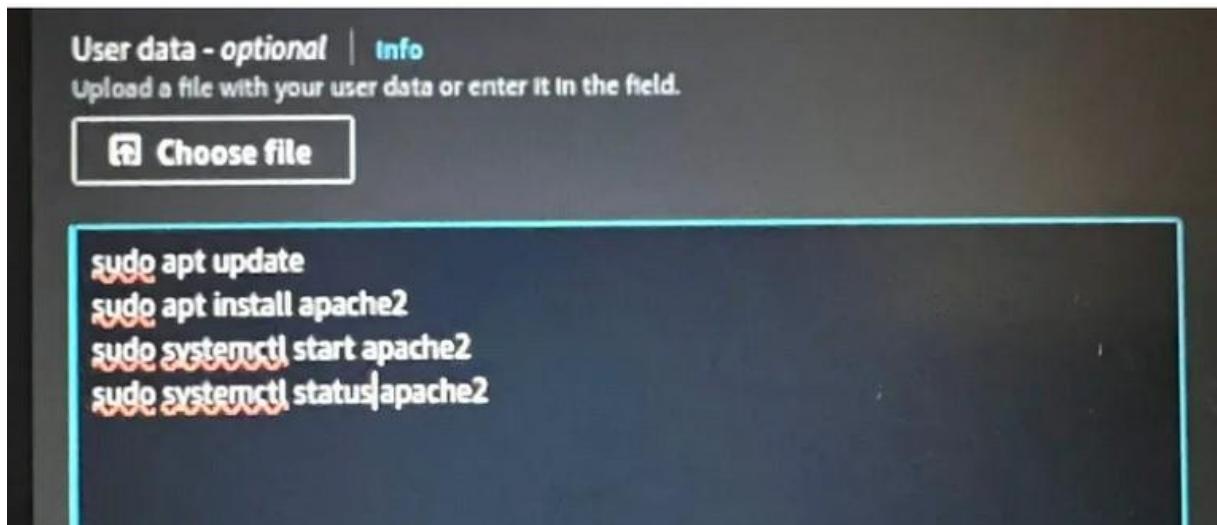
You can specify the storage size, typically defaulted to 10 GB, and further advanced details like user data, which allows for automated provisioning upon instance launch. Finally, review your configuration and launch the instance.

Upon launching, the instance may take a few moments to become available. You can monitor its status in the EC2 dashboard.

Next, establish a connection to the instance; in this example, I'm using EC2 Instance Connect.

There are two methods to install Apache2:

- By incorporating commands in the user data section before initiating the instance launch, as demonstrated below:



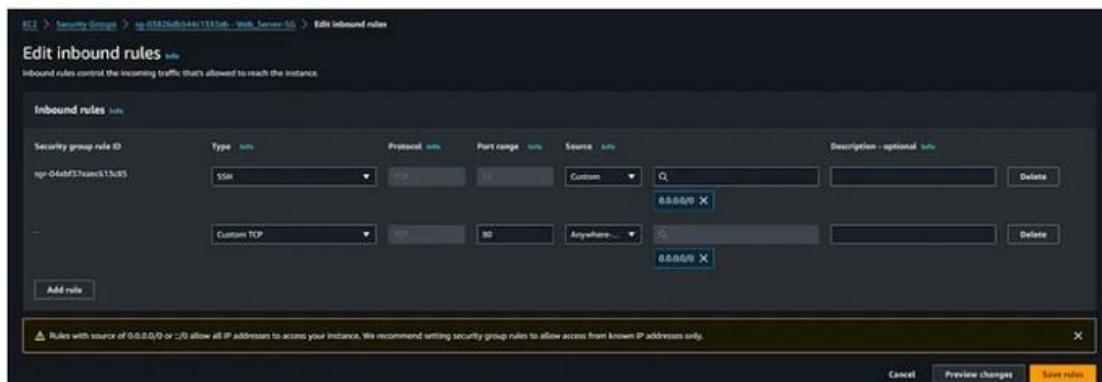
- Or, by executing the commands after connecting the ec2 instance via terminal.

- *sudo apt*
 - *update sudo*
 - *apt install*
 - *apache2 -y*
 - *sudo*
- systemctl*
- start apache2*
- sudo*
- systemctl*
- enable*
- apache2 sudo*

```
systemctl  
status  
apache2
```

Troubleshooting tip: If encountering connection issues, ensure the correct path to the SSH key or navigate to the directory containing the key before attempting to connect.

Once connected, you can verify the instance status, access web services like Apache, and configure firewall rules as needed to enable external access to your services. Additionally, click on the Security Group, edit inbound rules, add the following rule, and save the rules. It's important to note that allowing all IPs is solely for learning purposes and may vary depending on your project requirements.



Go to the browser and copy the public IP followed by ':80' to view the Apache page.

This is the default welcome page used to test the correct operation of the Apache2 server after installation on Ubuntu systems. It is based on the equivalent page on Debian, from which the Ubuntu Apache packaging is derived. If you can read this page, it means that the Apache HTTP server installed at this site is working properly. You should **replace this file** (located at `/var/www/html/index.html`) before continuing to operate your HTTP server.

If you are a normal user of this web site and don't know what this page is about, this probably means that the site is currently unavailable due to maintenance. If the problem persists, please contact the site's administrator.

Configuration Overview

Ubuntu's Apache2 default configuration is different from the upstream default configuration, and split into several files optimized for interaction with Ubuntu tools. The configuration system is **fully documented in `/usr/share/doc/apache2/README.Debian.gz`**. Refer to this for the full documentation. Documentation for the web server itself can be found by accessing the **manual** if the `apache2-doc` package was installed on this server.

The configuration layout for an Apache2 web server installation on Ubuntu systems is as follows:

```
/etc/apache2/
|-- apache2.conf
|   '-- ports.conf
|-- mods-enabled
|   '-- *.load
|   '-- *.conf
|-- conf-enabled
|   '-- *.conf
|-- sites-enabled
|   '-- *.conf
```

Elastic IP

- Powering Off the EC2 Instance:
- Take note of the current public and private IP addresses of the instance.
- Stop the instance using the AWS Management Console.
- Observing IP Address Changes:

After stopping the instance, observe that the public IP address is released while the private IP address remains unchanged.

- Refresh the console to verify the changes.

| Name | Instance ID | Instance state | Instance type | Status check | Alarm status | Availability Zone | Public IPv4 DNS | Public IPv4 | Elastic IP |
|------------|---------------------|----------------|---------------|--------------|--------------|-------------------|-----------------|-------------|------------|
| Web_Server | i-059ac545443e64c80 | Stopped | t2.micro | - | - | us-east-1c | - | - | - |

Instance: i-059ac545443e64c80 (Web_Server)

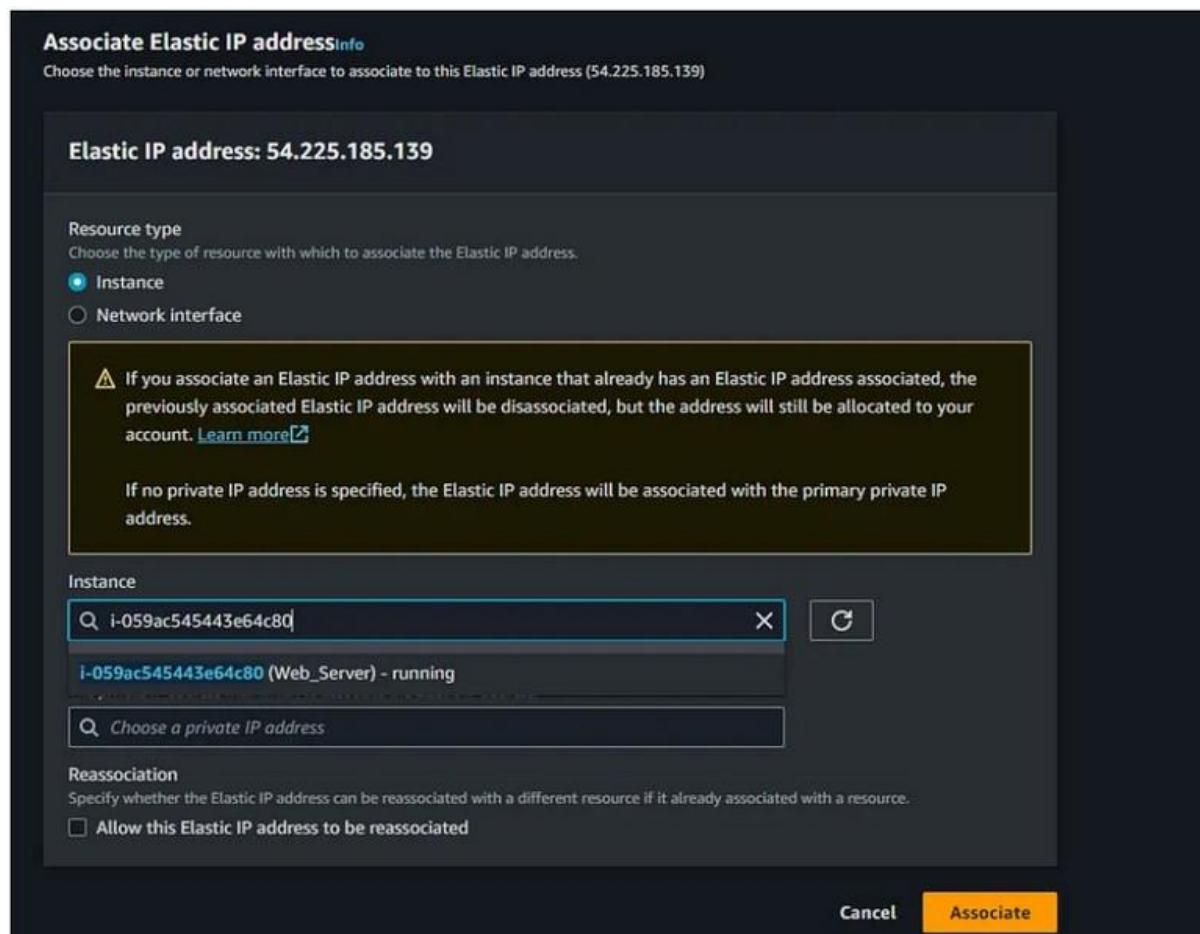
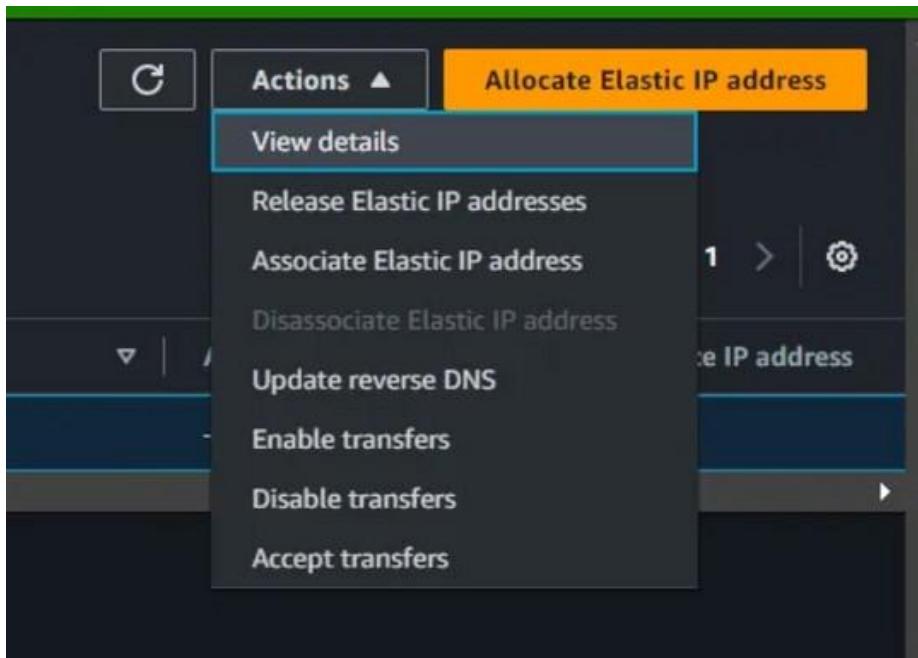
Details Status and alarms New Monitoring Security Networking Storage Tags

Instance summary Info

| Instance ID | Public IPv4 address | Private IPv4 address |
|----------------------------------|---------------------|----------------------|
| i-059ac545443e64c80 (Web_Server) | - | 172.31.26.160 |
| IPv6 address | Instance state | Public IPv4 DNS |
| - | Stopped | - |

Elastic IP Benefits:

- Elastic IP addresses offer the advantage of maintaining a static public IP address even in scenarios where EC2 instances are restarted or recreated.
- They provide a stable endpoint for applications or services that require a consistent public IP address.
- With Elastic IP, there's no need to update DNS records or reconfigure applications each time an instance is stopped or replaced.
- It simplifies infrastructure management and enhances reliability by ensuring uninterrupted connectivity to resources hosted on EC2 instances.
- Elastic IP addresses are particularly useful for scenarios requiring a fixed public IP address, such as hosting websites, APIs, or VPN gateways.
- Allocating an Elastic IP:
 - Allocate an Elastic IP address from the AWS pool.
 - Associate the Elastic IP with the desired EC2 instance.



EC2 Instance Components:

- In addition to the EC2 instance itself, other components include the network interface and volumes (virtual hard disks).

- The network interface manages connectivity and is associated with firewall rules, public IP addresses, and security groups.
- Volumes represent the virtual hard disks attached to the instance, storing data and operating system files.
- Understanding these components is essential for managing and configuring EC2 instances effectively.

| Name | Network Interface ID | Subnet ID | VPC ID | Availability Zone | Security group n... | Security group IDs | Interface Type |
|----------------------|--------------------------|----------------------|------------|-------------------|----------------------|----------------------|---------------------------|
| eni-0771a05a64031c49 | subnet-08d3ae9c060e0c459 | vpc-05c70ec7f453fb54 | us-east-1c | Web_Server-SG | sg-03826db544c139... | sg-03826db544c139... | Elastic network interface |

| Name | Volume ID | Type | Size | IOPS | Throughput | Snapshot | Created | Availability Zone | Volume state | Alarm stat |
|-----------------------|-----------|-------|------|------|------------|----------------|---------------------------|-------------------|--------------|------------|
| vol-05060cb3eeabd45bf | gp2 | 8 GiB | 100 | - | - | snap-057fd6... | 2024/04/12 10:55 GMT+5... | us-east-1c | In-use | No alarms |

Navigating the EC2 Dashboard:

- Encourage exploration of the EC2 dashboard to view information about resources like key pairs, security groups, and volumes.

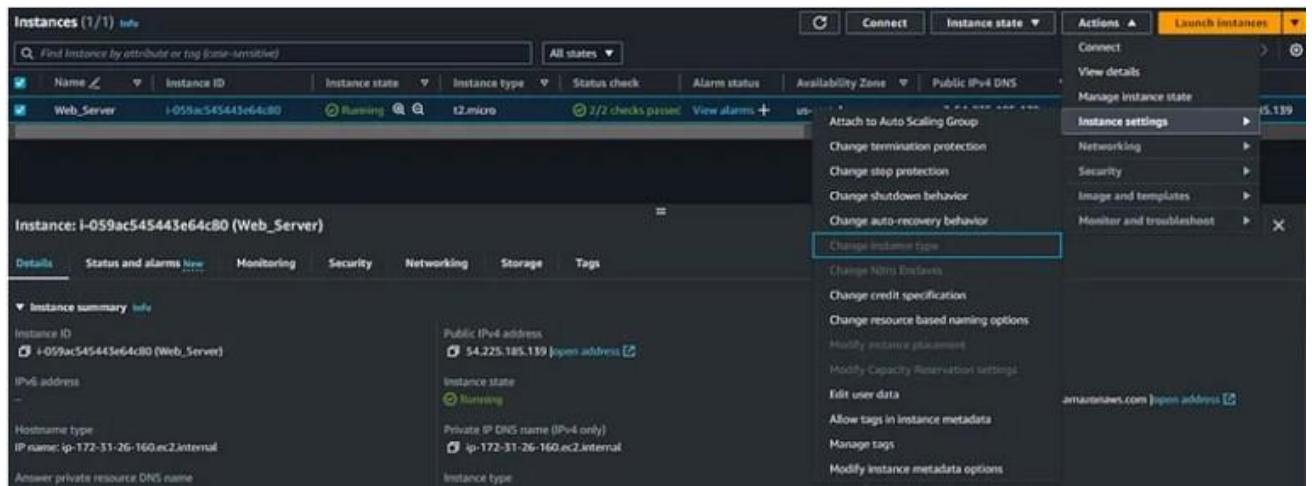
| Resources | | EC2 Global view | | | |
|---|---|---------------------|---|--|--|
| You are using the following Amazon EC2 resources in the US East (N. Virginia) Region: | | | | | |
| | | | | | |
| Instances (running) | 1 | Auto Scaling Groups | 0 | | |
| Elastic IPs | 1 | Instances | 1 | | |
| Load balancers | 0 | Placement groups | 0 | | |
| Snapshots | 0 | Volumes | 1 | | |
| | | | | | |

Adjusting Instance Settings:

- Demonstrate how to change instance type and other settings like Elastic IP disassociation, network interface attachment/detachment, and security group changes.

Reverting Instance Changes:

- Revert any instance changes made during the demonstration, such as changing the instance type, by powering on the instance and reverting the settings back to their original state.



Elastic Block Store (EBS)

EBS serves as the virtual hard disk for EC2 instances, offering both EBS volumes and snapshots for data storage and backup.

EBS Volumes:

- EBS volumes, akin to physical hard disks, house the operating system data of EC2 instances and can store various types of data such as databases, web servers, and files.
- During volume creation, users specify an availability zone, ensuring alignment with the instance's location.
- Data stored in EBS volumes is automatically replicated within the same availability zone, mitigating the risk of data loss due to local failures.

Snapshots:

- Snapshots, used for backing up entire volumes, provide an additional layer of data protection.
- These snapshots capture the state of the EBS volume at a specific point in time, facilitating data recovery and restoration.

Types of EBS Volumes:

- EBS volumes come in various types catering to different workload requirements and performance needs.
- General purpose volumes, SSD-based, are recommended for most workloads due to their balanced pricing and speed.

- Provisioned IOPS volumes offer enhanced performance and are ideal for large databases requiring high input-output operations per second (IOPS).
- Throughput optimized HDD volumes suit big data and data warehousing applications, while cold HDD volumes are cost-effective solutions for infrequently accessed data.
- Magnetic volumes, the cheapest option, are suitable for backups and archives but offer lower performance.

Choosing the Right Volume Type:

- The selection of an appropriate volume type depends on factors such as workload characteristics, performance requirements, and budget considerations.
- AWS documentation provides detailed information and use cases for each volume type, aiding users in making informed decisions.

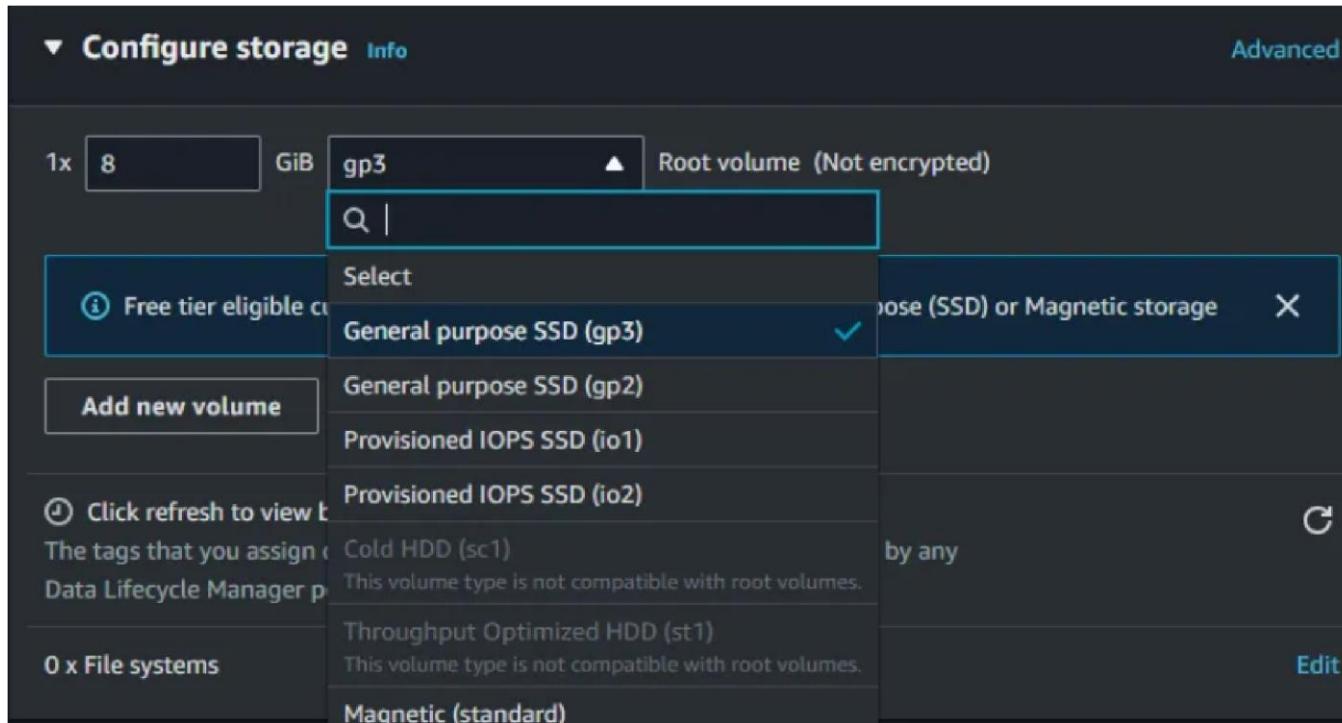
After establishing a connection to the instance, execute the following commands. Please note that these commands can also be added to the instance creation process in the user data area:

```
•  
•  
•  
•  
•  
systemctl  
status  
apache2 cd  
/tmp  
wget https://www.tooplate.com/zip-  
templates/2119_gymso_fitness.zip unzip -o  
2119_gymso_fitness.zip cp -r 2119_gymso_fitness/*  
/var/www/html/ systemctl restart apache2
```

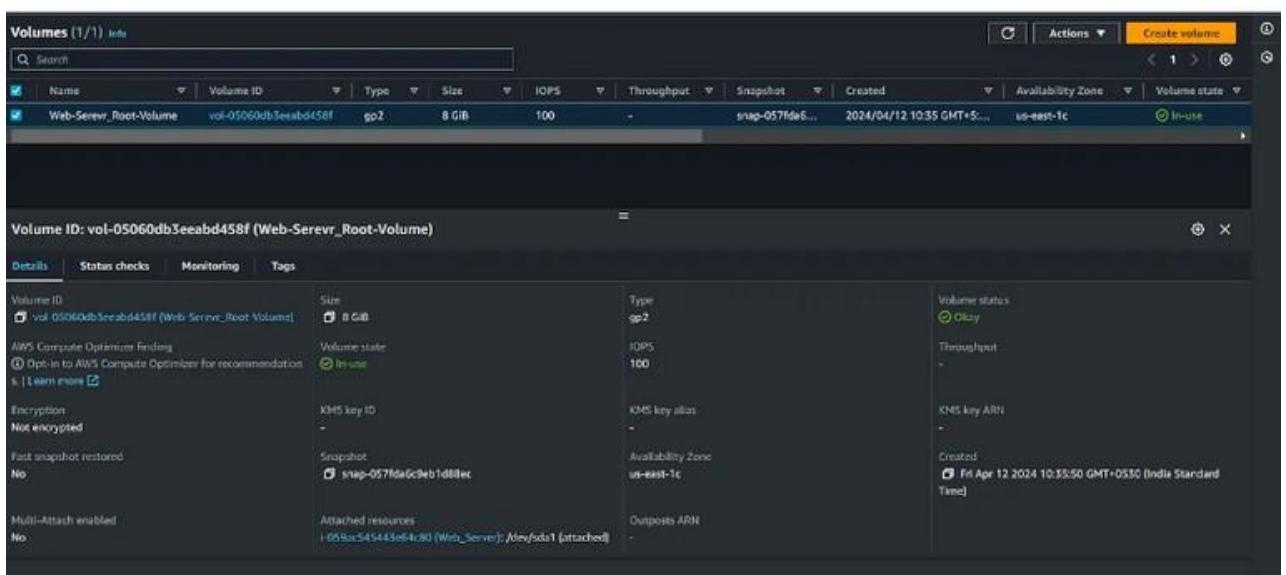
These commands are designed to check the status of the Apache web server, navigate to the temporary directory, download a ZIP file containing website templates from a specified URL, unzip the downloaded file, copy the contents to the Apache document root

directory, and finally restart the Apache service to apply the changes.

Here, we can observe the available volume types for creating an EC2 instance.



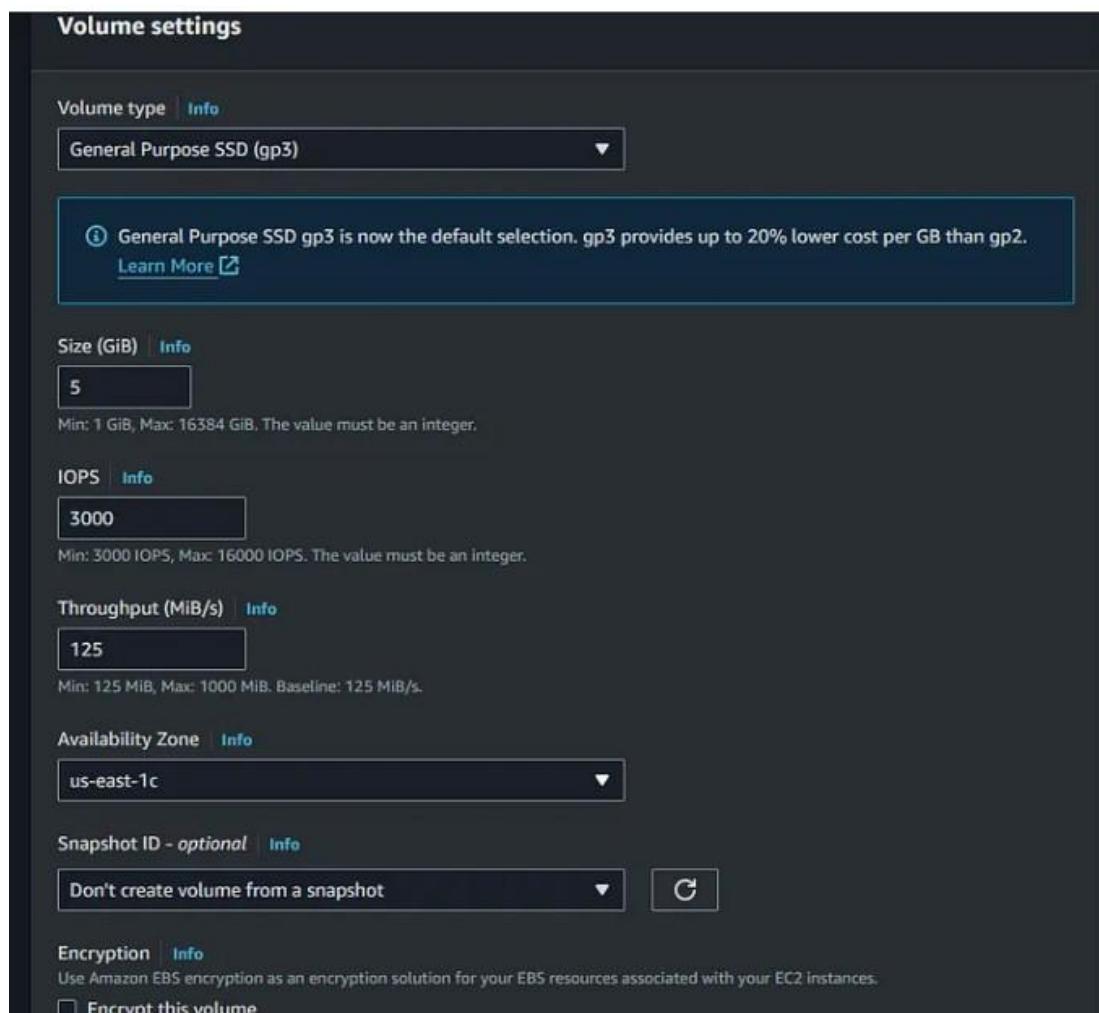
Navigate to the Volumes section located on the left-hand side of the EC2 instances page and update the name of the volume associated with the instance mentioned earlier.



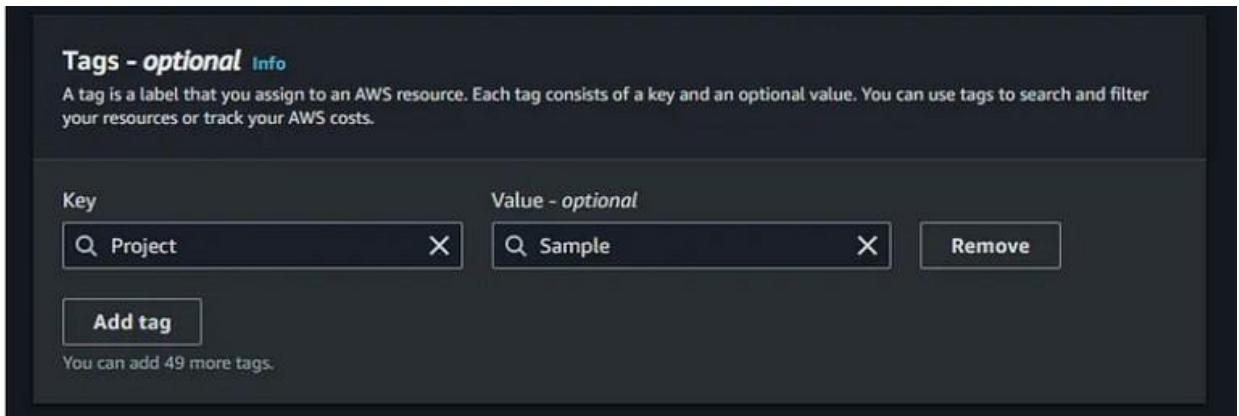
Adding extra storage to an EC2 instance, it's crucial to consider the availability zone. Both the instance and the volume need to be in the same zone to ensure connectivity.

When gathering requirements for additional storage, such as for web server images, it's important to adhere to the free tier limit of 30 GB for EBS volumes. For instance, if there's a need for five gigabytes of storage for web server images, it's imperative to ensure that the total storage remains within the free tier limit.

Click on create volume:

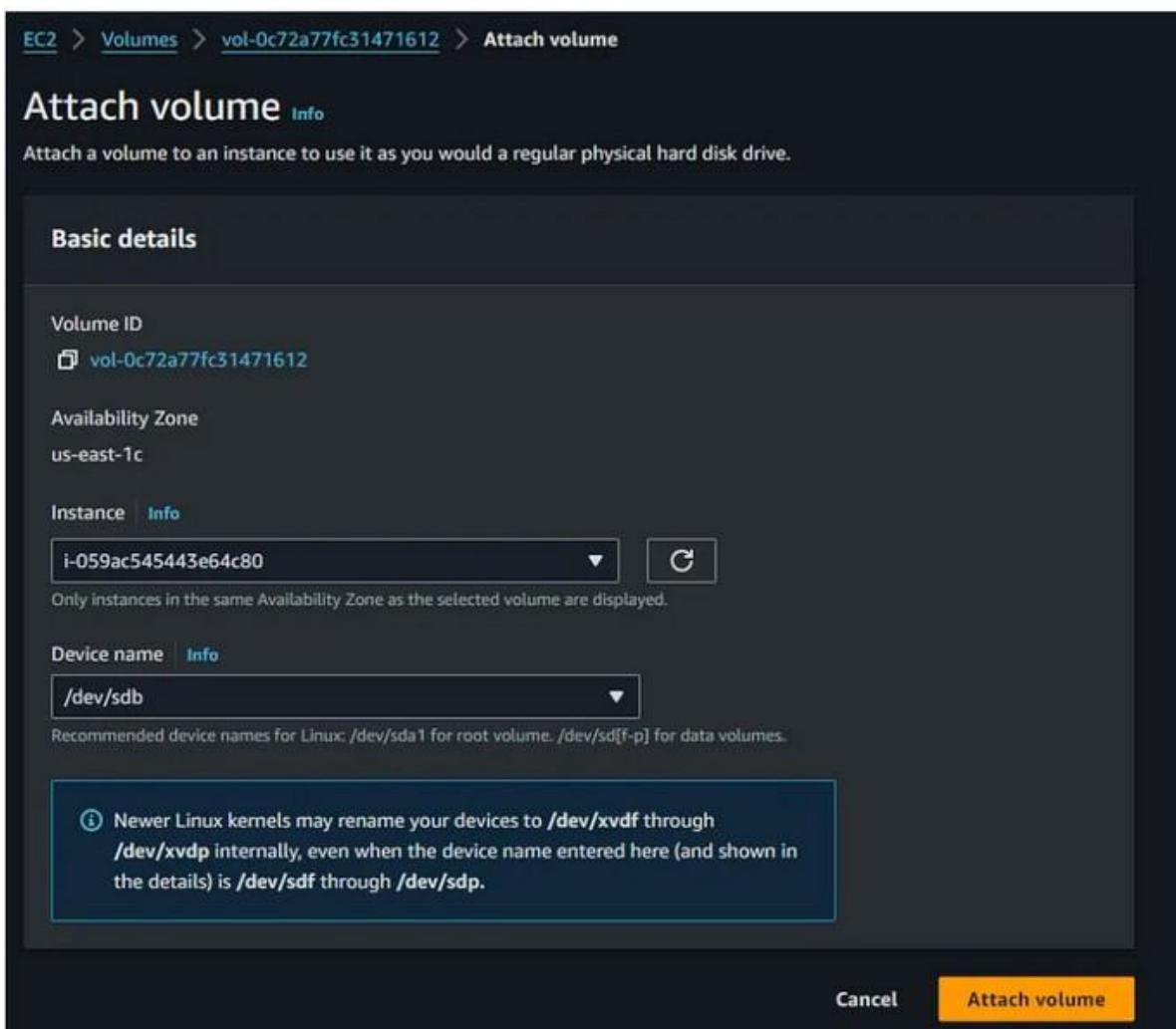


Tags, you can give upon your wish, then click on Create Volume



Choose the created volume, navigate to actions, and then opt for attaching the volume once it appears as available on the dashboard.

Select the instance , device name and click on attach volume.



(Please note I've selected /dev/xvdf in above screenshot as Device name). If we run the “fdisk -l” command, we could see the new volume listed:

•***fdisk -l***

List Disks: Start by running the command fdisk -l to list all disks. For instance, /dev/xvdf represents the root volume with a size of 8 GB, and it has a partition /dev/xvdaf1.

```
Disk /dev/xvdf: 5 GiB, 5368709120 bytes, 10485760 sectors
Units: sectors of 1 * 512 = 512 bytes
Sector size (logical/physical): 512 bytes / 512 bytes
I/O size (minimum/optimal): 512 bytes / 512 bytes
root@ip-172-31-26-160:/tmp#
```

- Understand Partitions: Linux assigns partitioning numbers (e.g., xvda1) rather than names. Verify the root volume by checking its mount point using df -h. If /dev/xvda1 is mounted to the root directory, it contains the operating system data.
- Create Partition: Identify the additional disk, such as /dev/xvdf, and use fdisk to create a partition. Choose the disk path (e.g., /dev/xvdf), follow the prompts to create a new partition, and write the changes using w.

```
Command (m for help): m

Help:

DOS (MBR)
  a  toggle a bootable flag
  b  edit nested BSD disklabel
  c  toggle the dos compatibility flag

Generic
  d  delete a partition
  F  list free unpartitioned space
  l  list known partition types
  n  add a new partition
  p  print the partition table
  t  change a partition type
  v  verify the partition table
  i  print information about a partition

Misc
  m  print this menu
  u  change display/entry units
  x  extra functionality (experts only)

Script
  I  load disk layout from sfdisk script file
  O  dump disk layout to sfdisk script file

Save & Exit
  w  write table to disk and exit
  q  quit without saving changes

Create a new label
  g  create a new empty GPT partition table
  G  create a new empty SGI (IRIX) partition table
  o  create a new empty DOS partition table
  s  create a new empty Sun partition table

Command (m for help): n
```

```
Command (m for help): n
Partition type
  p  primary (0 primary, 0 extended, 4 free)
  e  extended (container for logical partitions)
Select (default p): p
Partition number (1-4, default 1): 1
First sector (2048-10485759, default 2048):
Last sector, +/-sectors or +/-size{K,M,G,T,P} (2048-10485759, default 10485759):

Created a new partition 1 of type 'Linux' and of size 5 GiB.

Command (m for help): p
Disk /dev/xvdf: 5 GiB, 5368709120 bytes, 10485760 sectors
Units: sectors of 1 * 512 = 512 bytes
Sector size (logical/physical): 512 bytes / 512 bytes
I/O size (minimum/optimal): 512 bytes / 512 bytes
Disklabel type: dos
Disk identifier: 0xe3548d72

Device      Boot Start      End Sectors Size Id Type
/dev/xvdf1        2048 10485759 10483712   5G 83 Linux

Command (m for help): w
The partition table has been altered.
Calling ioctl() to re-read partition table.
Syncing disks.

root@ip-172-31-26-160:~#
```

- fdisk -l

```
Device      Boot Start      End Sectors Size Id Type
/dev/xvdf1        2048 10485759 10483712   5G 83 Linux
root@ip-172-31-26-160:~#
```

Format Partition: Use the appropriate formatting utility, such as mkfs.ext4, to format the partition. For example, run mkfs.ext4 /dev/xvdf1 to format the partition with ext4 format.

```
mkfs.cramfs  mkfs.ext2    mkfs.ext3    mkfs.ext4    mkfs.fat    mkfs.minix  mkfs.msdos
```

- mkfs.ext4 /dev/xvdf1

Multicloud Devops by veera nareshit

```
root@ip-172-31-26-160:~# mkfs.ext4 /dev/xvdf1
mke2fs 1.46.5 (30-Dec-2021)
Creating filesystem with 1310464 4k blocks and 327680 inodes
Filesystem UUID: 18c5893b-0f5d-49f8-bb29-167b8dbff181
Superblock backups stored on blocks:
    32768, 98304, 163840, 229376, 294912, 819200, 884736

Allocating group tables: done
Writing inode tables: done
Creating journal (16384 blocks): done
Writing superblocks and filesystem accounting information: done

root@ip-172-31-26-160:~#
```

Temporary Mount: Mount the partition temporarily using the mount command. For instance, mount /dev/xvdf1 /var/www/html/images. Verify the mount using *df -h*.

Create a directory: “ mkdir /tmp/image_backups”, and move all data from Images to newly created one

```
root@ip-172-31-26-160:/var/www/html#
root@ip-172-31-26-160:/var/www/html# mv images/* /tmp/image_backups/
root@ip-172-31-26-160:/var/www/html#
```

Then mount the volume and verify as below:

```
root@ip-172-31-26-160:~#
root@ip-172-31-26-160:~# mount /dev/xvdf1 /var/www/html/images/
root@ip-172-31-26-160:~#
root@ip-172-31-26-160:~# df -h
Filesystem      Size  Used Avail Use% Mounted on
/dev/root       7.6G  2.4G  5.3G  31% /
tmpfs          475M     0  475M   0% /dev/shm
tmpfs          190M  868K  190M   1% /run
tmpfs          5.0M     0  5.0M   0% /run/lock
/dev/xvda15     105M   6.1M   99M   6% /boot/efi
tmpfs          95M   4.0K   95M   1% /run/user/1000
/dev/xvdf1      4.9G   24K  4.6G   1% /var/www/html/images
root@ip-172-31-26-160:~#
root@ip-172-31-26-160:~#
```

Permanent Mount: Edit the /etc/fstab file to configure permanent mounting. Add an entry with the partition path, mount point, format type, and other options. Save the file and test the configuration using mount -a.

```
LABEL=cloudimg-rootfs   /           ext4   discard,errors=remount-ro   0 1
LABEL=UEFI      /boot/efi        vfat    umask=0077      0 1
/dev/xvdf1     /var/www/html/images  ext4    defaults        0 0
```

Save the file and exit.

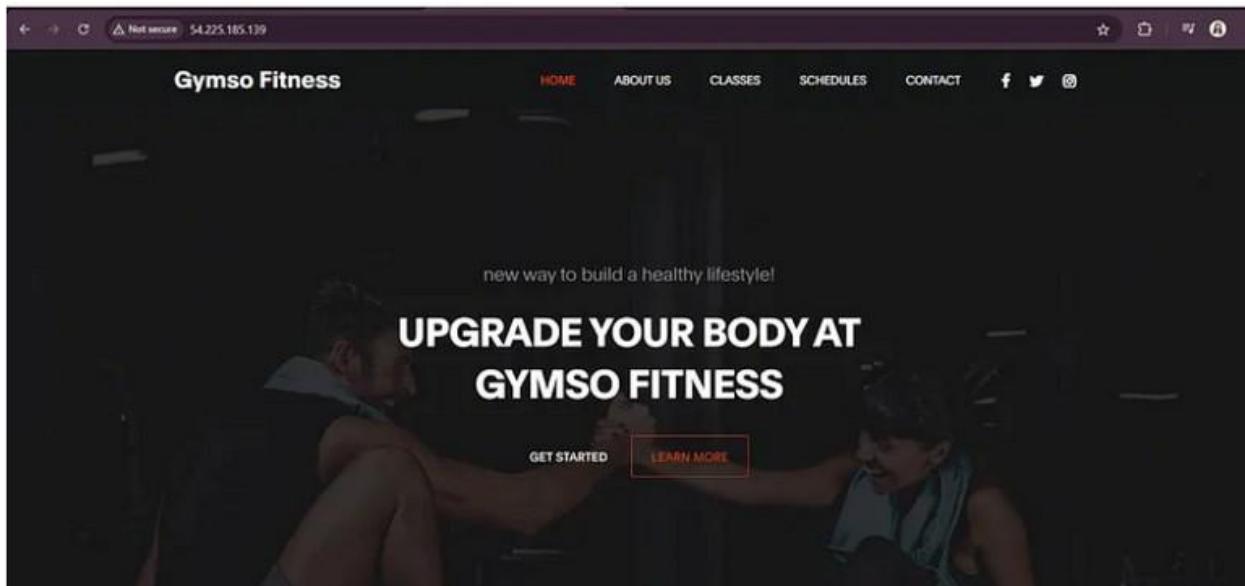
```
root@ip-172-31-26-160:~#  
root@ip-172-31-26-160:~# mount -a  
root@ip-172-31-26-160:~#
```

```
root@ip-172-31-26-160:~#  
root@ip-172-31-26-160:~# df -h  
Filesystem      Size  Used Avail Use% Mounted on  
/dev/root       7.6G  2.4G  5.3G  31% /  
tmpfs           475M    0  475M   0% /dev/shm  
tmpfs           190M  868K  190M   1% /run  
tmpfs           5.0M    0  5.0M   0% /run/lock  
/dev/xvda15     105M  6.1M  99M   6% /boot/efi  
tmpfs           95M  4.0K  95M   1% /run/user/1000  
/dev/xvdf1      4.9G  24K  4.6G   1% /var/www/html/images  
root@ip-172-31-26-160:~#
```

Move Data: Move data from the backup directory to the mounted partition. For example, use mv command to move data from /tmp/img-backups to /var/www/html/images., and restart the apache2 service

```
root@ip-172-31-26-160:~#  
root@ip-172-31-26-160:~# mv /tmp/image_backups/* /var/www/html/images/  
root@ip-172-31-26-160:~#  
root@ip-172-31-26-160:~#  
root@ip-172-31-26-160:~# systemctl restart apache2  
root@ip-172-31-26-160:~# systemctl status apache2  
● apache2.service - The Apache HTTP Server  
   Loaded: loaded (/lib/systemd/system/apache2.service; enabled; vendor preset: enabled)  
   Active: active (running) since Mon 2024-04-15 00:48:20 UTC; 13s ago  
     Docs: https://httpd.apache.org/docs/2.4/  
   Process: 3770 ExecStart=/usr/sbin/apachectl start (code=exited, status=0/SUCCESS)  
 Main PID: 3776 (apache2)  
    Tasks: 55 (limit: 1121)  
   Memory: 4.8M  
      CPU: 28ms  
     CGroup: /system.slice/apache2.service  
             └─3776 /usr/sbin/apache2 -k start  
                 ├─3777 /usr/sbin/apache2 -k start  
                 ├─3778 /usr/sbin/apache2 -k start  
  
Apr 15 00:48:20 ip-172-31-26-160 systemd[1]: Starting The Apache HTTP Server...  
Apr 15 00:48:20 ip-172-31-26-160 systemd[1]: Started The Apache HTTP Server.  
root@ip-172-31-26-160:~#  
root@ip-172-31-26-160:~#
```

Verify: Check if the data is accessible through the browser. If not, consider disabling SELinux temporarily by modifying the /etc/selinux/config file and rebooting the system. Take the public IP and paste in the browser: you could see the page.



EBS Snapshots

Unmount the volume we created above as below:

```
root@ip-172-31-26-160:~#  
root@ip-172-31-26-160:~# umount /var/www/html/images  
root@ip-172-31-26-160:~#  
root@ip-172-31-26-160:~# df -h  
Filesystem      Size  Used Avail Use% Mounted on  
/dev/root       7.6G  2.4G  5.3G  31% /  
tmpfs          475M    0   475M   0% /dev/shm  
tmpfs          190M  868K  190M   1% /run  
tmpfs          5.0M    0   5.0M   0% /run/lock  
/dev/xvda15     105M  6.1M   99M   6% /boot/efi  
tmpfs          95M  4.0K   95M   1% /run/user/1000  
root@ip-172-31-26-160:~#
```

Next, navigate to the AWS console, choose the volume that was created, and proceed to detach it by selecting the respective option.

Create a new volume same as shown above and attach it to the same ec2 instance.

Follow the same steps as mentioned above to create a new partition, format , and then create a directory

```
root@ip-172-31-26-160:~#  
root@ip-172-31-26-160:~# mkdir -p /var/lib/mysql  
root@ip-172-31-26-160:~#  
root@ip-172-31-26-160:~#
```

Add the new volume to the /etc/fstab.

```
root@ip-172-31-26-160:~# cat /etc/fstab
LABEL=cloudimg-rootfs      /          ext4    discard,errors=remount-ro      0 1
LABEL=UEFI                  /boot/efi    vfat    umask=0077      0 1
/dev/xvdbf1                /var/lib/mysql ext4    defaults      0 0

```

Then mount it.

```
root@ip-172-31-26-160:~# mount -a
root@ip-172-31-26-160:~#
root@ip-172-31-26-160:~#
root@ip-172-31-26-160:~# df -h
Filesystem      Size  Used Avail Use% Mounted on
/dev/root       7.6G  2.4G  5.3G  31% /
tmpfs           475M   0    475M   0% /dev/shm
tmpfs           190M  868K  190M   1% /run
tmpfs           5.0M   0    5.0M   0% /run/lock
/dev/xvda15     105M  6.1M  99M   6% /boot/efi
tmpfs           95M  4.0K  95M   1% /run/user/1000
/dev/xvdbf1     2.9G  24K  2.8G   1% /var/lib/mysql
root@ip-172-31-26-160:~#
root@ip-172-31-26-160:~#
```

Install mariadb-server package with the below steps:

- *apt*
 - *install*
 - *mariadb-*
 - *server*
- systemctl*
start
mariadb
systemctl
enable
mariadb
systemctl
status
mariadb

List the /var/lib/mysql contents:

- `ls /var/lib/mysql/`

```
[root@ip-172-31-26-140 ~]# ls /var/lib/mysql/
error.log.00000001  error_log_control  ib_recovery.log  debian-10.6.flag  ib_buffer_pool  ib_logfile0  ibdata1  ibtmp1  lost+found  multi-master.info  mysql  mysql_upgrade_info  performance_schema  sys
[root@ip-172-31-26-140 ~]#
```

Snapshot functionality primarily serves as a backup and restoration tool, crucial for mitigating data loss in the event of system failures. However, snapshots offer more than just backup capabilities.

In the event of data loss, the initial step involves unmounting the partition to prevent data overwriting. If the goal is to restore the existing partition, snapshots may not be applicable, as they provide a new volume containing the data. In such cases, unmounting the partition and employing recovery tools become essential.

When utilizing snapshots, unmounting the partition is the preliminary action, followed by detaching the corrupted volume. Subsequently, a new volume is created from the snapshot, retaining all data and partition configurations. Upon creation, the new volume is ready for attachment and mounting, effectively replacing the corrupted volume.

In essence, snapshots facilitate data restoration through volume replacement rather than repair, streamlining the recovery process and ensuring data integrity.

Take a Snapshot:

- Navigate to the Volume section in the AWS console.
- Select the relevant volume, then click on “Actions” and choose “Create Snapshot.” Provide a meaningful description and add tags for easy identification.

Click on “Create Snapshot” and wait for it to enter the completed state.

Multicloud Devops by veera nareshit

The image shows three screenshots of the AWS Management Console related to EBS volumes and snapshots.

Volumes (1/2) Info: A table showing two volumes. The first volume, "Web-Servr_Web-Volume", has a Volume ID of "vol-05060fb3eeab45af", Type "gp2", Size 8 GiB, IOPS 100, Throughput -, Snapshot "snap-057fda6...", and Created "2024/04/12 10:35 GMT+5:30". The second volume has a Volume ID of "vol-045649e72d6bf20e0", Type "gp3", Size 100 GiB, IOPS 1000, Throughput 125, Snapshot "-", and Created "2024/04/15 06:35 GMT+5:30". The Actions menu is open, showing options like "Create snapshot", "Create snapshot lifecycle policy", "Delete volume", "Attach volume", "Detach volume", "Force detach volume", "Manage auto-enabled I/O", "Manage tags", and "Fault injection".

Create snapshot: A modal dialog titled "Create snapshot" with "Info" link. It asks to "Create a point-in-time snapshot to back up the data on an Amazon EBS volume to Amazon S3." It contains fields for "Volume ID" (set to "vol-045649e72d6bf20e0"), "Description" (set to "db-volume-snapshot"), and "Encryption" (set to "Not encrypted"). Below the form is a "Tags" section with a note about tags being optional for cost tracking. It shows "No tags associated with the resource" and a "Add tag" button. At the bottom are "Cancel" and "Create snapshot" buttons.

Snapshots (1) Info: A table showing one snapshot. The snapshot has a Snapshot ID of "snap-01150462fb30338", Volume size 100 GiB, Description "db-volume-snapshot", Storage tier "Standard", Status "Completed", Started "2024/04/15 07:18 GMT+5:30", Progress "Available (100%)", and Encryption "Not encrypted". The Actions menu is open, showing options like "Recycle Bin", "Actions", "Create snapshot", and "Delete snapshot".

Delete Data:

- If you accidentally delete data or encounter corruption, proceed with data deletion.
- Access the directory containing the data you wish to remove, e.g., /var/lib/mysql. Use the command `rm -rf *` to remove all files and directories within the specified location.

Multicloud Devops by veera nareshit

Multicloud Devops by veera nareshit

```
mysql#  
+  
didi_recovery.log debian-10.4.flag ib_buffer_pool ib_logfile0 ibdata1 ibtmp1 lost+found multi-master.info mysql mysql_op  
+  
zm -xf *
```

Unmount:

- Stop any relevant services to prevent further data overwrites.
- Unmount the affected directory, e.g., /var/lib/mysql.

```
root@ip-172-31-26-160:/var/lib/mysql# zm -xf *  
root@ip-172-31-26-160:/var/lib/mysql# df -h  
Filesystem      Size  Used  Avail Use% Mounted on  
/dev/root       7.6G  2.5G  5.1G  33% /  
tmpfs          475M     0  475M   0% /dev/shm  
tmpfs          190M  880K  189M   1% /run  
tmpfs          5.0M    0  5.0M   0% /run/lock  
/dev/xvda15     105M  6.1M  99M   6% /boot/efi  
/dev/xvdb1      2.9G 124M  2.6G   1% /var/lib/mysql  
tmpfs          95M  4.0K  95M   1% /run/user/1000  
root@ip-172-31-26-160:/var/lib/mysql#  
root@ip-172-31-26-160:/var/lib/mysql#  
root@ip-172-31-26-160:/var/lib/mysql#  
root@ip-172-31-26-160:/var/lib/mysql# systemctl stop mariadb  
root@ip-172-31-26-160:/var/lib/mysql#  
root@ip-172-31-26-160:/var/lib/mysql# systemctl status mariadb  
● mariadb.service - MariaDB 10.6.16 database server  
   Loaded: loaded (/lib/systemd/system/mariadb.service; enabled; vendor preset: enabled)  
   Active: inactive (dead) since Mon 2024-04-15 01:51:02 UTC; 5s ago  
     Docs: man:mariadb(8)  
           https://mariadb.com/kb/en/library/systemd/  
   Process: 4733 ExecStart=/usr/sbin/mariadb -S /var/run/mariadb/mariadb.sock (code=exited, status=0/SUCCESS)  
   Main PID: 4733 (code=exited, status=0/SUCCESS)  
     Status: "MariaDB server is down"  
        CPU: 367ms
```

```
root@ip-172-31-26-160:~#  
root@ip-172-31-26-160:~# umount /var/lib/mysql/  
root@ip-172-31-26-160:~#  
root@ip-172-31-26-160:~#  
root@ip-172-31-26-160:~#
```

Detach Corrupted Volume:

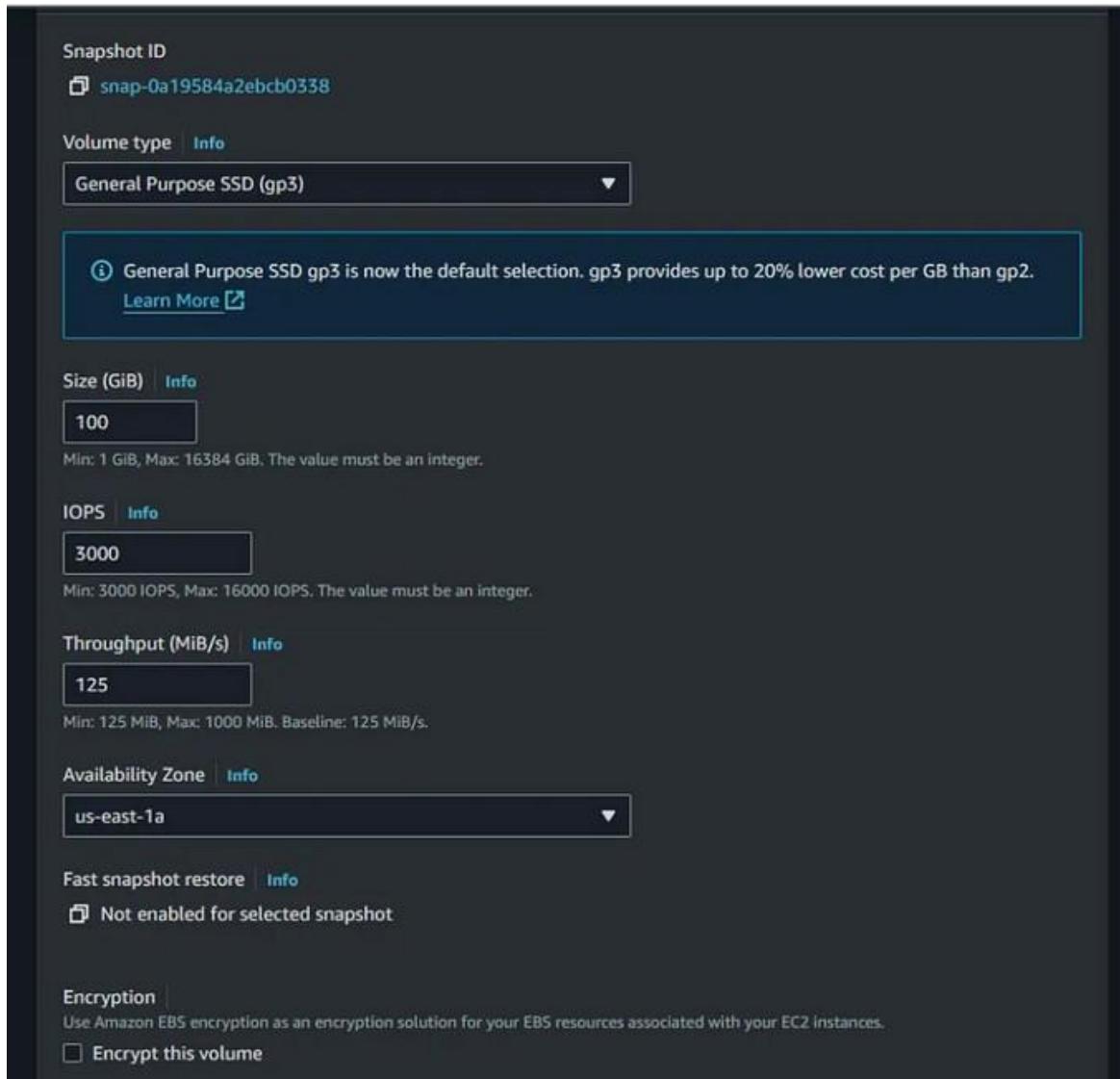
- Navigate to the Volumes section and locate the corrupted volume. Rename the volume for easy identification, e.g., “corrupted.”

Select “Actions” and choose “Detach Volume” to disconnect it.

| Volumes (1/2) <small>Info</small> | | | | | | | | | | | <small>Create volume</small> |
|-----------------------------------|-----------------------|------|--------|------|------------|-----------------|---------------------------|-------------------|--------------|------------------------|------------------------------|
| <small>Search</small> | | | | | | | | | | | <small>< 1 ></small> |
| Name | Volume ID | Type | Size | IOPS | Throughput | Snapshot | Created | Availability Zone | Volume state | Actions | |
| Web-Serve_Root-Volume | vol-05060db5e6ab45bf | gp2 | 8 GB | 100 | - | snap-057fdaf... | 2024/04/12 10:35 GMT+5... | us-east-1c | In-use | <small>Actions</small> | <small>⋮</small> |
| Corrupted | vol-0456499e72f6b720d | gp3 | 100 GB | 5000 | 125 | - | 2024/04/15 06:35 GMT+5... | us-east-1c | In-use | <small>Actions</small> | <small>⋮</small> |

Recover from Snapshot:

- Head to the Snapshots section and select the desired snapshot.
- Click on “Actions” and choose “Create snapshots from volume”.
- Customize volume settings such as type, size, and zone if needed.
- Add appropriate tags for clarity
- Click on “Create Volume” to finalize the recovery process.

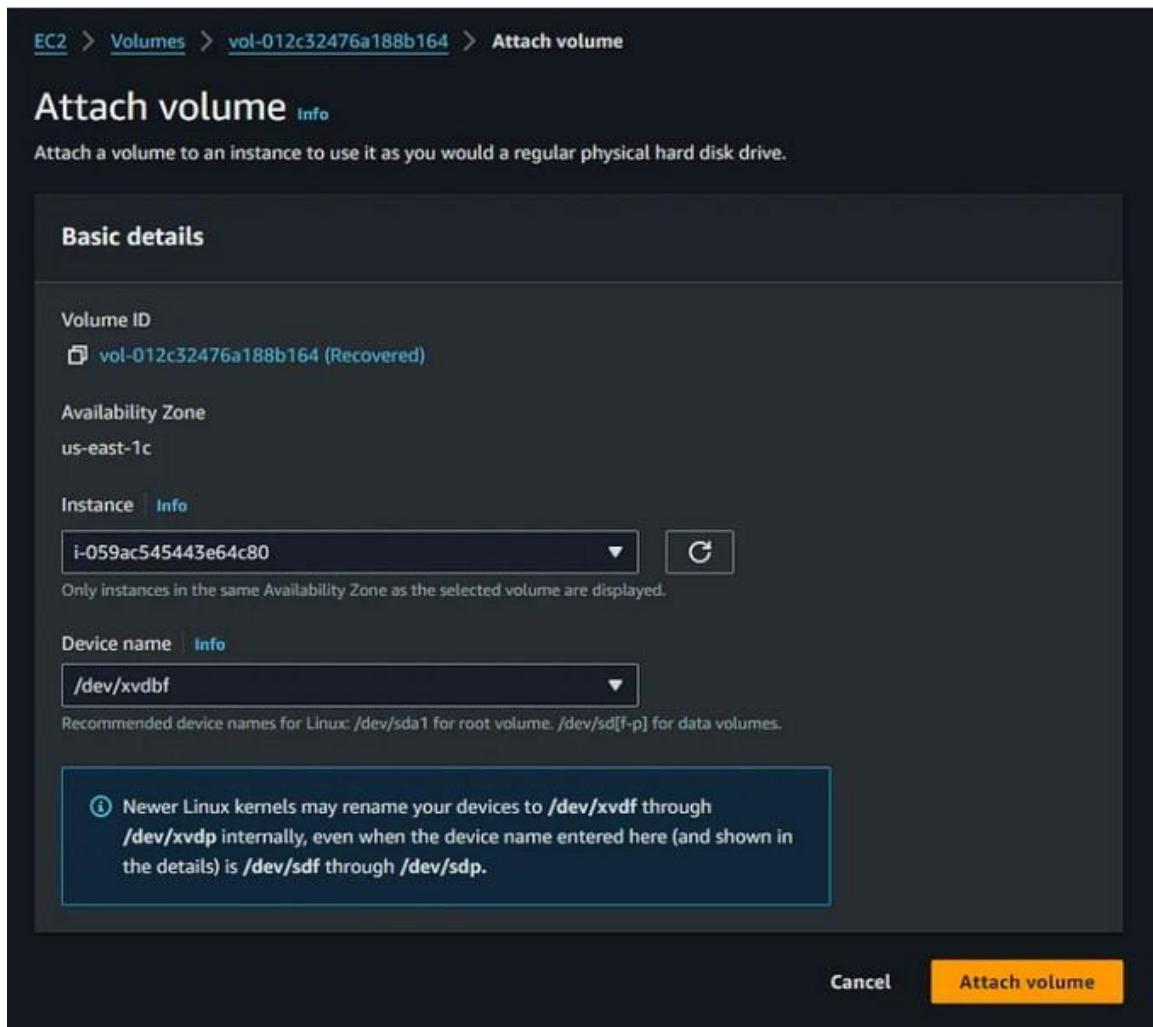


In the Volumes console. goto Volumes and rename the new volume as “recovered”.

| Volumes (1/3) Info | | | | | | | | | | | Actions | Create volume | ... |
|------------------------------------|-----------------------|------|---------|------|------------|-----------------|---------------------------|-------------------|--|--|-------------------------|-------------------------------|---------------------|
| <input type="text"/> Search | | | | | | | | | | | Actions | Create volume | ... |
| Name | Volume ID | Type | Size | IOPS | Throughput | Snapshot | Created | Availability Zone | Volume state | | | | |
| Web-Serve_Root-Volume | vol-05060db5eabbd45bf | gp2 | 8 GiB | 100 | - | snap-057fdaf... | 2024/04/12 10:35 GMT+5... | us-east-1c | <input checked="" type="radio"/> In-use | | | | |
| Corrupted | vol-045649e7268bf720d | gp3 | 100 GiB | 3000 | 125 | - | 2024/04/15 06:35 GMT+5... | us-east-1c | <input checked="" type="radio"/> Available | | | | |
| Recovered | vol-012c13476a1685f64 | gp3 | 100 GiB | 3000 | 125 | snap-0a19584... | 2024/04/15 07:50 GMT+5... | us-east-1c | <input checked="" type="radio"/> Available | | | | |

Multicloud Devops by veera nareshit

Goto actions and attach the volume



Below, we examined the disks and found that they were not mounted. Consequently, we executed the mount -a command to rectify the issue. Upon verification, the disks were automatically mounted, allowing us to recover the deleted data by executing ls /var/lib/mysql/.

```
root@ip-172-31-26-160:~# df -h
Filesystem      Size  Used Avail Use% Mounted on
/dev/eboot     7.6G  2.5G  5.1G  33% /
tmpfs          473M   6  473M   0% /dev/shm
tmpfs          195M  872K  189M  1% /run
tmpfs          5.0M    0  5.0M  0% /run/lock
/dev/sda15     105M  6.1K  99M  6% /boot/efi
tmpfs          95M  4.0K  95M  1% /run/user/1000
/dev/xvdbf      2.9G 122M  2.6G  5% /var/lib/mysql
root@ip-172-31-26-160:~# mount -a
root@ip-172-31-26-160:~# df -h
Filesystem      Size  Used Avail Use% Mounted on
/dev/eboot     7.6G  2.5G  5.1G  33% /
tmpfs          473M   6  473M   0% /dev/shm
tmpfs          195M  872K  189M  1% /run
tmpfs          5.0M    0  5.0M  0% /run/lock
/dev/sda15     105M  6.1K  99M  6% /boot/efi
tmpfs          95M  4.0K  95M  1% /run/user/1000
/dev/xvdbf      2.9G 122M  2.6G  5% /var/lib/mysql
root@ip-172-31-26-160:~# ls /var/lib/mysql/
aria_log_00000001  aria_log_control  dd1_recovery.log  debian-10.6.flag  ib_buffer_pool  ib_logfile0  ibdata1  ibtmp1  lost+found  multi-master.info  mysql  mysql_upgrade_info  performance_schema
```

In the snapshot console, you'll find various options beyond basic backup functionality. These include deleting snapshots, creating volumes from existing snapshots (which we've already covered), and managing fast snapshot restores for large volumes, albeit with additional charges.

Additionally, you can leverage snapshots to create images, known as Amazon Machine Images (AMIs). If the snapshot originates from a root volume, you can directly create an image from it.

Another useful feature is the ability to copy snapshots. By clicking the “copy” button, you can duplicate snapshots to different regions, which is particularly handy for multisite setups. Moreover, during the copying process, you have the option to enable encryption for added security.

Furthermore, the snapshot interface offers a “modify permissions” button, allowing you to adjust access settings. This enables you to make snapshots public or share them with specific AWS accounts by providing their account IDs.

In summary, snapshots offer a simple yet powerful toolset with various benefits, making them an indispensable asset for managing AWS resources effectively.

Elastic

-
- Lo
ad
Bal
an
cer
Lo
ad
bal
anc
ing

is a

c

rucial aspect of managing clusters of servers efficiently. It involves distributing incoming network traffic across multiple servers to ensure optimal resource utilization, reliability, and scalability.

- In AWS, Elastic Load Balancer (ELB) provides a scalable and reliable solution for load balancing, allowing users to distribute traffic across multiple targets seamlessly.

• Key Components of ELB:

■ Frontend Port and Backend Port:

- The frontend port is the port on which the load balancer listens for incoming traffic. For example, for HTTPS connections, the frontend port would typically be 443.
- The backend port is the port on which the load balancer forwards traffic to the backend servers. It corresponds to the port where the application/service is running on the backend servers. For instance, if a Tomcat server is running on port 8080 on backend instances, the backend port would be set to 8080.

■ Types of ELB:

- Classic Load Balancer: Operates at the network layer (Layer 4) of the OSI model. It provides basic load balancing capabilities and distributes incoming traffic to backend instances based on configured rules.
- Application Load Balancer (ALB): Operates at the application layer (Layer 7) of the OSI model. ALB supports advanced routing based on URL paths, hostnames, and other application-level attributes, making it ideal for HTTP and HTTPS traffic.
- Network Load Balancer (NLB): Operates at the transport layer (Layer 4) of the OSI model. NLB is designed to handle extremely high volumes of traffic and offers ultralow latency, making it suitable for TCP, UDP, and TLS traffic.

■ Load Balancer Targets:

- ELB distributes traffic to backend targets, which can include EC2 instances, containers, IP addresses, or Lambda functions.

- Targets are distributed across multiple availability zones to improve fault tolerance and high availability.

Benefits of ELB:

- Scalability: ELB automatically scales in response to incoming traffic, ensuring that applications remain responsive even during traffic spikes.
- Fault Tolerance: By distributing traffic across multiple targets and availability zones, ELB enhances fault tolerance and resiliency.
- Simplified Management: ELB abstracts away the complexities of managing load balancers, allowing users to focus on building and scaling their applications.

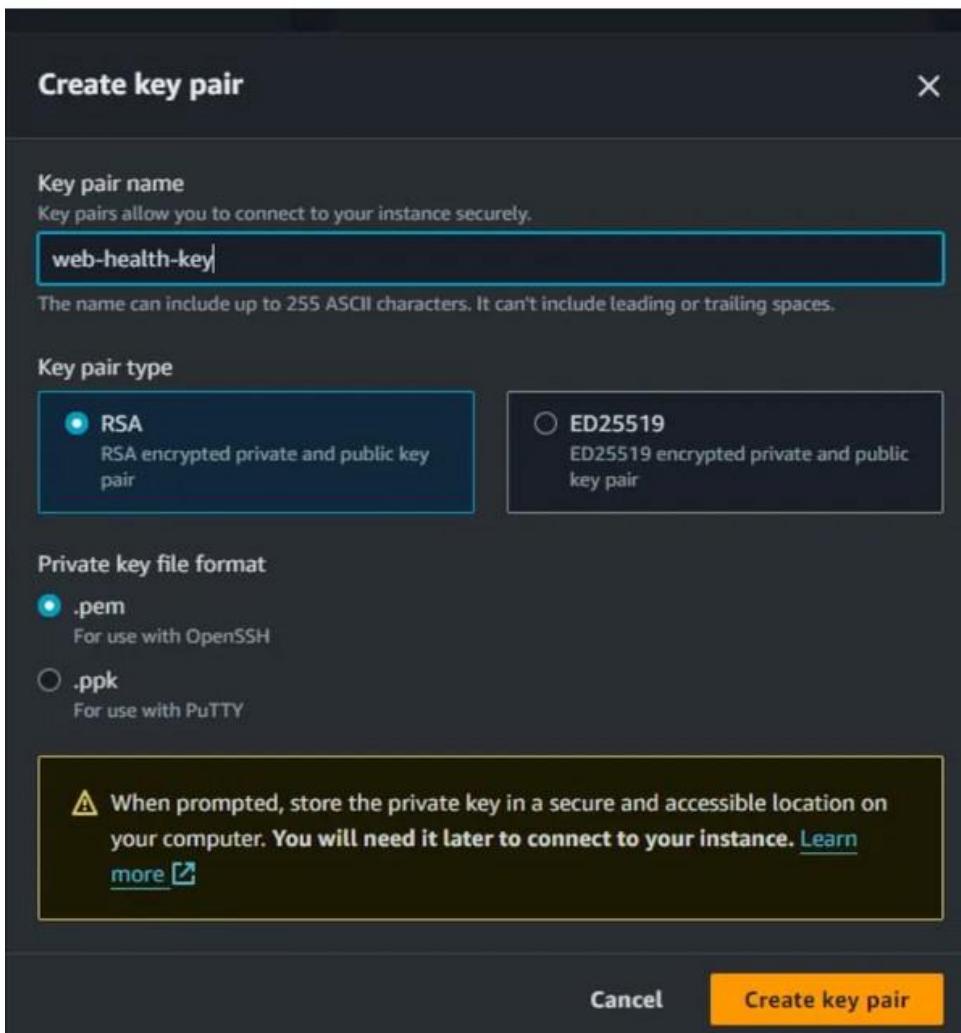
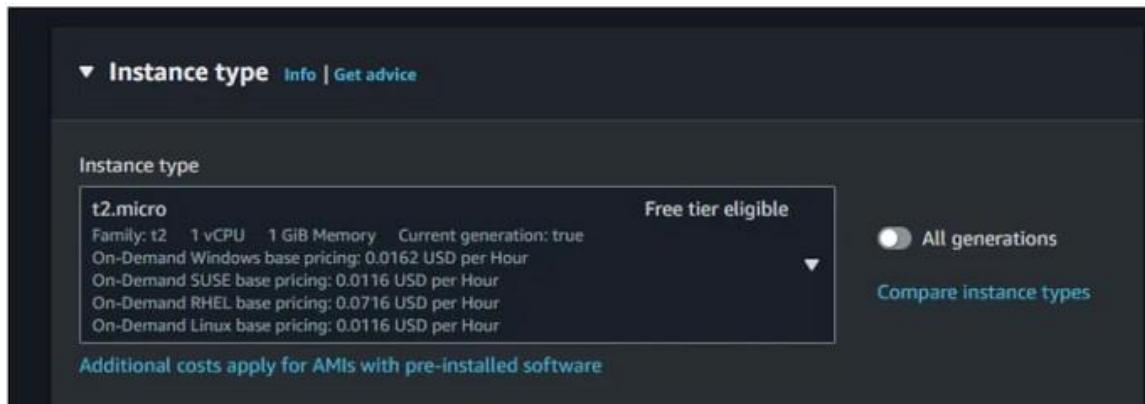
Let's proceed with launching an instance and naming it as "web-health".

You have the flexibility to choose between Amazon Linux, Ubuntu, or CentOS. The script provided in the resource section is compatible with both RPM-based and Debian-based operating systems. Since Amazon Linux is RPM-based, we'll stick with it for this instance.

Multicloud Devops by veera nareshit

The screenshot shows the 'Launch an instance' wizard in the AWS Management Console. The first step, 'Name and tags', has a 'Name' field containing 'web-health'. The second step, 'Application and OS Images (Amazon Machine Image)', includes a search bar and a 'Quick Start' section with icons for Amazon Linux, macOS, Ubuntu, Windows, Red Hat, and SUSE Linux. A 'Browse more AMIs' link is also present.

Let's opt for a t2.micro instance type and create a new key pair named "web-healthkey" for secure access.



Now, for the security group, let's create a new one named “web-health-sg”, allowing inbound traffic on port 22 for SSH from your IP, and an additional rule for port 80, allowing traffic only from your IP.

Multicloud Devops by veera nareshit

Security group name - required

web-health-sg

This security group will be added to all network interfaces. The name can't be edited after the security group is created. Max length is 255 characters. Valid characters: a-z, A-Z, 0-9, spaces, and _-./()#,@[]+=&:;"\$*

Description - required Info

web-health-sg

Inbound Security Group Rules

▼ Security group rule 1 (TCP, 22, 0.0.0.0/0)

Type Info Protocol Info Port range Info
ssh TCP 22

Source type Info Source Info Description - optional Info
Anywhere Add CIDR, prefix list or security e.g. SSH for admin desktop
0.0.0.0/0 X

▼ Security group rule 2 (TCP, 80, 0.0.0.0/0)

Type Info Protocol Info Port range Info
Custom TCP TCP 80

Source type Info Source Info Description - optional Info
Anywhere Add CIDR, prefix list or security e.g. SSH for admin desktop
0.0.0.0/0 X

Add security group rule

EC2 > Security Groups > sg-05e4c5ab0d4caeb - web-health-sg > Edit inbound rules

Edit inbound rules Info

Inbound rules control the incoming traffic that's allowed to reach the instance.

Inbound rules Info

| Security group rule ID | Type Info | Protocol Info | Port range Info | Source Info | Description - optional Info | Delete |
|------------------------|-----------|---------------|-----------------|-------------|-----------------------------|---------------------------------------|
| sg-r1bda5d0c8ae719c | HTTP | | Custom | Anywhere | 103.201.72.75/32 | <input type="button" value="Delete"/> |
| - | SSH | | Anywhere | | 0.0.0.0/0 | <input type="button" value="Delete"/> |

Add rule

Moving forward, let's proceed with the default storage settings and navigate to advanced details. Here, we'll utilize the script in the "User data" section under "Advanced details". Simply copy and paste it as instructed.

Script:

```
#!/bin/bash

# Variable Declaration
#PACKAGE="httpd wget unzip"
#SVC="httpd"
URL='https://www.tooplate.com/zip-templates/2098_health.zip'
ART_NAM

E='2098_'

health'

TEMPDIR=

"/tmp/we

bfiles"

yum --help

&>

/dev/null

if [ $? -eq 0 ]
then
# Set Variables for CentOS
PACKAGE="httpd wget
unzip" SVC="httpd"

echo "Running Setup on CentOS"
# Installing Dependencies
echo "#####
echo "Installing packages."
echo "#####
sudo yum install
$PACKAGE -y > /dev/null
echo

# Start & Enable Service
```

```
echo
#####
##### echo "Start & Enable HTTPD
Service"
echo
#####
##### sudo systemctl start $SVC
sudo systemctl enable $SVC echo

# Creating Temp Directory
echo #####
echo "Starting Artifact Deployment"
echo
#####
##### mkdir -p $TEMPDIR cd
$TEMPDIR echo

fi
wget $URL >
/dev/null unzip
$ART_NAME.zip >
/dev/null sudo cp -r
$ART_NAME/*
/var/www/html/
echo

# Bounce Service
echo
#####
##### echo "Restarting HTTPD service"
echo #####
systemc
tl
restart
$SVC
echo
```

```
# Clean Up
echo
#####
##### echo "Removing Temporary Files"
echo
#####
##### rm -rf $TEMPDIR echo

sudo
systemctl
status
$SVC ls
/var/www
w/html/

else
# Set Variables for Ubuntu
PACKAGE="apache2 wget
unzip" SVC="apache2"

echo "Running Setup on CentOS"
# Installing Dependencies
echo #####
echo "Installing packages."
echo
#####
##### sudo apt update
sudo apt install $PACKAGE
-y > /dev/null echo

# Start & Enable Service
echo
#####
##### echo "Start & Enable HTTPD
Service"
echo
#####
#####
```

```
#####" sudo systemctl start $SVC
sudo systemctl enable $SVC echo

# Creating Temp Directory
echo #####
echo "Starting Artifact Deployment"
echo
#####
mkdir -p $TEMPDIR cd
$TEMPDIR echo

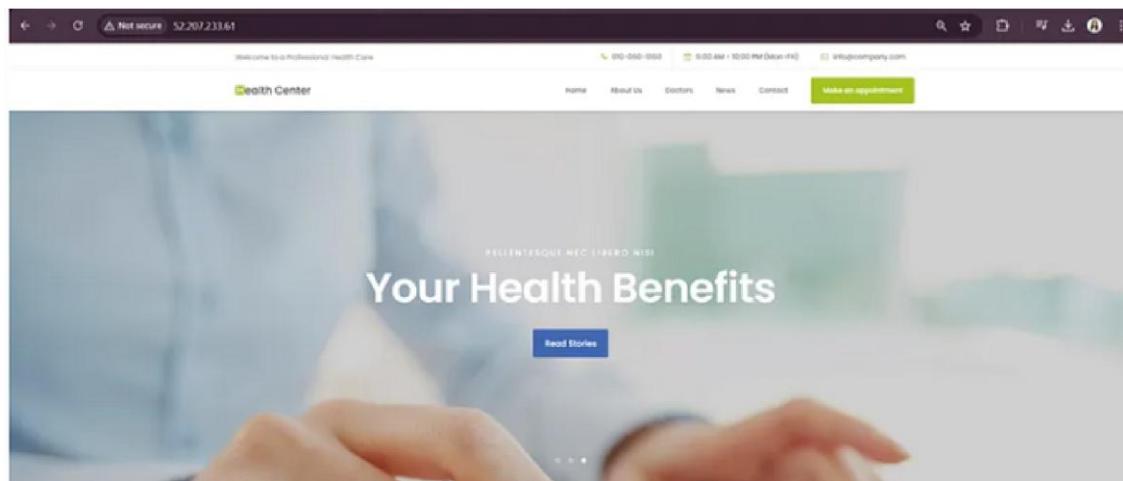
wget $URL >
/dev/null unzip
$ART_NAME.zip >
/dev/null sudo cp -r
$ART_NAME/*
/var/www/html/
echo

# Bounce Service
echo
#####
#echo "Restarting HTTPD service"
echo
#####
systemctl restart $SVC
echo

# Clean Up
echo
#####
#echo "Removing Temporary Files"
echo
#####
rm -rf $TEMPDIR echo
```

```
sudo  
systemctl  
status  
$SVC ls  
/var/www  
w/html/
```

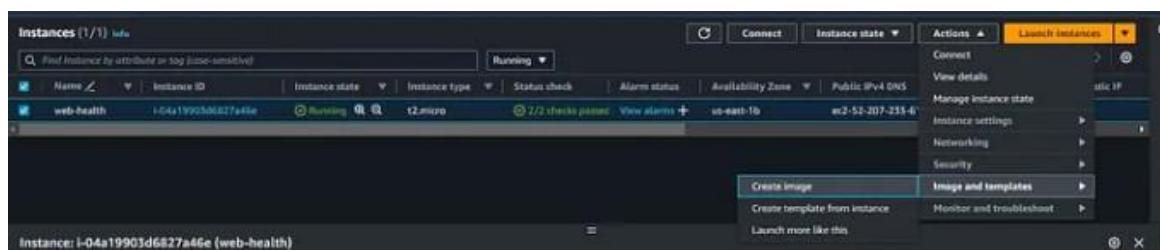
After launching the instance, it typically takes about 5 minutes to initialize. Once it's up, we can check the public IP address and access the website via a browser.



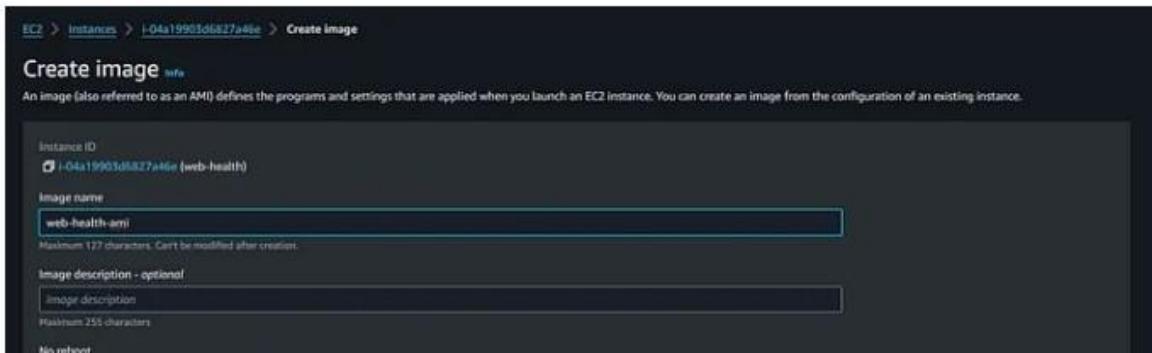
In case you encounter any issues accessing the website, ensure that the security group allows traffic from your IP and that the HTTP service is running on the instance.

Following the successful setup of our website on the EC2 instance, we'll proceed with creating an AMI and a launch template for future deployments.

Creating an AMI involves selecting the instance, navigating to Actions > Image and templates > Create image, and specifying a name for the AMI.

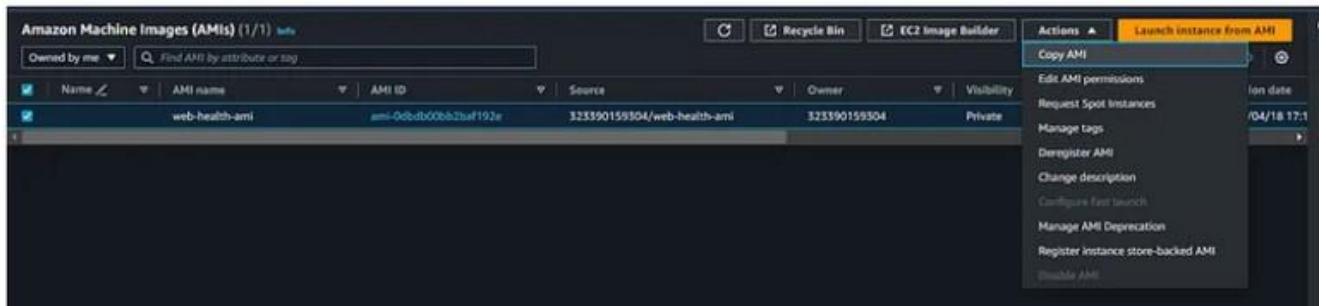


Provide a name for the AMI – Here I gave – “web-health-ami”.



Once the AMI is created, we can leverage it to launch instances with identical configurations, ensuring consistency across our cluster.

Additionally, we can explore options like copying the AMI to other regions or sharing it with other AWS accounts for collaboration.



Also, AWS offers a service known as EC2 Image Builder, allowing you to establish automated pipelines for image creation.

The image pipeline in Image Builder defines all aspects of the process to customize images. It consists of the image recipe, infrastructure configuration, distribution, and test settings.

With the AMI in place, we can streamline instance launches by creating a launch template. Given name as “web-health-template” and version as “V1”

The screenshot shows the 'Create launch template' wizard in the AWS Management Console. The current step is 'Launch template name and description'. The 'Launch template name' field is filled with 'web-health-template'. The 'Template version description' field contains 'V1'. Under 'Auto Scaling guidance', there is a checkbox labeled 'Provide guidance to help me set up a template that I can use with EC2 Auto Scaling'. Below the main form, there are two expandable sections: 'Template tags' and 'Source template'.

EC2 > Launch templates > Create launch template

Create launch template

Creating a launch template allows you to create a saved instance configuration that can be reused, shared and launched at a later time. Templates can have multiple versions.

Launch template name and description

Launch template name - *required*

web-health-template

Must be unique to this account. Max 128 chars. No spaces or special characters like '&', '"', '@'.

Template version description

V1

Max 255 chars

Auto Scaling guidance | [Info](#)

Select this if you intend to use this template with EC2 Auto Scaling

Provide guidance to help me set up a template that I can use with EC2 Auto Scaling

▶ [Template tags](#)

▶ [Source template](#)

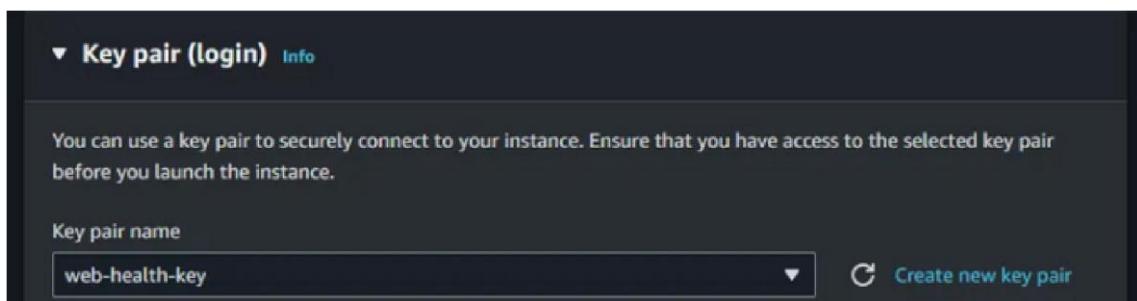
Then choose the “My AMIs” and select the AMI we created just before “web-healthami”

The screenshot shows the AWS Lambda console with the 'Application and OS Images (Amazon Machine Image)' section selected. A search bar at the top allows users to search for AMIs. Below the search bar are three tabs: 'Recents', 'My AMIs' (which is selected), and 'Quick Start'. Under the 'My AMIs' tab, there are three filter options: 'Don't include in launch template', 'Owned by me' (which is selected), and 'Shared with me'. To the right of these filters is a search icon and a link to 'Browse more AMIs', which includes AMIs from AWS, Marketplace, and the Community. Below the filters, a specific AMI is listed: 'web-health-ami' (AMI ID: ami-0dbdb00bb2baf192e, created: 2024-04-18T11:48:31.000Z, Virtualization: hvm, ENA enabled: true, Root device type: ebs). The 'Description' and 'Architecture' sections are also visible.

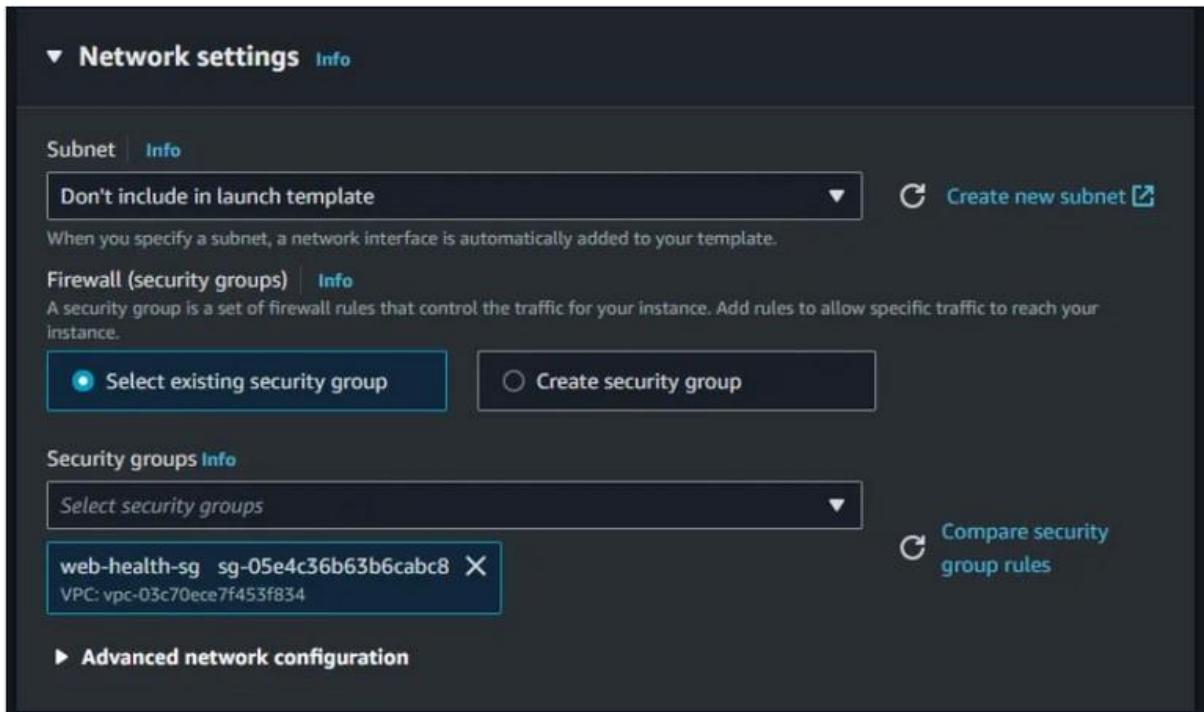
Select the instance type as t2.micro

The screenshot shows the 'Instance type' selection screen. The 't2.micro' instance type is selected, with its details displayed: Family: t2, 1 vCPU, 1 GiB Memory, Current generation: true, On-Demand Windows base pricing: 0.0162 USD per Hour, On-Demand SUSE base pricing: 0.0116 USD per Hour, On-Demand RHEL base pricing: 0.0716 USD per Hour, On-Demand Linux base pricing: 0.0116 USD per Hour. The 'Free tier eligible' status is shown. There are buttons for 'All generations' and 'Compare instance types'. A note at the bottom states: 'Additional costs apply for AMIs with pre-installed software'.

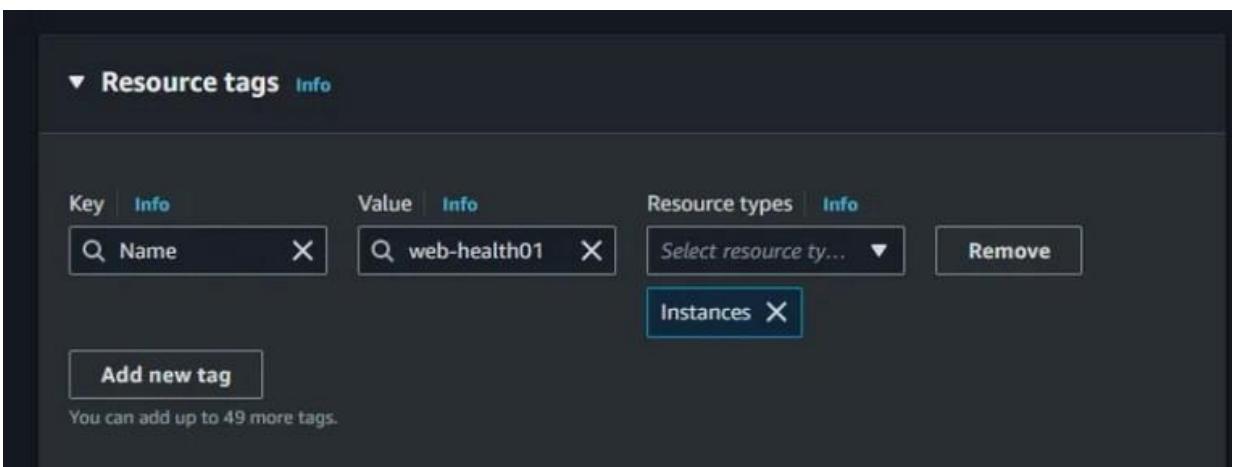
Select our key pair we created for the instance “web-health”

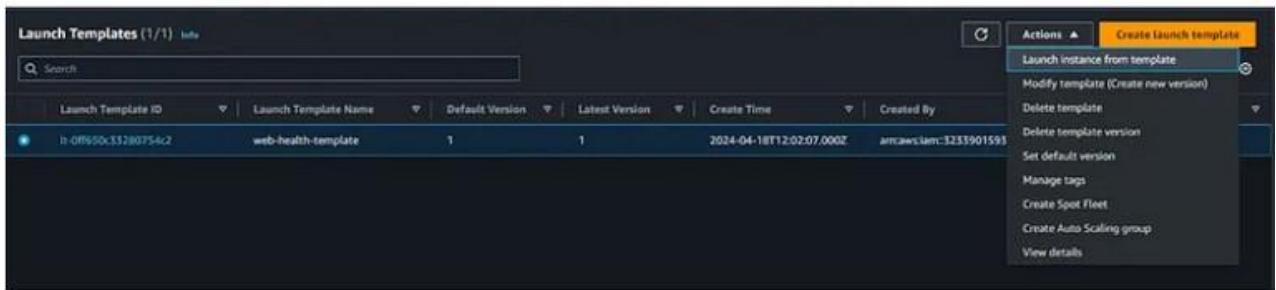


Select the security group as “web-health-sg” we created for our instance “web-health”



Just give as below: When launching an instance from this template, I'll modify the name to follow a format such as 2,3,4 etc.



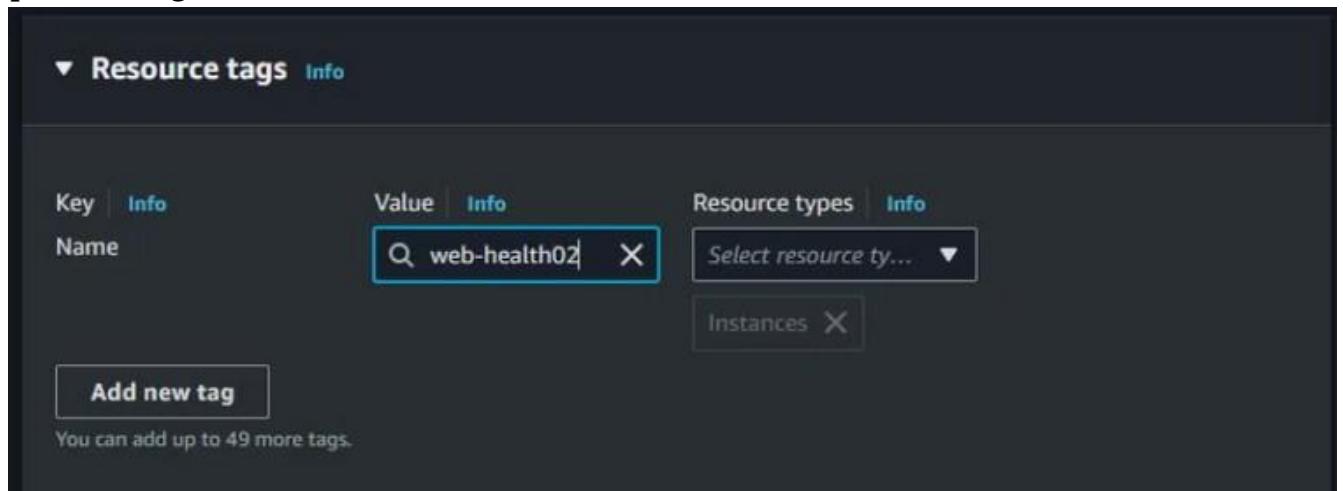


With the AMI in place, we can streamline instance launches by creating a launch template. This template encapsulates all the configuration details, allowing for quick and consistent instance deployments.

By utilizing launch templates, we can expedite the process of launching instances and maintain standardized configurations across our infrastructure.

After configuring our instances, AMI, and launch template, we'll proceed to set up a load balancer to distribute traffic efficiently across our web servers.

Simply navigate to the template section, locate your desired template, and select it. Then, choose the action “launch instance from template.” During the launch process, you’ll have the opportunity to make adjustments as needed. For example, you can modify the tag to “web-health01 to web-health02 two” before proceeding to launch the instance. That’s all there is to it.



Now we have 2 instances which we created from templates and we can create 'n' number of instances upon our business needs.

| Instances (2) Info | | | | | | | | | | | | | | | | | |
|------------------------------------|---------------------|--|---------------|--------------|-------------------------------|-------------------|-------------------------|---------------|------------|--|--|--|--|--|--|--|--|
| | | Find instance by attribute or tag (case-sensitive) | | Running | | | | | | | | | | | | | |
| Name | Instance ID | Instance state | Instance type | Status check | Alarm status | Availability Zone | Public IPv4 DNS | Public IPv4 | Elastic IP | | | | | | | | |
| web-health02 | i-024cd82254f908505 | Running | t2.micro | Initializing | View alarms + | us-east-1b | ec2-54-234-144-31.co... | 54.234.144.31 | - | | | | | | | | |
| web-health | i-04a19903d6827a46e | Running | t2.micro | Initializing | View alarms + | us-east-1b | ec2-52-207-233-61.co... | 52.207.233.61 | - | | | | | | | | |

We have successfully created an AMI and established a launch template utilizing this custom AMI.

Presently, we have two web servers operating here. These servers are not individually accessible by users. Instead, users must access them through a unified endpoint, which directs requests to either of these instances. This single endpoint is facilitated by a load balancer.

In the Load Balancer section, we will first navigate to target groups. Target groups function as collections of instances, each subject to health checks. Essentially, they are groups comprising instances. Let's proceed to create a target group. Starting with creating a target group, we'll define health checks and specify the instances to include.

Choose target type as "Instances"

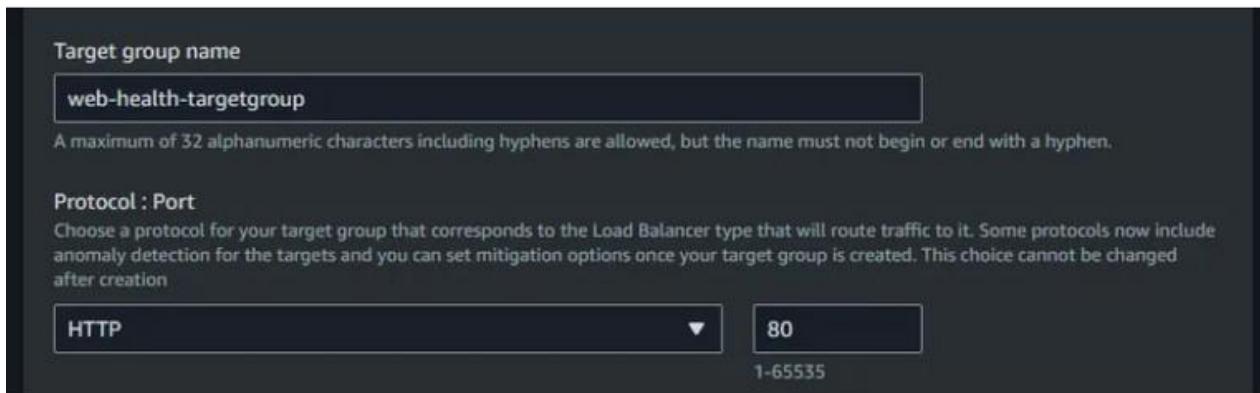
Basic configuration
Settings in this section can't be changed after the target group is created.

Choose a target type

Instances

- Supports load balancing to instances within a specific VPC.
- Facilitates the use of [Amazon EC2 Auto Scaling](#) to manage and scale your EC2 capacity.

Provide a name as below and Port number 80, protocol HTTP.



In our setup, we've implemented health checks within the target group. These health checks allow the target group to assess the status of each instance. If an instance is deemed unhealthy during these checks, the target group will refrain from routing requests to that specific instance.

Now, you might wonder, how does the target group determine if an instance is healthy or not? It's a similar process to how we, as humans, assess the health of a website. We simply open the website in a browser and check if it's functioning properly. If everything looks good, we consider the website to be healthy.

Similarly, the target group conducts health checks by sending HTTP requests to the instances. When configuring these health checks, we specify the protocol (HTTP) and the health check path (/).

Let's break down what the health check path (/) signifies. Imagine your website's URL is <http://IP:80>. The slash at the end of this URL represents the root directory of the website. It's where the homepage resides.

Now, if your website has subfolders, such as "videos" or "images," the health check path would be adjusted accordingly (e.g., /videos or /images).

In essence, by specifying the health check path as "/", we're instructing the target group to check the root directory of our website. When the target group receives a successful response (HTTP 200 OK) from this path, it considers the instance to be healthy.

The screenshot shows the 'Health checks' configuration page. At the top, it says 'Health checks' and 'The associated load balancer periodically sends requests, per the settings below, to the registered targets to test their status.' Below this, there's a section for 'Health check protocol' with a dropdown menu set to 'HTTP'. Under 'Health check path', there's a text input field containing '/' with the placeholder 'Up to 1024 characters allowed.' At the bottom of this section is a link '► Advanced health check settings'.

In our setup, traffic is expected to flow through the default port, which is port 80. However, if your website operates on a different port, you would need to specify it here.

To do so, you would select the “override” option and provide the appropriate port number, such as 8090. However, since we know that our website operates on port 80, we will leave it as is.

This screenshot shows the expanded 'Advanced health check settings' section. It includes a 'Restore defaults' button at the top right. Below it is a 'Health check port' section with a note about the default behavior and an 'Override' toggle switch. The 'Traffic port' option is selected, indicated by a blue dot next to the radio button.

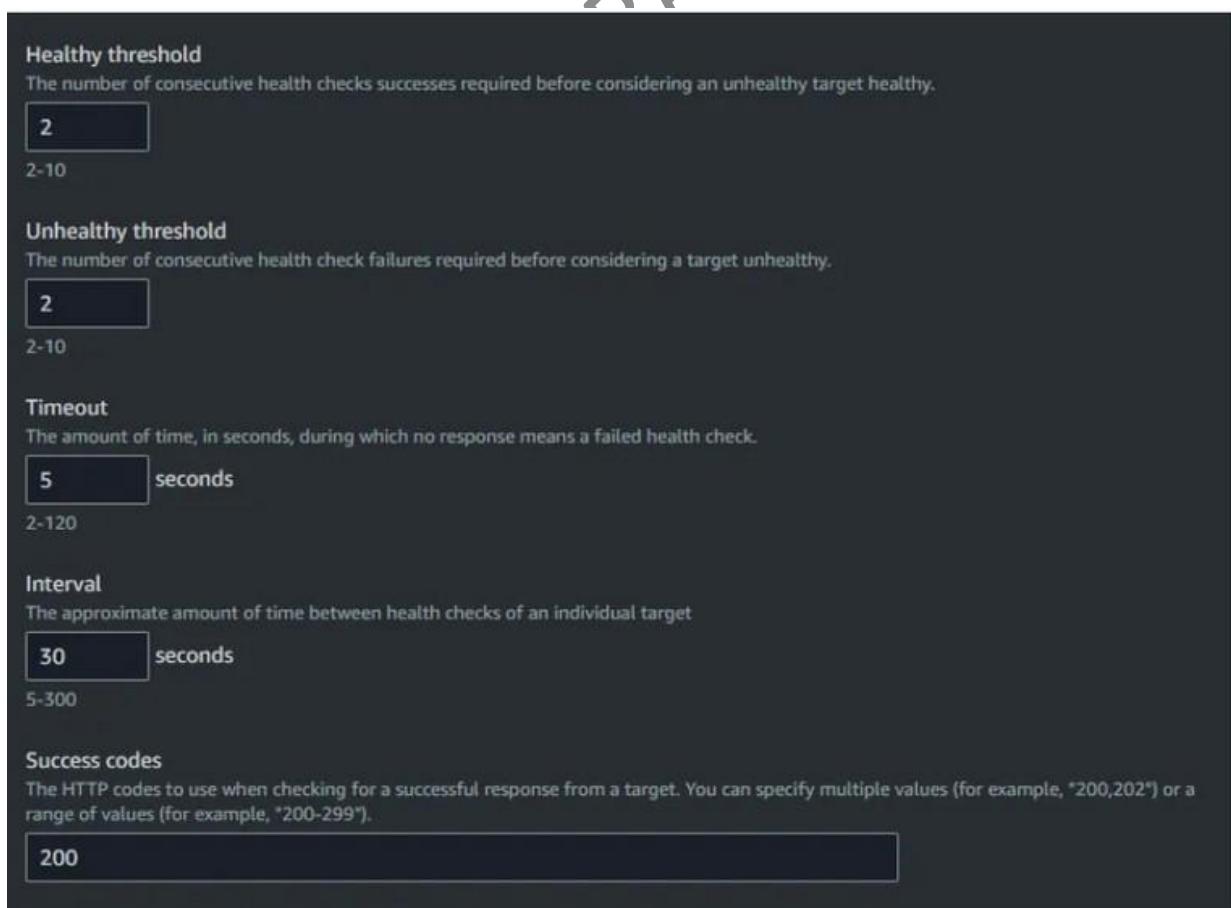
The “Healthy threshold” parameter determines the number of consecutive successful health checks required before declaring the instance as healthy. By default, it is set to five, but you can adjust it to two or ten as per your preference. With a threshold of five, the system checks the instance’s health five times, and if all checks pass successfully, the instance is marked as healthy.

In our case, I will reduce it to two. This means that the instance will be considered healthy if it passes two consecutive health checks. Similarly, the “Unhealthy threshold” specifies the number

of consecutive failed health checks required to declare the instance as unhealthy. If an instance returns an unhealthy status twice in a row, it will be marked as unhealthy.

The “Timeout” parameter determines how long the system waits for a response from the website during each health check. If the website does not respond within the specified timeout period (default is 5 seconds), the system moves on to the next check. Health checks, whether for determining the instance’s health or unhealthiness, occur at regular intervals, typically every 30 seconds. This interval can be adjusted within a range of 5 to 300 seconds.

When performing health checks, the system looks for a success exit code, typically HTTP 200, indicating that the website is functioning properly. While we can visually confirm this by checking the website in a browser, the target group relies on these exit codes to determine the instance’s health status.



A screenshot of a configuration interface for a target group. The interface has a dark background with white text and input fields. It shows the following settings:

- Healthy threshold:** Set to 2. A tooltip indicates: "The number of consecutive health checks successes required before considering an unhealthy target healthy." The input field has a value of 2 and a range of 2-10.
- Unhealthy threshold:** Set to 2. A tooltip indicates: "The number of consecutive health check failures required before considering a target unhealthy." The input field has a value of 2 and a range of 2-10.
- Timeout:** Set to 5 seconds. A tooltip indicates: "The amount of time, in seconds, during which no response means a failed health check." The input field has a value of 5 and a range of 2-120.
- Interval:** Set to 30 seconds. A tooltip indicates: "The approximate amount of time between health checks of an individual target." The input field has a value of 30 and a range of 5-300.
- Success codes:** Set to 200. A tooltip indicates: "The HTTP codes to use when checking for a successful response from a target. You can specify multiple values (for example, "200,202") or a range of values (for example, "200-299")." The input field has a value of 200.

Multicloud Devops by veera nareshit

So the below are the instances running and click on the “Include as pending below option”.

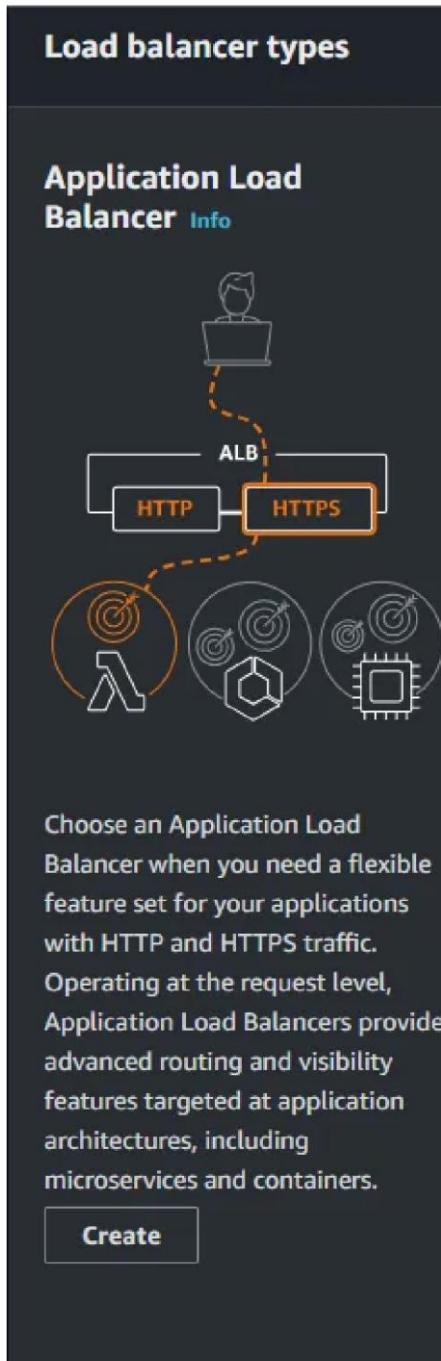
The screenshot shows the 'Register targets' interface. At the top, it says 'Available instances (2/2)'. Below is a table with columns: Instance ID, Name, State, Security groups, Zone, and Private IPv4 address. Two instances are selected: 'i-024c082254f908505' named 'web-health02' and 'i-04a19903d6827a46e' named 'web-health'. Both are in the 'Running' state. The 'Security groups' column shows 'web-health-sg' for both. The 'Zone' column shows 'us-east-1b' for both. The 'Private IPv4 address' column shows '172.31.93.171' for 'web-health02' and '172.31.86.181' for 'web-health'. Below the table, it says '2 selected'. Underneath, there's a section for 'Ports for the selected instances' with a text input field containing '80'. A note says 'Ports for routing traffic to the selected instances' and '1-65535 (separate multiple ports with comma)'. At the bottom right is a button labeled 'Include as pending below'.

And the targets are in running state.

The screenshot shows the 'Review targets' interface. At the top, it says 'Targets (2)'. Below is a table with columns: Instance ID, Name, Port, State, Security groups, Zone, Private IPv4 address, Subnet ID, and Launch time. The same two instances are listed: 'web-health02' and 'web-health'. Both have a port of '80' and are in the 'Running' state. The security group is 'web-health-sg', zone is 'us-east-1b', private IP is '172.31.93.171' and '172.31.86.181' respectively, and launch time is 'April 18, 2024, 17:33 (UTC+05:30)'. At the bottom left, it says '2 pending'. At the bottom right, there are buttons for 'Cancel', 'Previous', and 'Create target group'.

Once the target group is set up, we'll proceed to create an application load balancer (ALB), specifying its name, internet-facing nature, and availability zones.

We will go with application load balancer for http, https traffic.



Give a load balancer name “web-health-elb” and We want it to be available on the Internet, so we’ll keep it Internet facing.

The screenshot shows the 'Create Application Load Balancer' wizard in the AWS Management Console. The current step is 'Basic configuration'. The 'Load balancer name' field is set to 'web-health-elb'. The 'Scheme' section is set to 'Internet-facing'. The 'IP address type' section is set to 'IPv4'. A large watermark 'REVIEW' is diagonally across the page.

EC2 > Load balancers > Create Application Load Balancer

Create Application Load Balancer Info

The Application Load Balancer distributes incoming HTTP and HTTPS traffic across multiple targets such as Amazon EC2 instances, microservices, and containers, based on request attributes. When the load balancer receives a connection request, it evaluates the listener rules in priority order to determine which rule to apply, and if applicable, it selects a target from the target group for the rule action.

▶ How Application Load Balancers work

Basic configuration

Load balancer name

Name must be unique within your AWS account and can't be changed after the load balancer is created.

web-health-elb

A maximum of 32 alphanumeric characters including hyphens are allowed, but the name must not begin or end with a hyphen.

Scheme Info

Scheme can't be changed after the load balancer is created.

Internet-facing

An internet-facing load balancer routes requests from clients over the internet to targets. Requires a public subnet. [Learn more](#)

Internal

An internal load balancer routes requests from clients to targets using private IP addresses.

IP address type Info

Select the type of IP addresses that your subnets use.

IPv4

Includes only IPv4 addresses.

Dualstack

Includes IPv4 and IPv6 addresses.

When setting up the load balancer, it's necessary to specify the availability zones where the load balancer will operate. The system requires a minimum of two zones to ensure redundancy and high availability.

For our configuration, I will choose to select all available zones. This ensures that the load balancer operates across multiple zones, maximizing its availability and fault tolerance. This approach provides a high level of availability at the load balancer level, which is crucial for maintaining continuous service availability.

By selecting all zones, we distribute the load balancing functionality across multiple data centers or regions, reducing the risk of downtime due to failures in any single zone. This enhances the resilience and reliability of our application infrastructure, allowing it to withstand potential disruptions in individual zones.

Network mapping Info
The load balancer routes traffic to targets in the selected subnets, and in accordance with your IP address settings.

VPC Info
Select the virtual private cloud (VPC) for your targets or you can [create a new VPC](#). Only VPCs with an internet gateway are enabled for selection. The selected VPC can't be changed after the load balancer is created. To confirm the VPC for your targets, view your target groups.

vpc-03c70ece7f453fb34
IPv4 VPC CIDR: 172.31.0.0/16

Mappings Info
Select at least two Availability Zones and one subnet per zone. The load balancer routes traffic to targets in these Availability Zones only. Availability Zones that are not supported by the load balancer or the VPC are not available for selection.

- us-east-1a (use1-az1)
- us-east-1b (use1-az2)
- us-east-1c (use1-az4)
- us-east-1d (use1-az6)
- us-east-1e (use1-az3)
- us-east-1f (use1-az5)

Click on the “Create a new security group”

Security groups Info
A security group is a set of firewall rules that control the traffic to your load balancer. Select an existing security group, or you can [create a new security group](#).

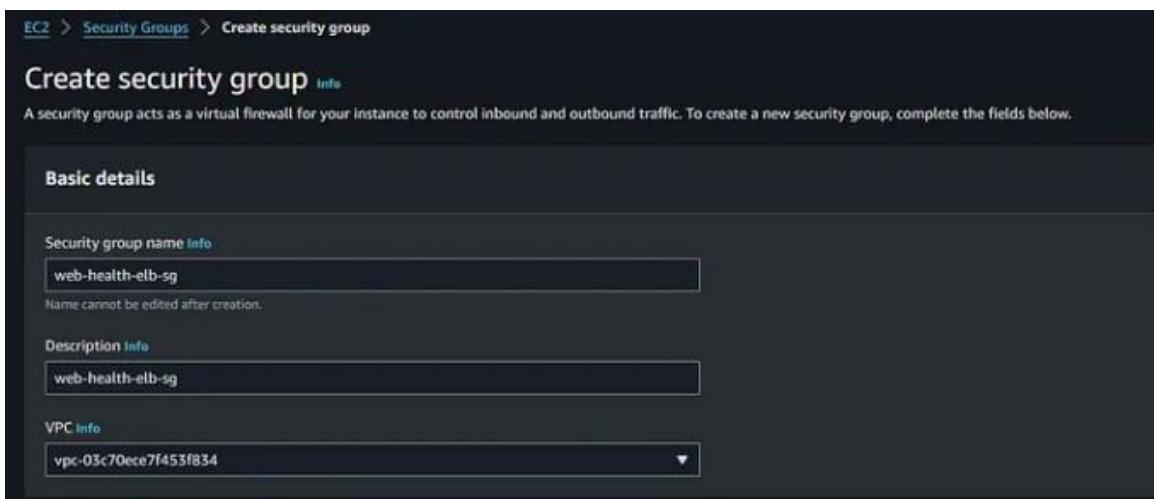
Security groups

Select up to 5 security groups

default sg-0d572019340b0c41e VPC: vpc-03c70ece7f453fb34

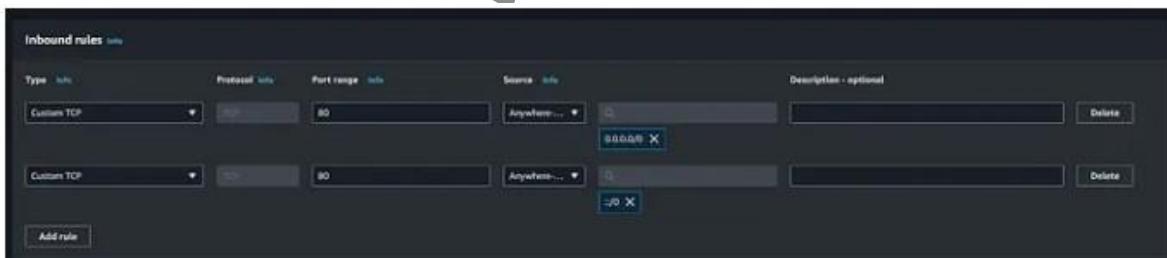
For the security configuration of our load balancer, we'll create a new security group to define its access permissions. It's important to name this security group appropriately for easy identification and management.

We'll name it “web-health-elb-sg” to indicate that it's the security group for our health project's load balancer. Using descriptive names helps maintain clarity and organization within our AWS environment.

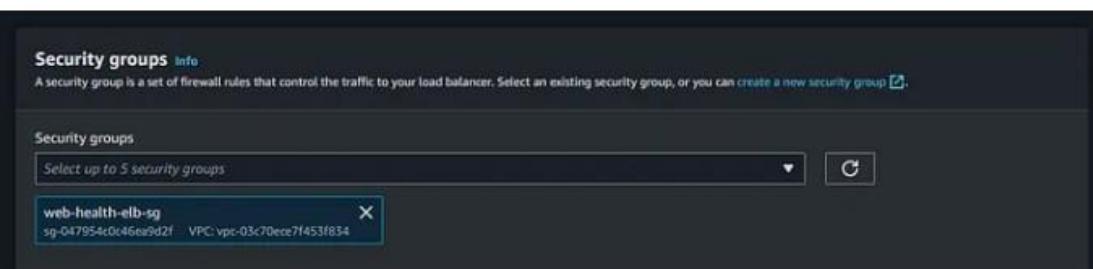


As for the inbound rule, we want to allow access on port 80 from anywhere, making the load balancer accessible to the public. We'll specify the protocol as IPv6 to accommodate modern internet standards, especially considering the increasing adoption of IPv6 by internet service providers, including mobile networks.

Ensuring broad accessibility while maintaining security is crucial for our load balancer's functionality and usability. Making these considerations and adopting best practices, such as proper naming conventions and access control, contributes to the overall robustness of our infrastructure.



Now select the security group we created just now:



Ensuring that the ALB's security group allows inbound traffic on port 80 from any source, we'll configure listeners to route traffic to the target group on port 80.

The screenshot shows the 'Listeners and routing' configuration for a listener named 'HTTP:80'. The 'Protocol' is set to 'HTTP' and the 'Port' is '80'. The 'Default action' is 'Forward to' a target group named 'web-health-targetgroup'. Below this, there's a 'Create target group' button. The section 'Listener tags - optional' includes a 'Add listener tag' button and a note about adding up to 50 more tags. At the bottom, there's a 'Add listener' button.

Now let's create a load balancer.

The screenshot shows the 'Load balancers' list page with one entry named 'web-health-elb'. The detailed view for this load balancer shows the following information:

| Details | | | |
|--|--|--|---|
| Load balancer type: Application | Status: Active | VPC: vpc-05c70eca7f453ff34 | IP address type: IPv4 |
| Scheme: Internet-facing | Hosted zone: Z555ADOTRQ7XK | Availability Zones: subnets-00a71ae94062e0c459, subnets-0fe51b4061e0b272, subnets-0918de0116ca5a16d, subnets-0c25aef08089e423, subnets-00ba9131930106ff, subnets-05001034e7a416295 | Date created: April 18, 2024, 17:58 (UTC+05:30) |
| Load balancer ARN: arn:aws:elasticloadbalancing:us-east-1:123456789012345678904:loadbalancer/app/web-health-elb/51ee1fb5a5e4fe | DNS name info: web-health-elb-387798721.us-east-1.elb.amazonaws.com (A Record) | | |

Upon viewing the load balancer, you'll notice that it initially enters the provisioning state, transitioning to the active state once fully provisioned. However, despite reaching the active state, attempting to access the website directly from the load balancer may not yield the expected results.

When encountering this issue, accessing the website via the load balancer's DNS name might result in a 504 Gateway Timeout error, indicating a failure to establish a connection between the load balancer and the instances. Despite verifying that the instances are healthy, the presence of this error suggests a potential obstacle.

preventing traffic from passing between the load balancer and the instances.

To address this challenge, we need to identify and resolve the underlying cause obstructing traffic flow between the load balancer and the instances. This discrepancy could be attributed to misconfigured security groups, routing issues, or other networkrelated factors.

By diagnosing and rectifying this barrier, we can ensure seamless connectivity between the load balancer and the instances, thereby resolving the 504 Gateway Timeout error and enabling users to access the website reliably through the load balancer.

Each instance is shielded by a security group, a vital layer of defense governing inbound and outbound traffic. Upon inspection, we observed that both instances share the same security group. However, upon reviewing its inbound rules, we discovered that port 80 was only permitted from our IP address, rather than from the load balancer.

To rectify this, we need to adjust the security group settings to allow traffic from the load balancer. But before proceeding, it's prudent to verify the status of the target group. The target group conducts health checks to determine an instance's well-being. If an instance is deemed unhealthy, the target group will cease routing traffic to it. In this scenario, instances were marked as unhealthy due to the inability to access port 80.

Returning to the security group settings, we locate the security group associated with the instances and modify its inbound rules. Specifically, we add a new rule permitting port 80 access from the security group of the load balancer. By specifying the source as the load balancer's security group, we ensure that any member of the load balancer's group can communicate with the instances over port 80.

Additionally, it's beneficial to include a descriptive label for this rule, such as "allow port 80 from ELB," to enhance clarity and

documentation. Descriptive labels streamline management and aid in understanding the purpose of each rule, promoting efficient administration and troubleshooting.

Make the security groups inbound rules as below of “web-health”.

Inbound rules

| Security group rule ID | Type | Protocol | Port range | Source | Description - optional |
|------------------------|------|----------|------------|--------|------------------------|
| sg-0232f1630eb9d54x9 | HTTP | Custom | 80 | Custom | allow port 80 from ELB |
| sg-0ef1bda5008a8719c | HTTP | Custom | 80 | Custom | 103.203.72.75/32 |
| sg-0fe5a2140927a782d | SSH | Custom | 22 | Custom | 0.0.0.0/0 |

Add rule

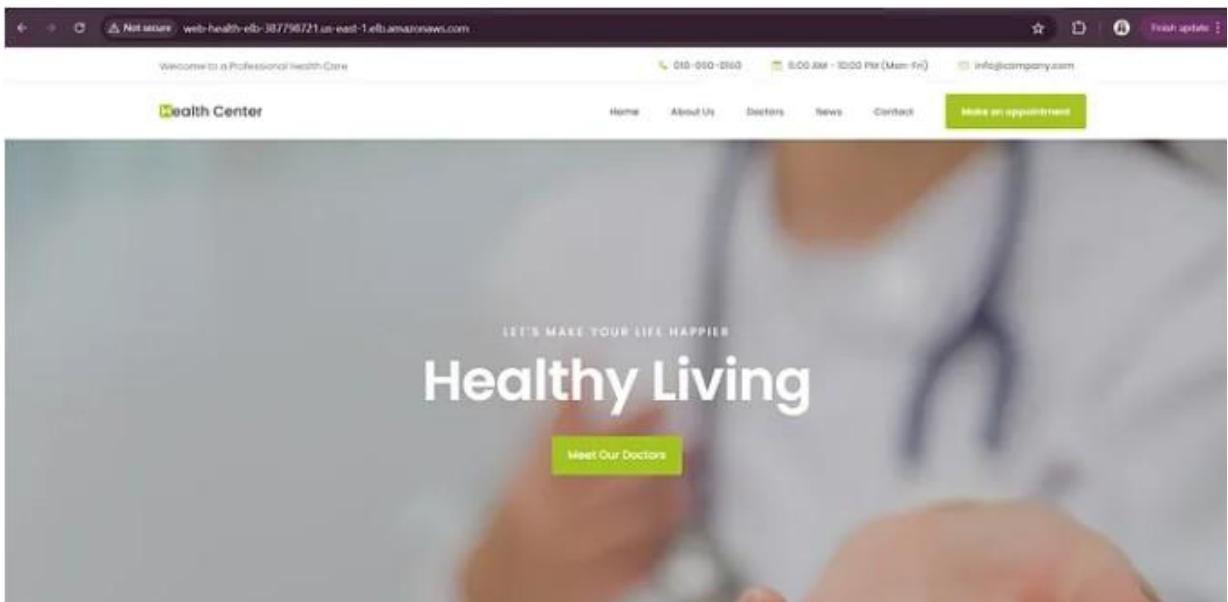
Save the rule.

Make sure that the targets are healthy.

Target group: web-health-targetgroup

| Instance ID | Name | Port | Zone | Health status | Launch time | Anomaly detection |
|---------------------|--------------|------|------------|---------------|-----------------------------------|-------------------|
| i-024082254ff906305 | web-health02 | 80 | us-east-1b | Healthy | April 18, 2024, 17:33 (UTC+05:30) | Normal |
| i-04a19903d8827a6fa | web-health | 80 | us-east-1b | Healthy | April 18, 2024, 17:05 (UTC+05:30) | Normal |

Now check it on the browser:



Managing instances within target groups is a critical aspect of optimizing the performance and availability of applications deployed on AWS. Target groups serve as collections of instances that are subject to health checks, ensuring that traffic is routed only to healthy instances. Here are some key points to consider when managing instances within target groups:

1. **Health Monitoring:** Target groups continuously monitor the health of instances by conducting health checks based on configured parameters. These health checks assess the instance's ability to respond to requests and determine its overall health status.
2. **Registration and Deregistration:** Instances can be registered or deregistered from target groups dynamically based on operational requirements. This flexibility allows for seamless scaling and maintenance operations without disrupting application availability.
3. **Maintenance and Scaling:** Registering or deregistering instances from target groups is a common practice during maintenance activities or when scaling the application. By removing unhealthy or underperforming instances from the target group, traffic can be effectively redirected to healthy instances, ensuring consistent application performance.
4. **Load Balancing:** Target groups play a crucial role in load balancing by distributing incoming traffic across multiple instances based on predefined rules and health check results. This ensures optimal resource utilization and enhances application reliability and scalability.

5. Automation: Managing instances within target groups can be automated using AWS services such as Auto Scaling, which automatically adjusts the number of instances in response to changing demand. This automation helps optimize resource usage and ensures that the application can handle fluctuations in traffic effectively.

Finally, we'll familiarize ourselves with managing instances within the target group, including registering and deregistering instances as needed for maintenance or scaling purposes.

By grasping these concepts and experimenting with the setup, we can gain a comprehensive understanding of deploying and managing web applications on AWS.

AWS CloudWatch

Introduction to CloudWatch:

- AWS CloudWatch is a comprehensive monitoring service offered by Amazon Web Services.
- Originally focused on monitoring, CloudWatch has expanded its capabilities to include logging, events, and more.
- It serves as a central hub for monitoring the performance and health of AWS environments.

Key Features of CloudWatch:

■ Monitoring Service:

- CloudWatch serves as a monitoring solution for AWS resources, tracking performance metrics such as CPU utilization, disk I/O, and network traffic.
- It automatically generates standard metrics for various AWS services used in a region.

■ Logging Solution:

- In addition to monitoring, CloudWatch functions as a logging solution, allowing users to collect, store, and analyze log data generated by AWS resources and applications.
- Logs from services like EC2 instances can be streamed to CloudWatch for centralized log management.

■ Event Monitoring:

- CloudWatch captures real-time events within the AWS environment, such as instance launches, terminations, or volume creations.
- Users can set triggers and notifications based on these events, often integrated with AWS Lambda functions.
- Standard and Custom Metrics:
 - CloudWatch provides both standard and custom metrics for monitoring AWS resources.
 - Standard metrics cover common performance indicators like CPU utilization, network traffic, and disk operations.
- Users can define custom metrics tailored to their specific monitoring needs.
- Alarms and Notifications:
 - Users can set alarms on CloudWatch metrics to trigger notifications when predefined thresholds are breached.
 - Notifications can be sent via email or integrated with Amazon SNS for broader alerting capabilities.
- Integration with AWS Services:
 - CloudWatch seamlessly integrates with various AWS services, including EC2 instances and EBS volumes.
 - Metrics and logs from these services are collected and monitored by CloudWatch, providing insights into resource performance.

Practical Use Cases:

- CloudWatch simplifies monitoring by automatically collecting metrics for AWS resources.
- Users can customize monitoring settings and set up alarms to receive timely notifications of any performance anomalies.
- Practical examples include setting alarms for CPU utilization exceeding a certain threshold, which triggers email notifications via SNS.

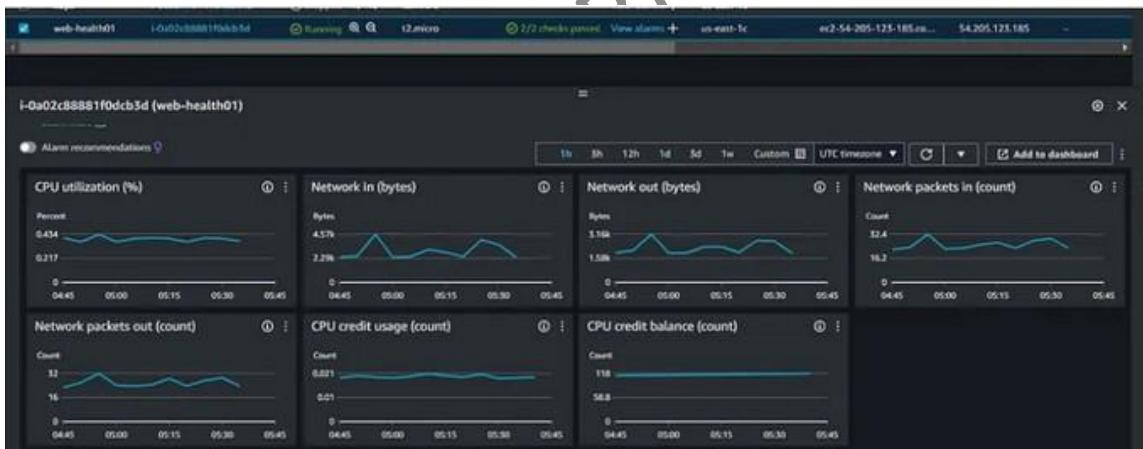
If you've already reviewed the previous blog post, we've provided a template for launching EC2 instances. Kindly review it, or proceed to create an EC2 instance. Select the EC2 instance and click on monitoring:

Before proceeding further, take note of the various metric names available, such as CPU utilization, which indicates the percentage of CPU being used over time. Additionally, observe other metrics like status checks, network in and out (in bytes), network packets in and out (count), and disk read operations.

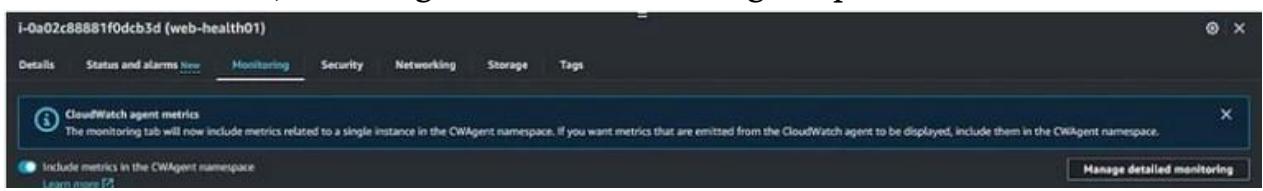
These metrics are automatically generated by CloudWatch when you launch an EC2 instance.

If you require additional metrics, such as RAM or disk utilization, you'll need to create custom metrics. However, for the purpose of this session, we'll focus solely on CPU utilization, as it's one of the most critical metrics to monitor.

By default, CloudWatch checks and updates these metrics every 5 minutes, populating the corresponding graphs. If you prefer more frequent updates, you can enable detailed monitoring, though it's important to note that this option isn't free.

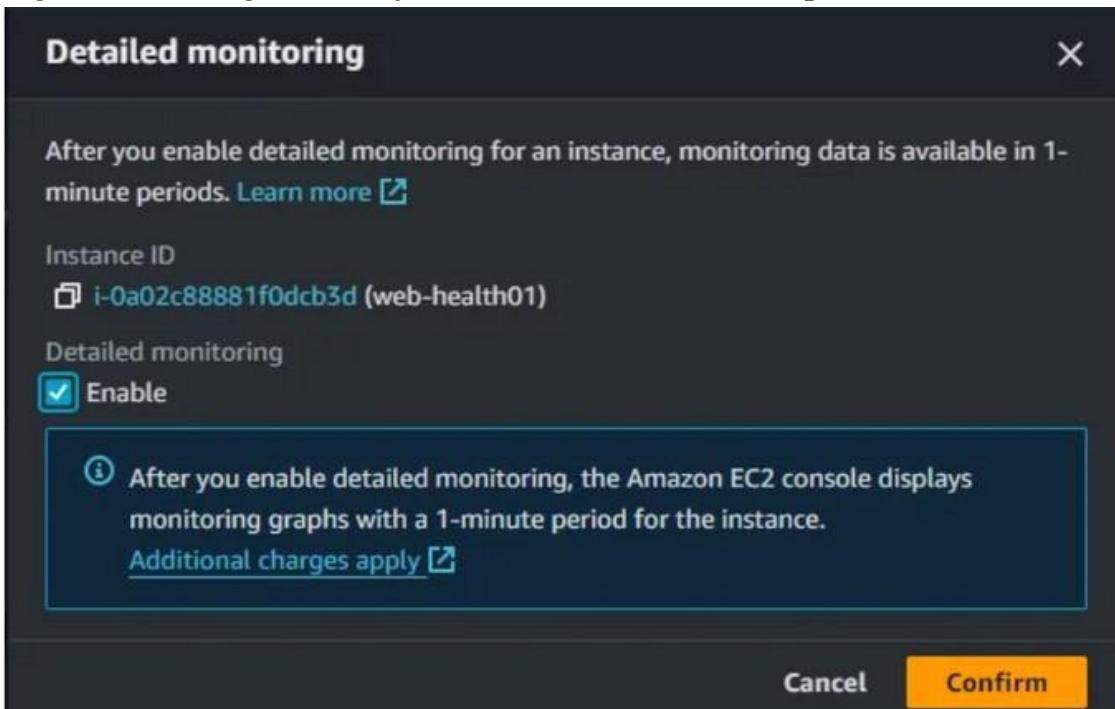


For this hands-on, enabling detailed monitoring is optional.



1. Select the option to “Enable” detailed monitoring, indicating that CloudWatch should monitor metrics every minute.

- 2.Keep in mind that detailed monitoring incurs additional charges compared to the default 5-minute monitoring interval.
- 3.Despite the additional cost, enabling detailed monitoring offers more granular insights into your AWS environment's performance.



Log in to the terminal of the EC2 instance and switch as root user.

Then install the stress tool along with its dependencies on your instance.

Stress is a versatile tool designed to test the limits of your Linux operating system, particularly focusing on CPU performance and other system metrics. Simply running the stress command allows you to stress the CPU. You can specify parameters such as the number of CPUs to stress and even utilize functions like the square root function. Additionally, stress can be used to stress other components like IO and RAM, but for our purpose, we're concentrating on CPU stress testing.

```
yum install  
stress -y
```

```
[root@ip-172-31-29-118 ~]# nohup stress -c 3 -t 100 &  
[2] 66003  
[root@ip-172-31-29-118 ~]# nohup: ignoring input and appending output to 'nohup.out'
```

```
nohup  
stress -c 3 -  
t 100 &
```

This command executes the stress utility in the background using the nohup command, which allows the process to continue running even after the terminal session is terminated. Let's break down the command:

- nohup: Prevents the following command from being terminated when the terminal session ends.
- stress: The stress testing utility.
- -c 3: Specifies that 3 worker threads will be used to stress the CPU cores.
- -t 100: Specifies that the stress test will run for 100 seconds. &: Runs the command in the background.

So, this command will run the stress test with 3 worker threads for 100 seconds, and the process will continue running even if the terminal session is closed. Execute the top command.

top

You'll notice four stress processes running, indicating the CPU utilization at 100 or close to it, reflected by the load average increment.

Continuously stressing the instance like this will be monitored by CloudWatch, updating the graph every minute.

Repeat this process a few times, running for intervals like 100 seconds and 200 seconds. After a few minutes, you'll observe a pattern forming on the graph.

```
[root@ip-172-31-29-118 ~]# top - 06:50:03 up 16:21, 1 user,  load average: 1.18, 0.30, 0.10
Tasks: 107 total,   4 running, 103 sleeping,   0 stopped,   0 zombie
%Cpu(s): 99.7 us,  0.0 sy,  0.0 ni,  0.0 id,  0.0 wa,  0.0 hi,  0.0 si,  0.3 st
MiB Mem :  949.6 total,   538.1 free,   158.2 used,   253.3 buff/cache
MiB Swap:      0.0 total,      0.0 free,      0.0 used.   646.4 avail Mem

PID USER      PR  NI    VIRT    RES    SHR S %CPU %MEM TIME+ COMMAND
65905 root      20   0  3512   116     0 R 33.2  0.0  0:11.25 stress
65903 root      20   0  3512   116     0 R 32.9  0.0  0:11.24 stress
65904 root      20   0  3512   116     0 R 32.9  0.0  0:11.24 stress
  1 root      20   0 170708  16540  10084 S  0.0  1.7  0:07.29 systemd
  2 root      20   0      0     0     0 S  0.0  0.0  0:00.01 kthreadd
  3 root      0 -20      0     0     0 I  0.0  0.0  0:00.00 rcu_gp
  4 root      0 -20      0     0     0 I  0.0  0.0  0:00.00 rcu_par_gp
  5 root      0 -20      0     0     0 I  0.0  0.0  0:00.00 slab flushwq
```

Let's create a script for this process.

```
vim stress.sh
```

```
sleep 60 && stress -c 4 -t 60 && sleep 60 && stress -c 4 -t 60 &&
sleep 60 && stress -c 4 -t 30 && sleep 60 && stress -c 4 -t 100 &&
sleep 30 && stress -c 4 -t 200 :wq!
```

This sequence of commands progressively varies the duration of the stress tests, interspersed with periods of rest. It's a way to simulate different levels of CPU load on the system over time. designed to stress test the system's CPU using the stress utility. Let's break down each part:

- **sleep 60:** This command pauses execution for 60 seconds before proceeding to the next command. It introduces a delay in the sequence.
- **stress -c 4 -t 60:** This command runs the stress utility with the following options:
 - **-c 4:** Specifies that 4 worker threads will be used to stress the CPU cores.
 - **-t 60:** Specifies that the stress test will run for 60 seconds.
- **sleep 60:** Another pause of 60 seconds.
- **stress -c 4 -t 60:** Similar to the second command, this runs a stress test for another 60 seconds.
- **sleep 60:** Another pause of 60 seconds.
- **stress -c 4 -t 30:** This command runs a shorter stress test for 30 seconds.
- **sleep 60:** Another pause of 60 seconds.

- **stress -c 4 -t 100:** This command runs a longer stress test for 100 seconds.
- **sleep 30:** A shorter pause of 30 seconds.
- **stress -c 4 -t 200:** This command runs an even longer stress test for 200 seconds.

Execute this script:

```
./stress.sh
```

```
[root@ip-172-31-29-118 ~]# nohup ./stress.sh &
[1] 66549
[root@ip-172-31-29-118 ~]# nohup: ignoring input and appending output to 'nohup.out'
```

We generate a graph by running the command. I simply executed it using nohup stress.sh &.

Ensure your script

is executable.

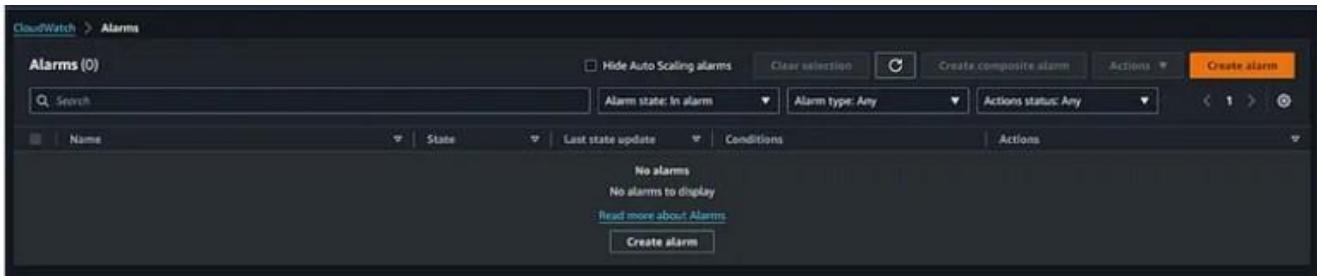
```
chmod +x stress.sh
```

Now, when observing the top output, you may intermittently notice the stress command based on the ongoing operations within your script.



Now, let's navigate to the CloudWatch service to configure an alarm for CPU utilization on this instance.

First, head over to the “All alarms” section. Next, click on the “Create Alarm” button.

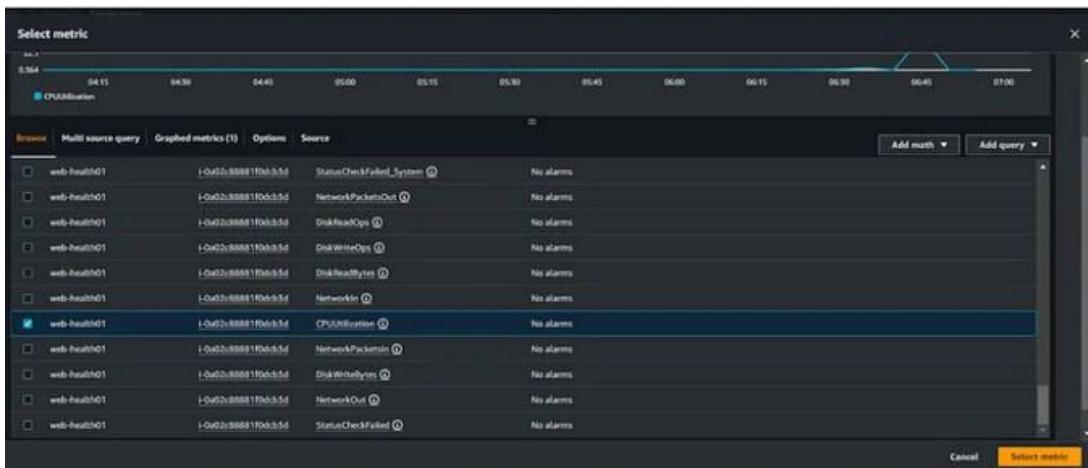


To select the CPU utilization metric, start by navigating to the EC2 service.

Then, locate the “Per-instance metrics” section and find your instance. If you don’t immediately see your instance, wait for a few moments for the information to load.

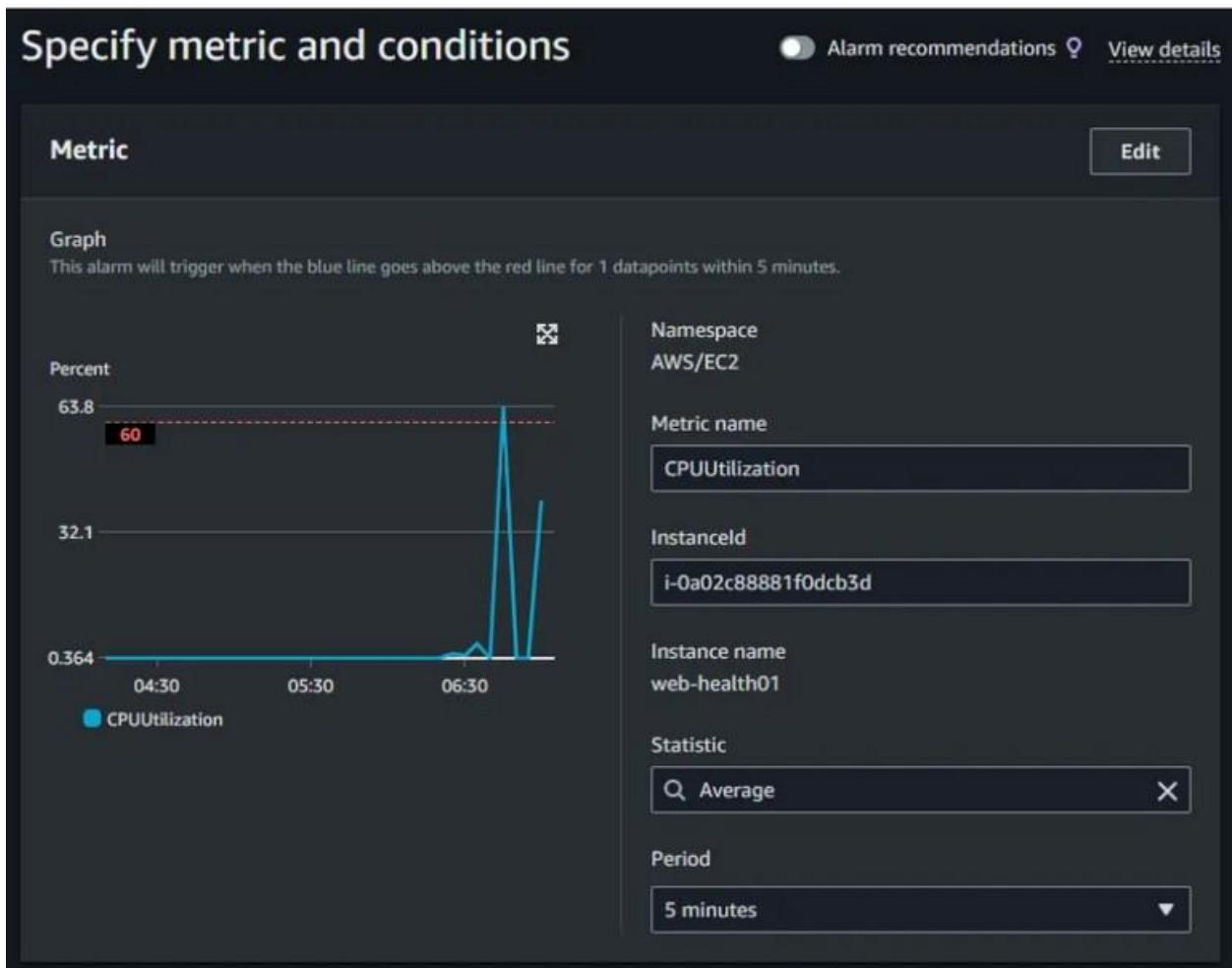
Once you’ve found your instance, look for the CPU utilization metric.

Select the CPU utilization metric by clicking on it.



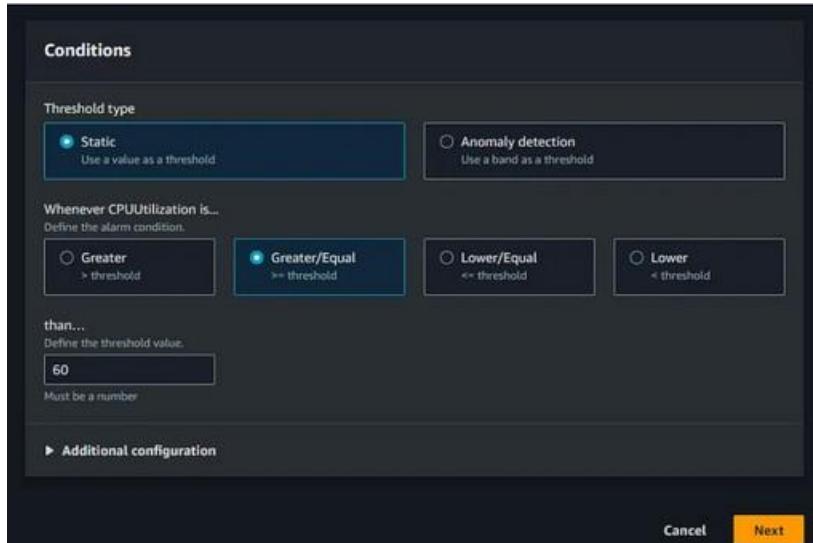
In this step, we can choose the period for which the alarm will be evaluated. For this demonstration, I’ll keep the period set to 5 minutes.

Next, we specify the condition for triggering the alarm based on CPU utilization.



For example, we can set the condition to trigger the alarm if the CPU utilization is greater than or equal to 60 for a period of 5 minutes.

Once we've defined the condition, we proceed to the next step.



If you don't find your desired topic listed, you can click on "Create Topic." Provide a name for the topic, enter your email address, and click on "Create Topic."

However, since we already created a topic for billing alarms earlier, I'll select the same topic, which contains my email address.

This topic is designated for notifications. So, when the alarm state is triggered, it will send a notification to this topic, ultimately resulting in an email notification being sent to me.

Configure actions

Notification

Alarm state trigger
Define the alarm state that will trigger this action.

In alarm
The metric or expression is outside of the defined threshold.

OK
The metric or expression is within the defined threshold.

Insufficient data
The alarm has just started or not enough data is available.

Send a notification to the following SNS topic
Define the SNS (Simple Notification Service) topic that will receive the notification.

Select an existing SNS topic

Create new topic

Use topic ARN to notify other accounts

Send a notification to...

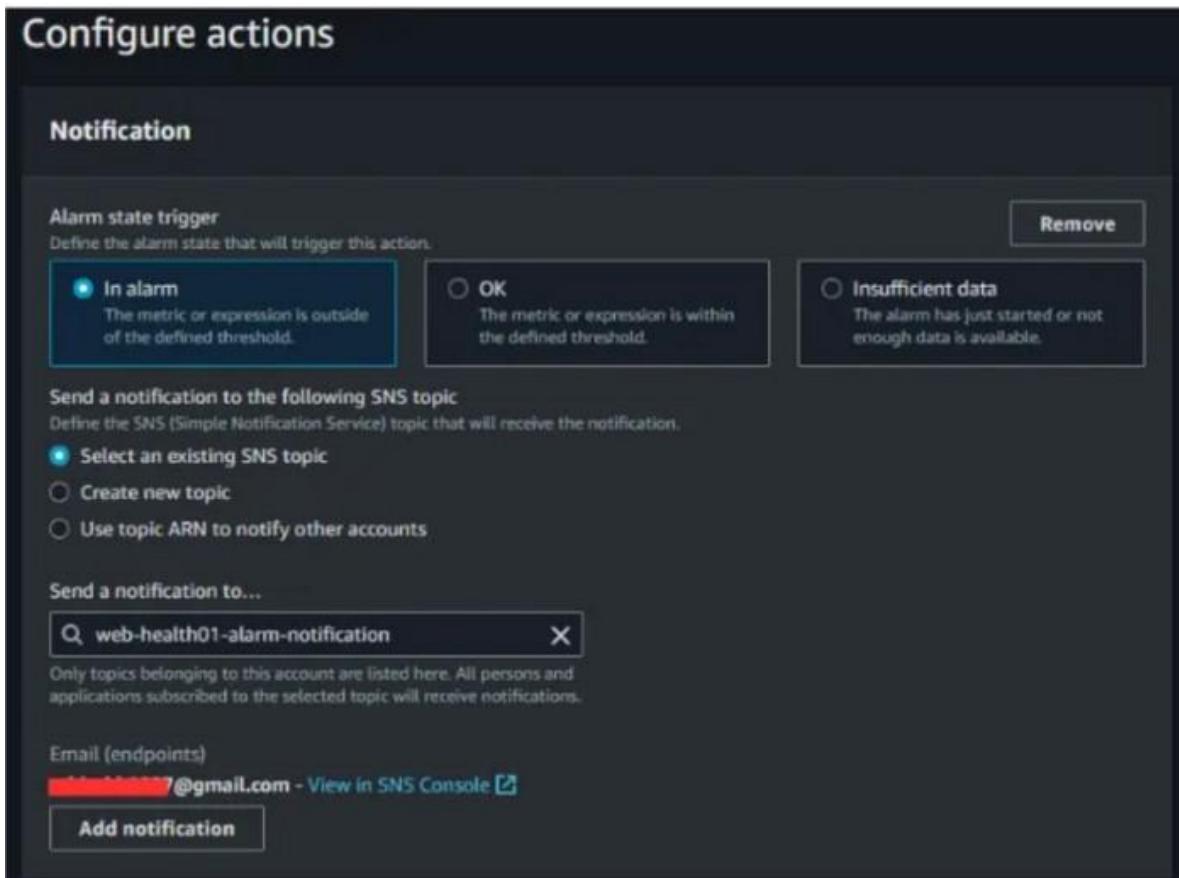
X

Only topics belonging to this account are listed here. All persons and applications subscribed to the selected topic will receive notifications.

Email (endpoints)

[\[REDACTED\]@gmail.com - View in SNS Console](#)

[Add notification](#)



There are several other actions you can take, such as EC2 actions. For instance, you might want to stop, terminate, or reboot the instance if it triggers an alarm.

In some cases, a high CPU utilization might prevent you from logging into the instance via SSH. Rebooting the instance could be a temporary solution to address this issue.

However, for the purpose of this demonstration, we will skip these additional actions.

We will stick to email notifications only and proceed to the next step.

The screenshot shows the 'EC2 action' configuration page. At the top, it says 'Alarm state trigger' and 'Define the alarm state that will trigger this action.' There are three options: 'In alarm' (selected), 'OK', and 'Insufficient data'. A 'Remove' button is located in the top right corner of the trigger section. Below this, the section 'Take the following action...' is titled 'Define what will happen to the EC2 instance with the Instance ID i-0a02c88881f0dcb3d when this alarm is triggered.' It lists four actions: 'Recover this instance' (selected), 'Stop this instance', 'Terminate this instance', and 'Reboot this instance'. Each action has a detailed description below it. At the bottom left is a 'Add EC2 action' button.

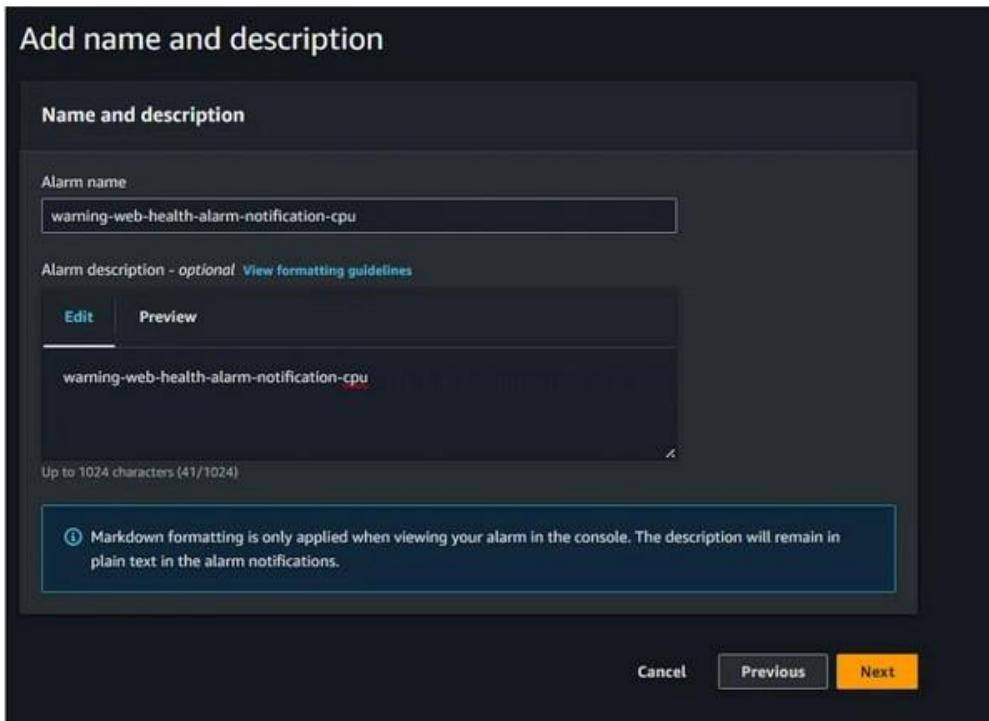
Let's give the alarm a descriptive name.

The naming convention will be: "warning-web-health-alarm-notification-cpu" for the specific instance.

In our organization, a warning is triggered when the utilization exceeds 60%. We could set up additional alarms for critical situations, such as when it surpasses 80%.

However, for now, we'll proceed with just one alarm.

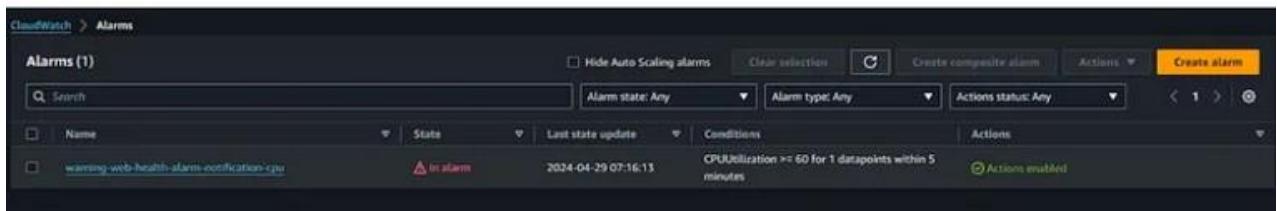
Now, let's move on to the next step.



After a while, CloudWatch will collect the data and display whether the instance is in an alarm state or if it's okay.

Now, this process aligns with the basic principles of any monitoring tool. You have metrics or checks, alarms triggered by those metrics or checks, and actions associated with those alarms, such as sending email notifications.

Whether it's Prometheus, Nagios, Icinga, Zenos, or others, the concept remains similar. However, with CloudWatch, monitoring is already configured, and you only need to set up alarms. In contrast, with other tools, such as those you set up yourself, the entire monitoring system needs configuration, typically by the monitoring or administration team.



Currently, the graph fluctuates because I've executed the stress command multiple times. Remember, you need to run the stress

Multicloud Devops by veera nareshit

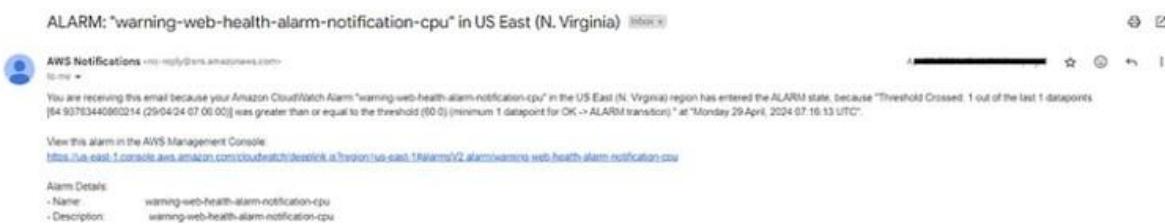
command for at least 5 minutes or longer for it to exceed the threshold and trigger the email notification.

After running the command for 5 minutes, wait for the notification to arrive in your inbox.

Moreover, you can create a reverse alarm to monitor the instance's OK state. For instance, if the CPU utilization is below 40, it's considered OK. You can set up alarms for different thresholds and states, and you have the flexibility to create up to ten alarms under the free tier.



Email notification:



I encourage you to experiment with these alarms to familiarize yourself with their functionalities. Set alarms for various data points, both in alarm and OK states, using different conditions like greater than or less than. Once you're comfortable, remember to terminate the instance and delete the alarms before moving on to the next section, particularly the auto-scaling group section.

Auto Scaling Groups

Auto Scaling Groups (ASG) seamlessly integrates with CloudWatch to monitor specific metrics, such as CPU utilization. If a metric, like CPU utilization, exceeds a predefined threshold, ASG adjusts capacity by either adding or removing instances in the group. This adjustment is based on CloudWatch alarms, which trigger actions to maintain performance or control costs.

When creating an auto scaling group, you'll configure launch configurations or launch templates, similar to those used in EC2 instances for ELB exercises. These configurations provide the necessary information for launching instances within the ASG.

To dynamically adjust capacity based on metrics, ASG employs scaling policies. You can define these policies, such as step scaling, to add or remove instances based on workload fluctuations. For instance, you could specify that if CPU usage exceeds 60%, launch two additional instances, or if it exceeds 80%, launch four instances.

When setting up an auto scaling group, you specify parameters like minimum, desired, and maximum capacity. The minimum size ensures a baseline level of instances that cannot be terminated, while desired capacity represents the typical number of running instances. The maximum size limits the total number of instances the group can scale up to, helping to control costs.

In practice, ASG operates by creating a group powered by launch configurations. When a metric breach triggers a CloudWatch alarm, scaling policies kick in to launch or terminate instances accordingly, ensuring optimal performance and resource utilization.

Now, let's dive into the practical implementation of these concepts.

Before diving into the auto scaling group setup, let's ensure we have all the necessary components in place.

Firstly, we'll need the launch template, which should already be created as per the instructions provided in the load balancer section. If the launch template has been removed or needs to be recreated,

refer back to the load balancer section for detailed steps on how to create it.

Now, let's proceed to the Target Group setup.

- Navigate to the Target Group section.
- Click on the “Create target group” button.

The screenshot shows the 'Specify group details' page for creating a target group. At the top, a note says: 'Your load balancer routes requests to the targets in a target group and performs health checks on the targets.' Below this is a section titled 'Basic configuration' with a note: 'Settings in this section can't be changed after the target group is created.' The main area is titled 'Choose a target type' and contains four options:

- Instances
 - Supports load balancing to instances within a specific VPC.
 - Facilitates the use of [Amazon EC2 Auto Scaling](#) to manage and scale your EC2 capacity.
- IP addresses
 - Supports load balancing to VPC and on-premises resources.
 - Facilitates routing to multiple IP addresses and network interfaces on the same instance.
 - Offers flexibility with microservice based architectures, simplifying inter-application communication.
 - Supports IPv6 targets, enabling end-to-end IPv6 communication, and IPv4-to-IPv6 NAT.
- Lambda function
 - Facilitates routing to a single Lambda function.
 - Accessible to Application Load Balancers only.
- Application Load Balancer
 - Offers the flexibility for a Network Load Balancer to accept and route TCP requests within a specific VPC.
 - Facilitates using static IP addresses and PrivateLink with an Application Load Balancer.

Configure the target group settings, such as name. For example, let's name it “webhealth-tg”.

- Keep the other configurations as default and proceed by clicking “Next”.
- Since we don't have any instances yet, select the option to create an empty target group.
- Complete the target group creation process by following the prompts.

Multicloud Devops by veera nareshit

Target group name
A maximum of 32 alphanumeric characters including hyphens are allowed, but the name must not begin or end with a hyphen.

Protocol : Port
Choose a protocol for your target group that corresponds to the Load Balancer type that will route traffic to it. Some protocols now include anomaly detection for the targets and you can set mitigation options once your target group is created. This choice cannot be changed after creation.

HTTP 80

IP address type
Only targets with the indicated IP address type can be registered to this target group.

IPv4
Each instance has a default network interface (eth0) that is assigned the primary private IPv4 address. The instance's primary private IPv4 address is the one that will be applied to the target.

IPv6
Each instance you register must have an assigned primary IPv6 address. This is configured on the instance's default network interface (eth0). [Learn more](#)

Indeed, the Auto Scaling group will automatically add instances to this target group as needed based on your scaling policies and the metrics you've configured, ensuring seamless scalability and efficient resource management.



| Target groups [1] | | | | | | | Actions | Create target group |
|-------------------|------------------------------|------|----------|-------------|---------------|------------------------|---------|---------------------|
| Name | ARN | Port | Protocol | Target type | Load balancer | VPC ID | | |
| web-health-tg | arn:aws:elasticloadbalanc... | 80 | HTTP | Instance | web-health-eb | vpc-0f6a608f533427a214 | | |

We'll proceed by creating a load balancer for this target group. Let's select "Create load balancer" and opt for the application load balancer, maintaining consistency with our previous setup.

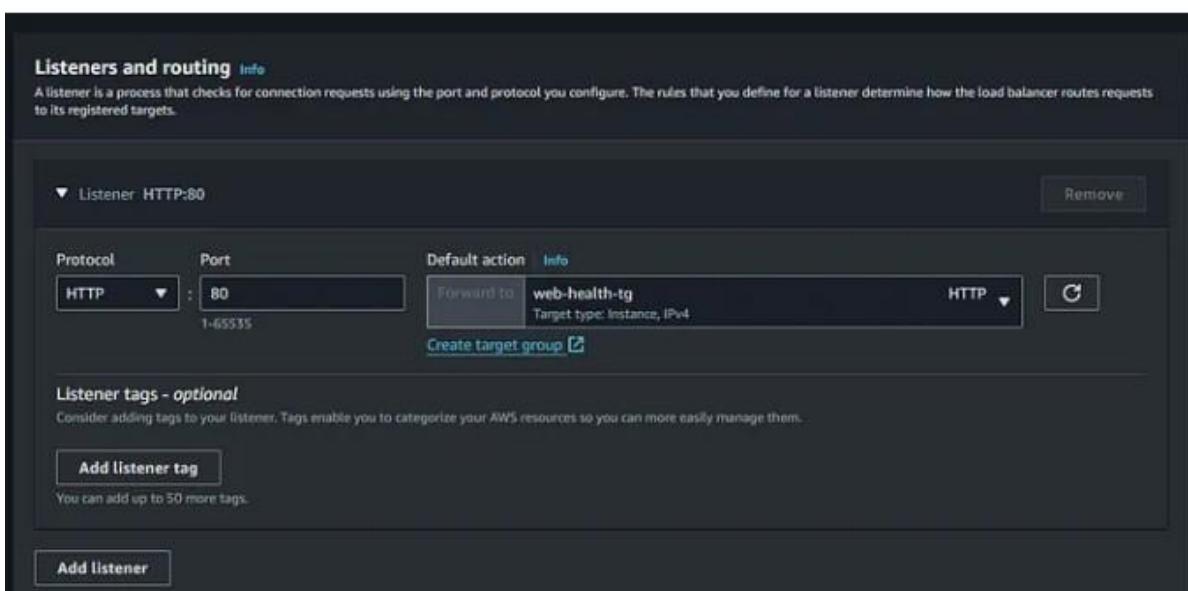
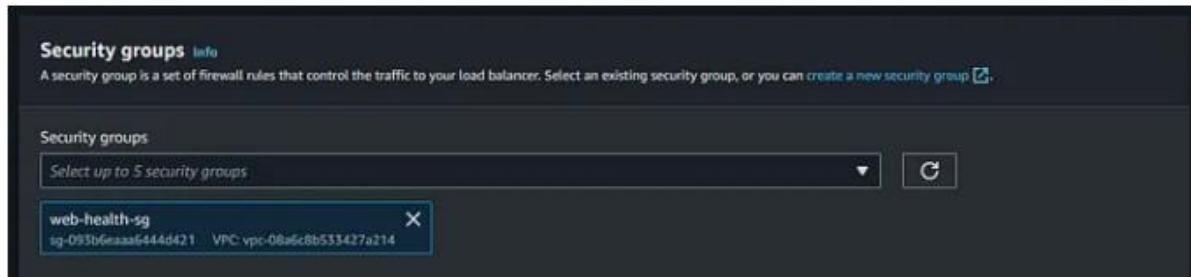
Multicloud Devops by veera nareshit

The screenshot shows the 'Create Application Load Balancer' wizard in the AWS Management Console. The current step is 'Basic configuration'. It includes fields for 'Load balancer name' (set to 'web-health-elb'), 'Scheme' (set to 'Internet-facing'), and 'IP address type' (set to 'IPv4'). A note indicates that the load balancer must be unique within the account and can't be changed after creation.

We'll name it “health-elb” and choose the same availability zones and security groups used previously.

The screenshot shows the 'Network mapping' configuration screen. It lists three subnets under a single VPC selection: 'vpc-08a6cBb533427a214' (IPv4 CIDR: 172.31.0.0/16). The 'Mappings' section shows three availability zones (az1, az3, az2) each associated with one of these subnets, all selected with checkboxes. Each zone entry includes an 'IPv4 address' field and an 'Assigned by AWS' note.

Remember to remove the Default Security Group. Next, we'll configure it to forward traffic to the target group we just created.



If you already have these resources set up, you can reuse them; otherwise, go ahead and create them now. Here's our load balancer.



Let's revisit the launch template. This template is essential as it defines the configuration for the instances that will be launched by the auto scaling group.

It's important to verify that the required AMI is available and accessible. If you've already prepared the AMI, it should be readily available for selection within the launch template settings.

| Amazon Machine Images (AMIs) (1) | | | | | | |
|----------------------------------|-----------------------|-----------------------------|--------------|---------|------------|--------|
| Name | AMI name | AMI ID | Source | Owner | Visibility | Status |
| web-health-ami | ami-0b1d84427ca798310 | 323390159304/web-health-ami | 323390159304 | Private | Available | |

The launch template specifies details such as the Amazon Machine Image (AMI) to be used for the instances.

| Launch Templates (1) | | | | | | |
|----------------------|----------------------|-----------------|----------------|--------------------------|--------------------------------|---------|
| Launch Template ID | Launch Template Name | Default Version | Latest Version | Create Time | Created By | Actions |
| ii-0334384058ff110 | web-health-template | 1 | 1 | 2024-05-06T20:03:14.000Z | arn:aws:iam::323390159304:root | |

Now, let's navigate to the auto scaling groups section.

Let's initiate the creation process by assigning a name to our auto scaling group. For simplicity, let's name it "web-health-asg".

Choose launch template or configuration Info

Specify a launch template that contains settings common to all EC2 instances that are launched by this Auto Scaling group. If you currently use launch configurations, you might consider migrating to launch templates.

Name

Auto Scaling group name
Enter a name to identify the group.

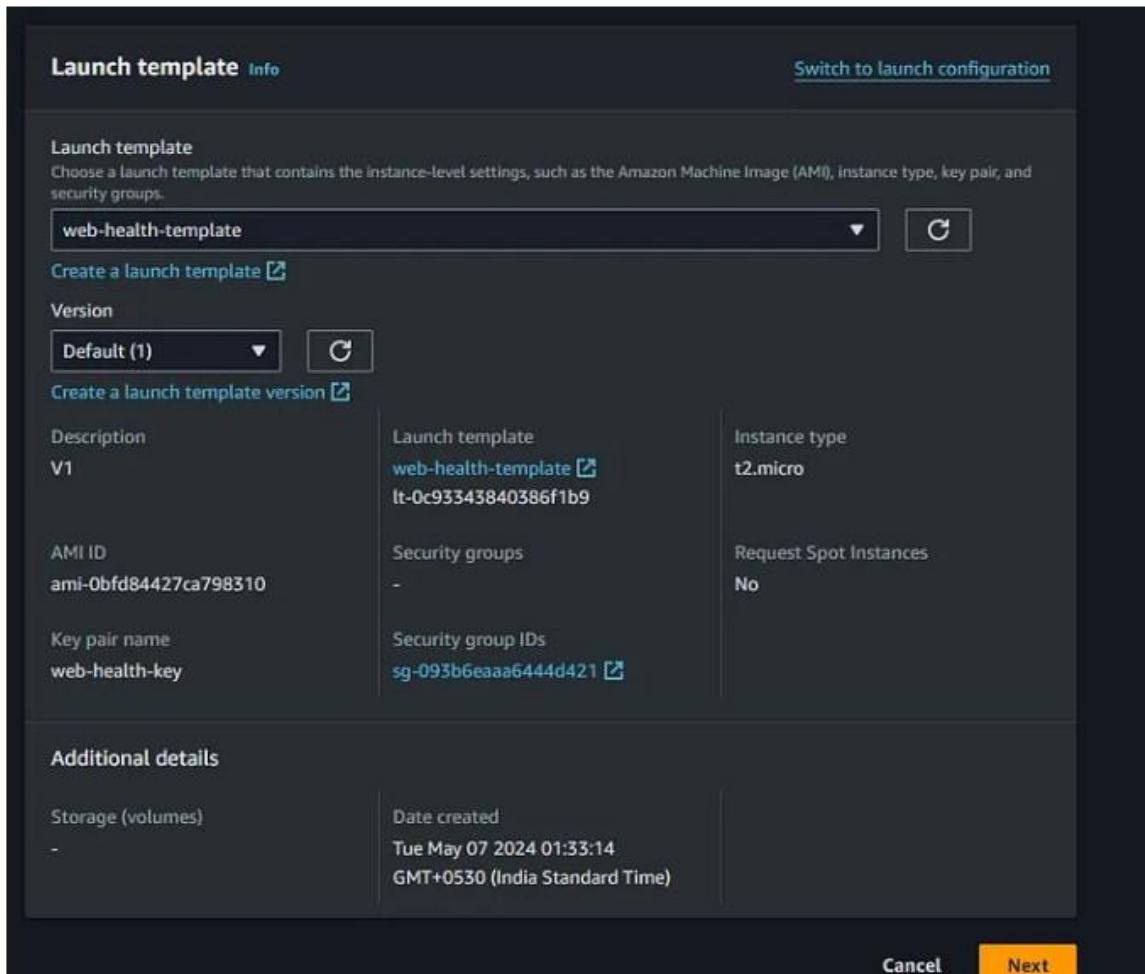
Must be unique to this account in the current Region and no more than 255 characters.

Next, we'll select our launch template for the auto scaling group.

You might also notice the option to use a launch configuration, which serves a similar purpose to a launch template. However, AWS now recommends using launch templates due to their added flexibility.

One notable advantage of launch templates is their editability. Unlike launch configurations, which cannot be modified once created, launch templates allow for easy edits and the creation of multiple versions. This versatility enables you to tailor your configurations to different environments or instance types seamlessly.

As we proceed, we'll opt for the launch template, leveraging its benefits for our auto scaling group. You'll notice the option to select different versions of the launch template, providing flexibility in configuration management.

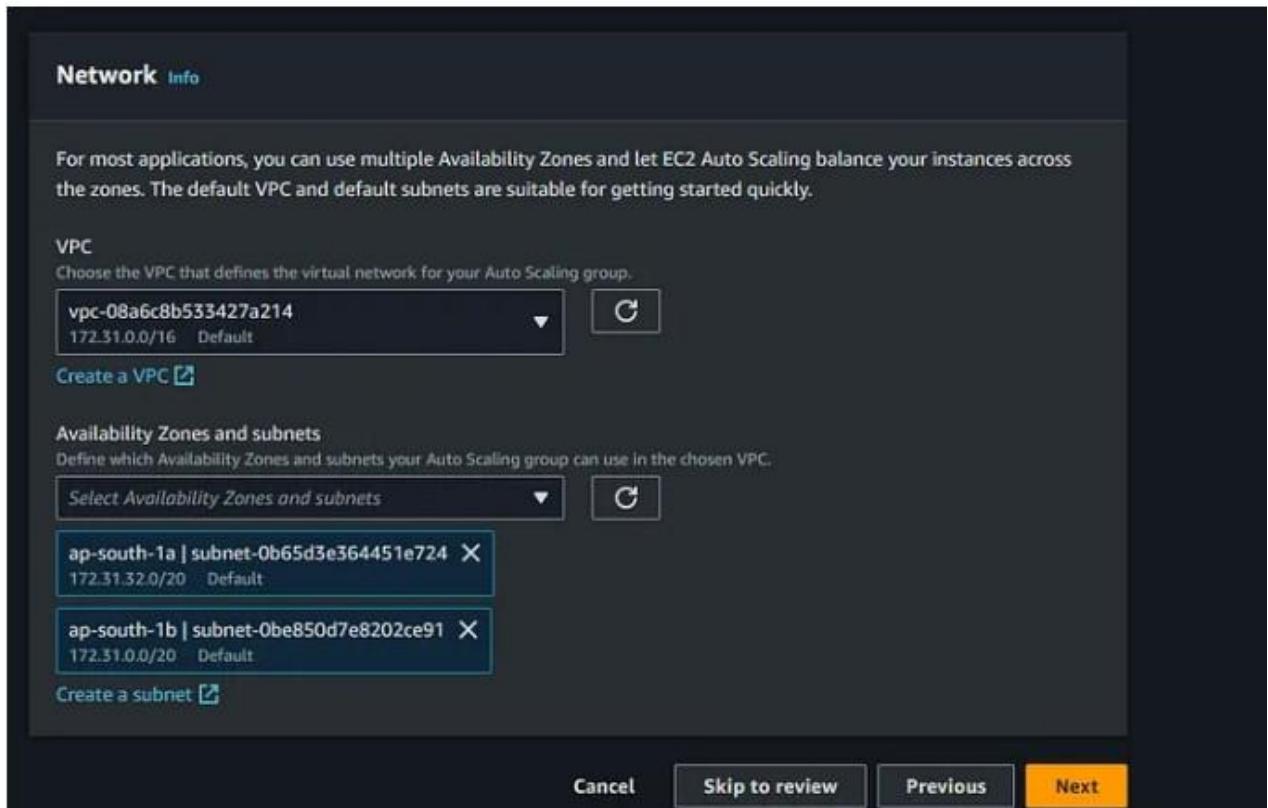


Choose instance launch options:

In this step, we need to specify the availability zones where we want our instances to be launched.

You have the option to select specific zones or allow AWS to choose based on availability. It's generally recommended to choose at least two zones for redundancy and fault tolerance.

Let's proceed by selecting all zones and moving to the next step.



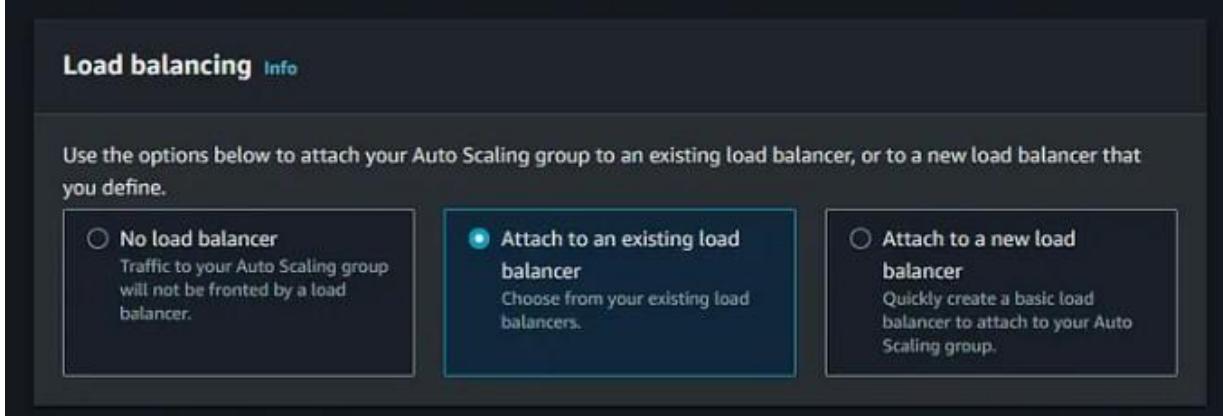
Configure advanced options

In this step, we have the option to attach our auto scaling group to an existing load balancer.

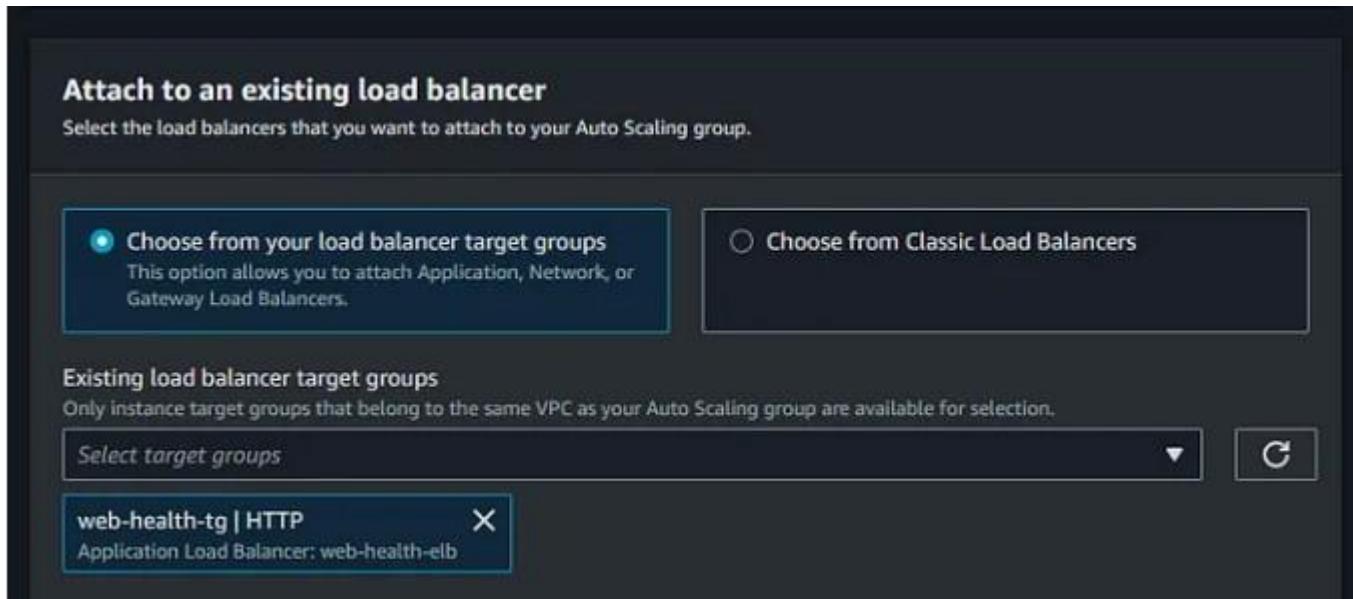
Since our website will benefit from a load balancer, especially with multiple instances, we'll choose to attach it to an existing load balancer.

Configure advanced options - optional Info

Choose a load balancer to distribute incoming traffic for your application across instances to make it more reliable and easily scalable. You can also set options that give you more control over health check replacements and monitoring.



We'll then select the target group that we previously created. This ensures that the instances launched by our auto scaling group are registered with the specified target group, allowing the load balancer to distribute traffic effectively across these instances.



In this step, we configure the health check settings for our instances within the auto scaling group.

The primary responsibility of the auto scaling group is to maintain the desired number of healthy instances, ensuring optimal performance and reliability.

By default, the auto scaling group performs basic EC2 health checks, which include hardware and virtual machine checks. However, these checks only verify the instance's basic functionality and do not guarantee the health of your application or website.

To enable more comprehensive health checks, we opt for ELB (Elastic Load Balancer) health checks. This leverages the target group health check functionality discussed in the load balancer section.

The target group conducts health checks every 300 seconds, verifying the health of instances based on specified criteria such as port availability and process responsiveness. If an instance fails this

health check, indicating a potential issue with the hosted application or service, the target group declares it as unhealthy.

Upon detecting an unhealthy instance, the auto scaling group takes corrective action by terminating the problematic instance and launching a replacement. This ensures that the overall system remains resilient and capable of handling fluctuations in workload and application health effectively.

The screenshot shows the 'Health checks' section of the AWS Auto Scaling configuration. It includes:

- EC2 health checks:** Status is **Always enabled**.
- Additional health check types - optional:** An option to **Turn on Elastic Load Balancing health checks** is selected, marked as **Recommended**. A note states: "Elastic Load Balancing monitors whether instances are available to handle requests. When it reports an unhealthy instance, EC2 Auto Scaling can replace it on its next periodic check."
- A callout box provides information: "EC2 Auto Scaling will start to detect and act on health checks performed by Elastic Load Balancing. To avoid unexpected terminations, first verify the settings of these health checks in the [Load Balancer console](#)." There is a close button (X) in the top right corner of the box.
- Health check grace period:** Set to **300 seconds**.

Configure group size and scaling

Now, let's proceed to the next step.

Here, we specify the capacity settings for our auto scaling group. We define the desired, minimum, and maximum number of instances that the group should maintain.

- **Desired capacity:** This indicates the number of instances we want to have in our auto scaling group. For example, setting it to 2 means we aim to have two instances running at all times.
- **Minimum capacity:** This represents the minimum number of instances that must be running in the auto scaling group,

regardless of the workload or scaling policies. In our case, we set it to 1 to ensure there is always at least one instance available.

- **Maximum capacity:** This defines the upper limit on the number of instances that can be launched by the auto scaling group. Even if the workload demands more instances, the group will not exceed this limit. Here, we set it to 8 to cap the maximum number of instances at eight.

Configure group size and scaling - optional Info

Define your group's desired capacity and scaling limits. You can optionally add automatic scaling to adjust the size of your group.

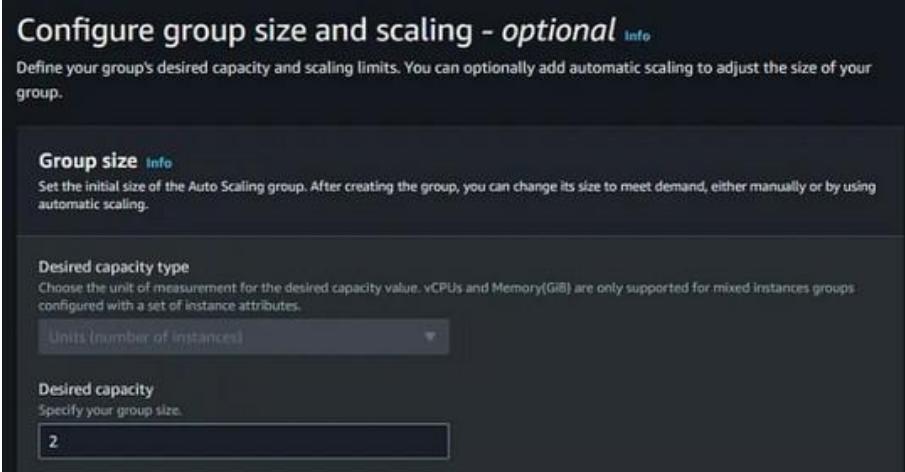
Group size Info

Set the initial size of the Auto Scaling group. After creating the group, you can change its size to meet demand, either manually or by using automatic scaling.

Desired capacity type
Choose the unit of measurement for the desired capacity value. vCPUs and Memory(GiB) are only supported for mixed instances groups configured with a set of instance attributes.

Units (number of instances)

Desired capacity
Specify your group size.



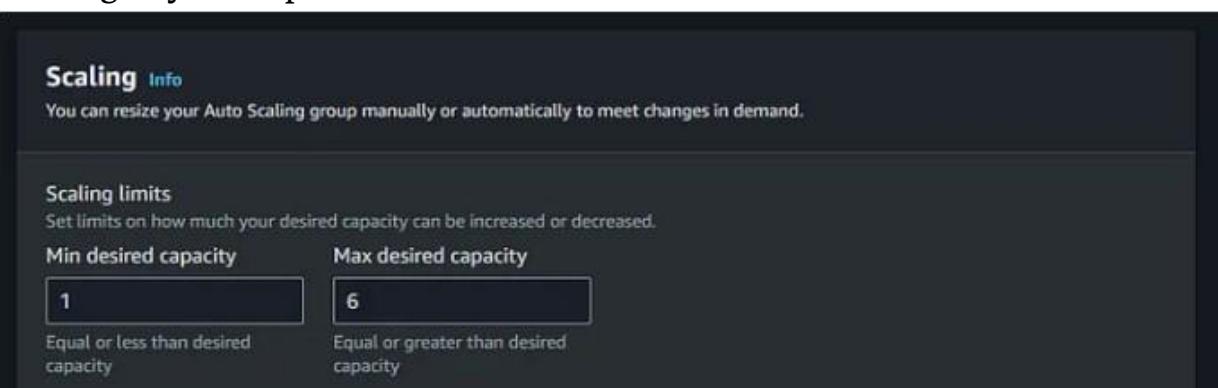
Understanding the distinction between desired, minimum, and maximum capacities is crucial. While desired capacity determines the ideal number of instances, minimum capacity ensures a baseline level of availability, and maximum capacity prevents excessive scaling beyond a predetermined threshold.

Scaling Info

You can resize your Auto Scaling group manually or automatically to meet changes in demand.

Scaling limits
Set limits on how much your desired capacity can be increased or decreased.

| | |
|---|---|
| Min desired capacity <input type="text" value="1"/> | Max desired capacity <input type="text" value="6"/> |
| Equal or less than desired capacity | Equal or greater than desired capacity |



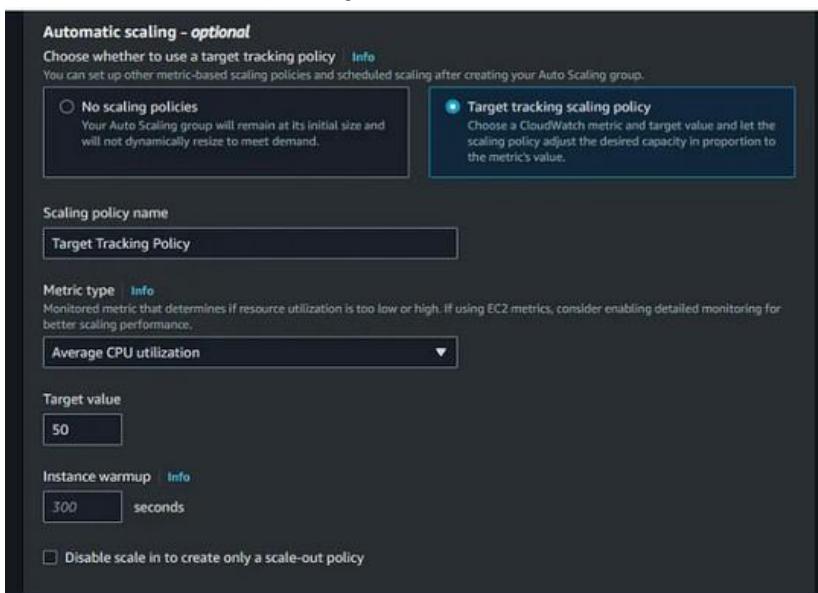
Next, we delve into scaling policies, which dictate how the auto scaling group responds to changes in workload or system metrics.

- **Scaling policies:** These policies govern when and how the auto scaling group should scale out (add instances) or scale in (remove instances) based on predefined conditions. We can choose from

various scaling policy types, such as target tracking scaling policies or simple scaling policies.

For instance, we opt for a target tracking scaling policy based on CPU utilization. If the combined CPU utilization across instances exceeds 50%, the auto scaling group will add instances to handle the increased workload. Conversely, if CPU utilization drops below 50%, instances may be removed to optimize resource usage.

Additionally, we have the option to enable scaling protection, which safeguards newly launched instances from immediate termination. This feature can be beneficial in scenarios where instances require time to initialize or synchronize data.



By configuring these capacity and scaling policy settings, we ensure that our auto scaling group operates efficiently, maintaining optimal performance while dynamically adjusting to changing demand.

Instance maintenance policy Info

Control your Auto Scaling group's availability during instance replacement events. This includes health checks, instance refreshes, maximum instance lifetime features and events that happen automatically to keep your group balanced, called rebalancing events.

Choose a replacement behavior depending on your availability requirements

Mixed behavior (Selected)

Prioritize availability

Control costs

Flexible

Instance scale-in protection

Scale-in protection prevents newly launched instances from being terminated by scaling activities. Make sure to remove scale-in protection for the group or individual instances when instances are ready to be terminated.

Enable instance scale-in protection

Cancel Skip to review Previous Next

Add notifications

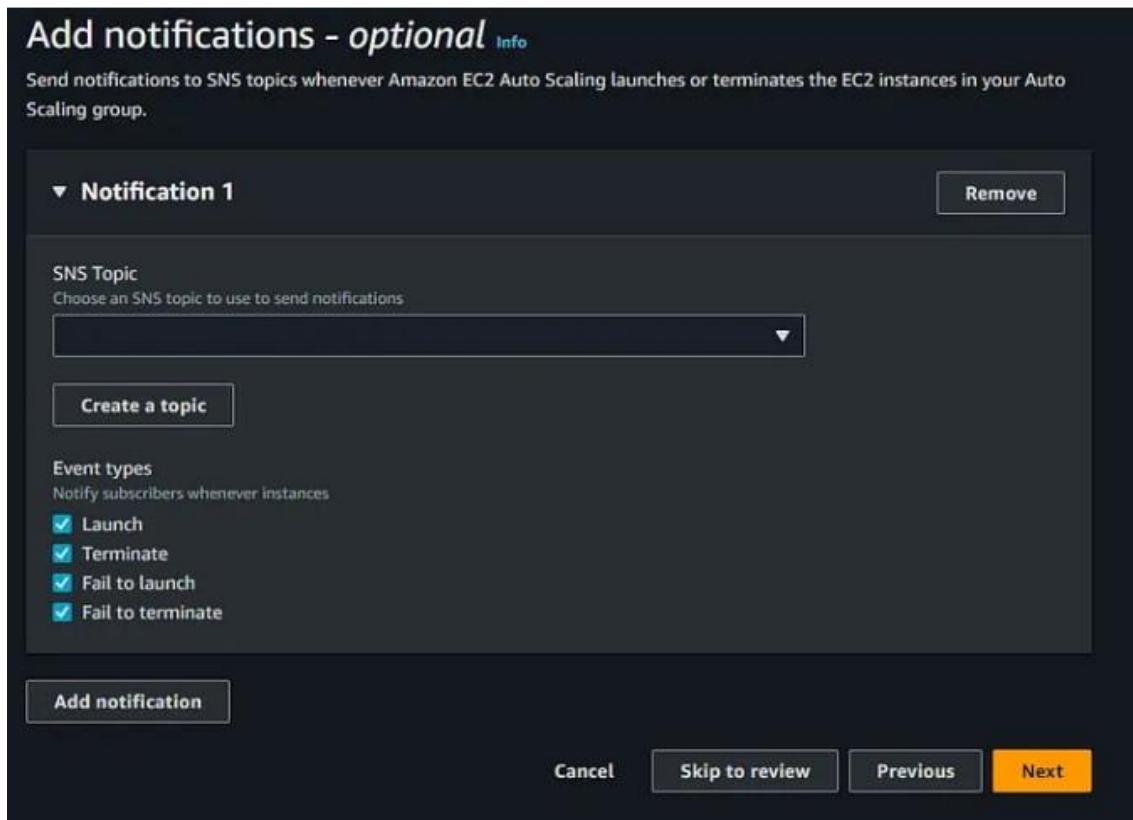
Now, let's proceed to the next step.

Here, we have the option to configure event notifications for various actions performed by the auto scaling group.

- **Notification settings:** We can specify a notification topic to receive notifications for specific events related to the auto scaling group's lifecycle. These events include instance launches, terminations, failures to launch, and failures to terminate.

By defining a notification topic, we ensure that relevant stakeholders or automated systems are promptly informed of any changes or issues affecting the auto scaling group's operation. This facilitates proactive monitoring and troubleshooting, allowing us to maintain the reliability and availability of our infrastructure.

After configuring the notification settings as needed, we can proceed to the next stage of the auto scaling group setup process.



Add tags

Let's add tags for identification purposes.

- Tag Name: “web-health”

We'll keep the naming simple and generic since the number of instances can vary dynamically based on scaling policies.

The screenshot shows the 'Add tags - optional' step of creating an Auto Scaling group. A tooltip explains that you can add tags to instances and EBS volumes by specifying tags in the launch template. The main area shows one tag named 'Name' with the value 'web-health'. A checkbox labeled 'Tag new instances' is checked. At the bottom, there are 'Cancel', 'Previous', and 'Next' buttons.

Now, let's proceed by clicking on "Next" and review all the provided information before finalizing the creation of the auto scaling group. Once verified, we can proceed with the creation process.

Currently, we don't have any instances running. As per our configuration, the desired capacity is set to two instances. So, the auto scaling group will start creating instances to meet this desired capacity.

You can observe that two instances are being created to fulfill the desired capacity. Additionally, there's a scaling policy in place. If the CPU utilization drops below 50%, it will terminate instances accordingly. However, it will always maintain a minimum of one instance.

In the instance management section, you'll find all the instances within the auto scaling group. This section is particularly useful for troubleshooting purposes. Here, you can perform various actions such as detaching an instance from the auto scaling group, setting it to standby mode, or bringing it back into service.

The screenshot shows the AWS Auto Scaling group activity history for 'weh-health-asg'. The 'Activity' tab is active. The 'Activity notifications' section is empty. The 'Activity history' section shows two entries:

- At 2024-05-06T20:19:27Z, a user request created an AutoScalingGroup changing the desired capacity from 0 to 2. An instance was started in response to a difference between desired and actual capacity, increasing the capacity from 0 to 2.
- At 2024-05-06T20:19:42Z, a user request created an AutoScalingGroup changing the desired capacity from 0 to 2. An instance was started in response to a difference between desired and actual capacity, increasing the capacity from 0 to 2.

You have the option to set scale-in protection for instances. This means that these instances won't be terminated by the auto-scaling process. This feature is particularly useful during troubleshooting or maintenance activities when you don't want instances to be terminated unexpectedly.

We'll wait for a few minutes to ensure that all instances are in a healthy state. Once they are marked as healthy, we'll explore some additional options. Currently, both instances are running and healthy. You should also see them updated in the target group.

The screenshot shows the AWS Auto Scaling group instance management page for 'weh-health-asg'. The 'Instances' section lists two instances:

- 1d9fb054faefb0: Instance type t2.micro, Launch template web-health-template, Availability Zone ap-south-1a, Health status Healthy.
- 1d9fb054faefb1: Instance type t2.micro, Launch template web-health-template, Availability Zone ap-south-1a, Health status Healthy.

Context menus are open for both instances, showing options like 'Set to Standby', 'Set to Service', 'Remove node in protection', and 'Remove node in operation'.

When managing instances within an auto-scaling group, it's crucial to remember that these instances are dynamic. They are automatically created and deleted by the autoscaling group based on the defined policies and metrics. Therefore, manual changes to instances should be avoided to prevent unexpected behavior.

Any changes to instances should be made through the launch template. This includes changes such as instance type, security group, application upgrades, or software installations. If you need to

make changes, create a new launch template or a new version of the existing launch template with the desired configurations.

To apply changes to instances, you can edit the auto-scaling group and select the new launch template. Once the changes are saved, you must initiate an instance refresh to apply the changes to existing instances. During instance refresh, instances are terminated and recreated based on the updated launch template.

The screenshot shows the 'Start instance refresh' page for the 'web-health-asg' launch configuration. It includes sections for availability settings, instance replacement method, and a healthy percentage input field.

Availability settings
Use the instance replacement method and instance warmup to control capacity and availability. The amount of capacity available will vary between the specified minimum and maximum values.

Instance replacement method
Choose the best option for the availability of your application

- Prioritize availability**
 Launch before terminating
Launch new instances and wait for them to be ready before terminating others. This allows you to go above your desired capacity by a given percentage and may temporarily increase costs.
- Control costs**
 Terminate and launch
Terminate and launch instances at the same time. This allows you to go below your desired capacity by a given percentage and may temporarily reduce availability.
- Flexible**
 Custom behavior
Set custom values for the minimum and maximum amount of available capacity. This gives you greater flexibility in setting how far below and over your desired capacity EC2 Auto Scaling goes when replacing instances.

Set healthy percentage
Set the minimum percentage of the desired capacity that must be healthy and ready for use for EC2 Auto Scaling to proceed with replacing instances.

Min: 90 Max: 100 % of 2 instances

Note: This range is currently within your group's scaling limits.

When initiating an instance refresh, ensure that you specify a minimum healthy percentage. This ensures that a certain percentage of instances remain healthy during the refresh process, preventing disruptions to your application. It's recommended to set a higher healthy percentage to maintain application stability.

Once changes are applied and instance refresh is initiated, monitor the auto-scaling group activity to ensure that new instances are launched successfully and the application remains available. Avoid

manual termination of instances unless necessary for troubleshooting or maintenance purposes.

Please note — instance: i-02e53d55300f6659a I stopped purposefully.

| Activity history (5) | | | |
|---------------------------------|---|---|------------------------------------|
| Status | Description | Cause | Start time |
| Not yet in service | Launching a new EC2 instance: i-0fa1d8665ea7326b | At 2024-05-06T20:27:42Z an instance was launched in response to an unhealthy instance needing to be replaced. | 2024 May 07, 01:57:45 AM +05:30 |
| Connection draining in progress | Terminating EC2 instance: i-0fb5453300f6659a - Waiting For ELB Connection Draining. | At 2024-05-06T20:27:41Z an instance was taken out of service in response to an EC2 health check indicating it has been terminated or stopped. | 2024 May 07, 01:57:41 AM +05:30 |
| Successful | Launching a new EC2 instance: i-02e53d5300f6659a | At 2024-05-06T20:25:36Z an instance was launched in response to an unhealthy instance needing to be replaced. | 2024 May 07, 01:55:38 AM +05:30 |
| Connection draining in progress | Terminating EC2 instance: i-0db3b0954540e2bb - Waiting For ELB Connection Draining. | At 2024-05-06T20:25:36Z an instance was taken out of service in response to an ELB system health check failure. | 2024 May 07, 01:55:36 AM +05:30 |
| Successful | Launching a new EC2 instance: i-0ca960e3ad2a8ff | At 2024-05-06T20:19:27Z a user request created an AutoScalingGroup changing the desired capacity from 0 to 2. At 2024-05-06T20:19:42Z an instance was started in response to a difference between desired and actual capacity, increasing the capacity from 0 to 2. | 2024 May 07, 01:49:44 AM +05:30 |
| Successful | Launching a new EC2 instance: i-0dbb5b0954540e2bb6 | At 2024-05-06T20:19:27Z a user request created an AutoScalingGroup changing the desired capacity from 0 to 2. At 2024-05-06T20:19:42Z an instance was started in response to a difference between desired and actual capacity, increasing the capacity from 0 to 2. | 2024 May 07, 01:49:44 AM +05:30 |

We can see 2 instances are running:

| Instances (2) Info | | | | | | | | | | |
|--|----------------------|-------------------------|---------------|-------------------|-------------------------------|-------------------|-------------------------|---------------|-------------------------|----------------------------------|
| Find instance by attribute or tag [more options] | | | | | | | | | | |
| Running | | | | | | | | | | |
| Name | Instance ID | Instance state | Instance type | Status check | Alarm status | Availability Zone | Public IPv4 DNS | Public IPv4 | Actions | Launch Instances |
| web-health | i-0dbb5b0954540e2bb6 | Running | t2.micro | 3/3 checks passed | View alarms + | ap-south-1b | ec2-13-235-69-162.ap... | 13.235.69.162 | Actions | Launch Instances |
| web-health | i-0ca960e3ad2a8ff | Running | t2.micro | 3/3 checks passed | View alarms + | ap-south-1a | ec2-3-110-214-146.ap... | 3.110.214.146 | Actions | Launch Instances |

By following these best practices, you can effectively manage instances within your auto-scaling group while ensuring the reliability and scalability of your application.

In conclusion, auto-scaling groups in AWS offer a powerful solution for managing the scalability and availability of your application infrastructure. By dynamically adjusting the number of EC2 instances based on predefined policies and metrics, auto-scaling groups enable you to efficiently handle fluctuating traffic loads while maintaining optimal performance and cost-effectiveness.

Key points to remember include:

1. Dynamic Nature: Auto-scaling groups automatically create and terminate EC2 instances based on configured scaling policies,

ensuring that your application can handle varying levels of demand without manual intervention.

- 2.Launch Templates: Changes to instance configurations should be made through launch templates, allowing for easy management of instance settings such as instance type, security groups, and software configurations.
- 3.Instance Refresh: When updating configurations, initiate an instance refresh to apply changes to existing instances. Specify a minimum healthy percentage to ensure application stability during the refresh process.
- 4.Avoid Manual Changes: Manual changes to instances within auto-scaling groups should be avoided to prevent conflicts with scaling policies and unexpected behavior. Instead, utilize launch templates for any necessary updates.
- 5.Monitoring and Maintenance: Regularly monitor the auto-scaling group activity to ensure that instances are being created and terminated as expected. Conduct troubleshooting and maintenance activities carefully to avoid disruptions to your application.

Amazon S3



Amazon S3 (Simple Storage Service):

Amazon S3 is a scalable object storage service provided by AWS. It allows users to store and retrieve any amount of data from anywhere on the internet. S3 is widely used for its durability, scalability, and ease of integration with other AWS services.

How S3 Works:

1. Object-based Storage: S3 stores data as objects, which include the data itself, metadata, and a unique identifier. These objects are stored in "buckets."
2. Bucket Structure: A bucket is a container for objects. Users can upload objects to buckets, and each object is identified by a unique key.
3. Global Access: S3 allows data access over the internet using API calls, SDKs, or the AWS Management Console.
4. High Durability & Availability: S3 is designed for 99.99999999% (11 nines) durability, with data automatically replicated across multiple availability zones.

Use Cases of S3:

1. Data Storage & Backup: Regular backups for applications, systems, and databases.
2. Big Data & Analytics: Storing large datasets for analysis and processing.
3. Website Hosting: Hosting static websites (HTML, CSS, JavaScript).
4. Content Distribution: Storing media files for distribution (e.g., videos, software).
5. Disaster Recovery: Storing critical data for quick recovery

Difference between EBS and S3:

EBS (Elastic Block Store):

- Type: Block storage for EC2 instances.
- Use Case: Low-latency storage for databases, file • •



- Attachment: Only attached to one EC2 instance at a time.
- Type: Object storage for unstructured data (images, backups).
- Use Case: Highly scalable and accessible from anywhere.
- Attachment: Not tied to any single instance, can be accessed globally

Cold Data vs. Hot Data:

- Hot Data: Frequently accessed data requiring lowlatenc storage (e.g., S3 Standard or EBS).
- Cold Data: Rarely accessed data stored in lower-cost solutions (e.g., S3 Glacier) for long-term retention.

Object Storage Class vs. Block Storage Class:

- Object Storage (S3): Stores data as objects with unique identifiers, suited for unstructured data (documents, media).
- Block Storage (EBS): Stores data in blocks, similar to a hard drive, designed for structured data (databases, file systems).



S3 Storage Classes:

1.S3 Standard: General-purpose storage for frequently accessed data. All Bookmarks

Try Premium for 10

2.S3 Intelligent-Tiering: Automatically moves data between frequent and infrequent access tiers to save costs.

3.S3 Standard-IA (Infrequent Access): Lower cost for infrequently accessed data.

4.S3 One Zone-IA: Similar to Standard-IA but stores data in a single availability zone.

5.S3 Glacier: Low-cost archival storage for data rarely accessed.

6.S3 Glacier Deep Archive: Lowest-cost storage for long-term, rarely accessed data.

| S3 Standard | S3 Intelligent-Tiering | S3 Standard-IA | S3 One Zone-IA | S3 Glacier | S3 Glacier Deep Archive |
|--|---|---|---|---|--|
| Frequent | | | Access frequency | | |
| <ul style="list-style-type: none">- Active, frequently accessed data- Milliseconds access- ≥ 5 AZ- \$0.0210/GB- Monitoring fee per object- Min storage duration | <ul style="list-style-type: none">- Data with changing access patterns- Milliseconds access- ≥ 5 AZ- \$0.0210 to \$0.0225/GB- Monitoring fee per object- Min object size | <ul style="list-style-type: none">- Infrequently accessed data- Milliseconds access- ≥ 5 AZ- \$0.0125/GB- Retrieval fee per GB- Min storage duration- Min object size | <ul style="list-style-type: none">- Re-creatable, less accessed data- Milliseconds access- 1 AZ- \$0.0190/GB- Retrieval fee per GB- Min storage duration- Min object size | <ul style="list-style-type: none">- Archive data- Select minutes or hours- ≥ 5 AZ- \$0.0043/GB- Retrieval fee per GB- Min storage duration- Min object size | <ul style="list-style-type: none">- Long-term archive data- Select hours- ≥ 5 AZ- \$0.0086/GB- Retrieval fee per GB- Min storage duration |

Advantages of S3:

- Scalability: Automatically scales to accommodate any amount of data.
- Durability: 99.99999999% durability with automatic replication across multiple locations.
- Security: Supports encryption, access control policies, and integrates with AWS Identity and Access Management (IAM).
- Cost-Effective: Pay for what you use, with multiple storage classes to optimize cost based on data access frequency.

Disadvantages of S3:

- Access Latency: Not suitable for low-latency, high-performance applications like databases (use EBS for that).
- Complexity in Management: Managing large datasets or setting up lifecycle policies may become complex for certain use cases.

- Data Retrieval Costs: Higher retrieval costs in certain storage classes like S3 Glacier or Glacier Deep Archive.

Features of S3:

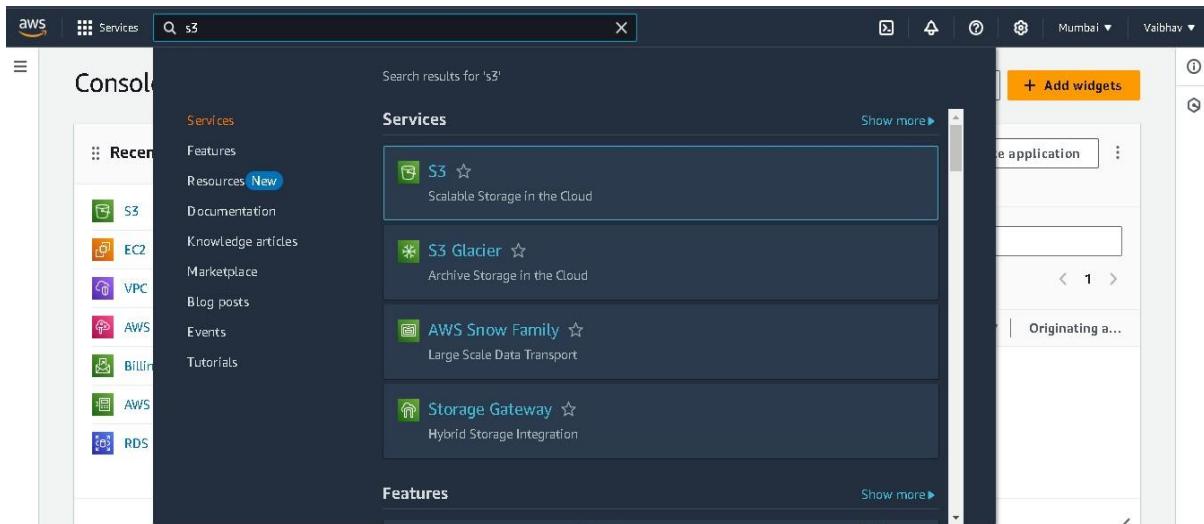
1. Versioning: Tracks changes to objects, allowing you to revert to previous versions.
2. Lifecycle Management: Automates the movement of data to lower-cost storage classes.
3. Cross-Region Replication: Automatically replicates data across AWS regions for better durability.
4. Encryption: Supports server-side and client-side encryption for securing data.
5. Access Control: Fine-grained control over who can access or modify data using IAM policies, bucket Policies, and access control lists (ACLs).



Step-by-Step Guide for
Uploading a File to Amazon S3:

Step 1: Log in to the AWS Management console

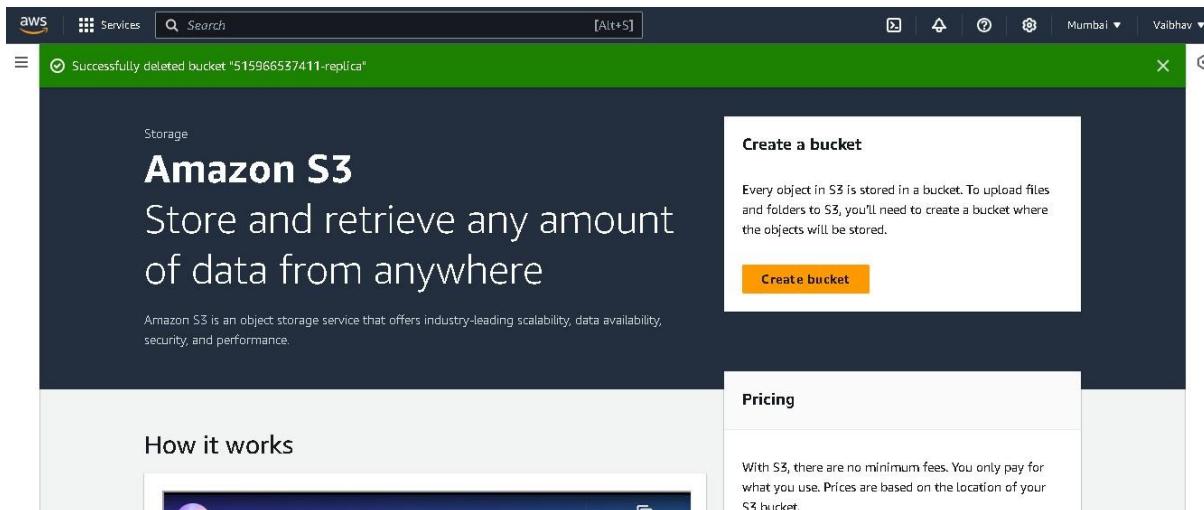
1. Go to the [AWS Management Console](#) and sing in with your credentials.
2. In the search bar, type S3 and click on the Amazon S3 service.



Step 2: Create a New S3 Bucket

1. Once in the S3 dashboard, click Create Bucket.
2. Enter a Bucket Name (it must be unique globally across all AWS users)

VEERA



The screenshot shows the AWS Create Bucket wizard. The first step, 'General configuration', is displayed. It includes fields for 'AWS Region' (Asia Pacific (Mumbai) ap-south-1), 'Bucket name' (myawsbucket), and a note about naming rules. The second step, 'Object Ownership', is partially visible below it.

General configuration

AWS Region
Asia Pacific (Mumbai) ap-south-1

Bucket name [Info](#)
myawsbucket

Bucket name must be unique within the global namespace and follow the bucket naming rules. [See rules for bucket naming](#)

Copy settings from existing bucket - *optional*
Only the bucket settings in the following configuration are copied.
[Choose bucket](#)

Format: s3://bucket/prefix

Object Ownership [Info](#)

Control ownership of objects written to this bucket from other AWS accounts and the use of access control lists (ACLs). Object ownership determines who can specify access to objects.

3. Select the AWS Region where you want to store your data (choose a region closest to your location for better performance.)
4. Configure bucket settings.

- Block Public Access: You can keep public access blocked unless you want to make your files accessible to the public.
- Bucket Versioning: You can enable or disable versioning. This feature keeps multiple versions of an object.

Multicloud Devops by veera nareshit

Object Ownership [Info](#)

Control ownership of objects written to this bucket from other AWS accounts and the use of access control lists (ACLs). Object ownership determines who can specify access to objects.

ACLs disabled (recommended)

All objects in this bucket are owned by this account. Access to this bucket and its objects is specified using only policies.

ACLs enabled

Objects in this bucket can be owned by other AWS accounts. Access to this bucket and its objects can be specified using ACLs.

Object Ownership

Bucket owner enforced

Block Public Access settings for this bucket

Public access is granted to buckets and objects through access control lists (ACLs), bucket policies, access point policies, or all. In order to ensure that public access to this bucket and its objects is blocked, turn on Block all public access. These settings apply only to this bucket and its access points. AWS recommends that you turn on Block all public access, but before applying any of these settings, ensure that your applications will work correctly without public access. If you require some level of public access to this bucket or objects within, you can customize the individual settings below to suit your specific storage use cases. [Learn more](#)

VEERA

Bucket Versioning

Versioning is a means of keeping multiple variants of an object in the same bucket. You can use versioning to preserve, retrieve, and restore every version of every object stored in your Amazon S3 bucket. With versioning, you can easily recover from both unintended user actions and application failures. [Learn more](#)

Bucket Versioning

Disable

Enable

Block Public Access settings for this bucket

Public access is granted to buckets and objects through access control lists (ACLs), bucket policies, access point policies, or all. In order to ensure that public access to this bucket and its objects is blocked, turn on Block all public access. These settings apply only to this bucket and its access points. AWS recommends that you turn on Block all public access, but before applying any of these settings, ensure that your applications will work correctly without public access. If you require some level of public access to this bucket or objects within, you can customize the individual settings below to suit your specific storage use cases. [Learn more](#)

Block all public access
Turning this setting on is the same as turning on all four settings below. Each of the following settings are independent of one another.

- Block public access to buckets and objects granted through new access control lists (ACLs)**
S3 will block public access permissions applied to newly added buckets or objects, and prevent the creation of new public access ACLs for existing buckets and objects. This setting doesn't change any existing permissions that allow public access to S3 resources using ACLs.
- Block public access to buckets and objects granted through any access control lists (ACLs)**
S3 will ignore all ACLs that grant public access to buckets and objects.
- Block public access to buckets and objects granted through new public bucket or access point policies**
S3 will block new bucket and access point policies that grant public access to buckets and objects. This setting doesn't change any existing policies that allow public access to S3 resources.
- Block public and cross-account access to buckets and objects through any public bucket or access point policies**
S3 will ignore public and cross-account access for buckets or access points with policies that grant public access to buckets and objects.

5. Scroll down and click Create Bucket.

Encryption type [Info](#)

Server-side encryption with Amazon S3 managed keys (SSE-S3)
 Server-side encryption with AWS Key Management Service keys (SSE-KMS)
 Dual-layer server-side encryption with AWS Key Management Service keys (DSSE-KMS)
Secure your objects with two separate layers of encryption. For details on pricing, see DSSE-KMS pricing on the Storage tab of the [Amazon S3 pricing page](#).

Bucket Key
Using an S3 Bucket Key for SSE-KMS reduces encryption costs by lowering calls to AWS KMS. S3 Bucket Keys aren't supported for DSSE-KMS.
[Learn more](#)

Disable
 Enable

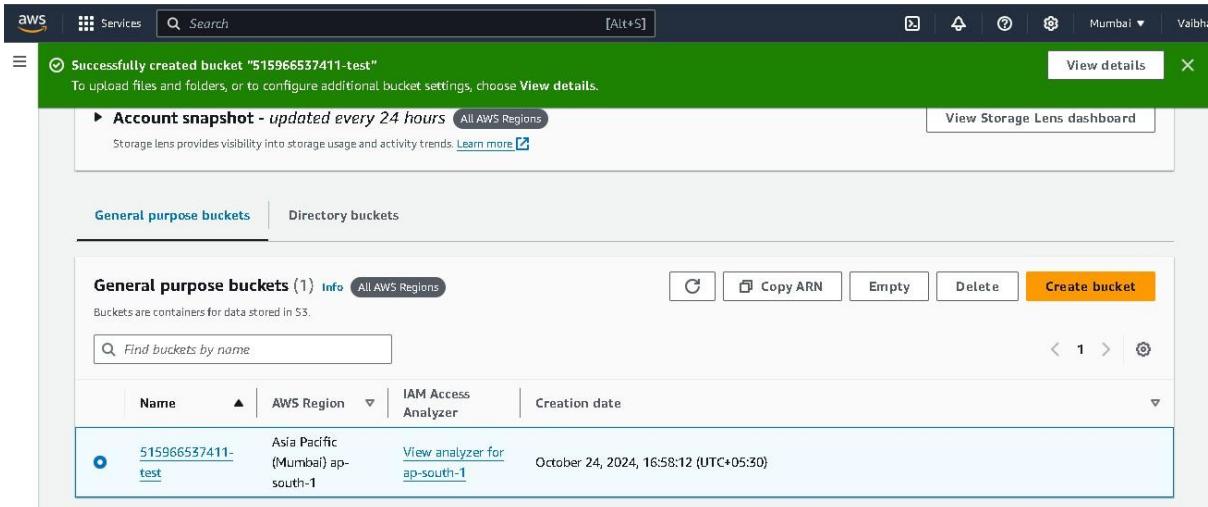
Advanced settings

Info After creating the bucket, you can upload files and folders to the bucket, and configure additional bucket settings.

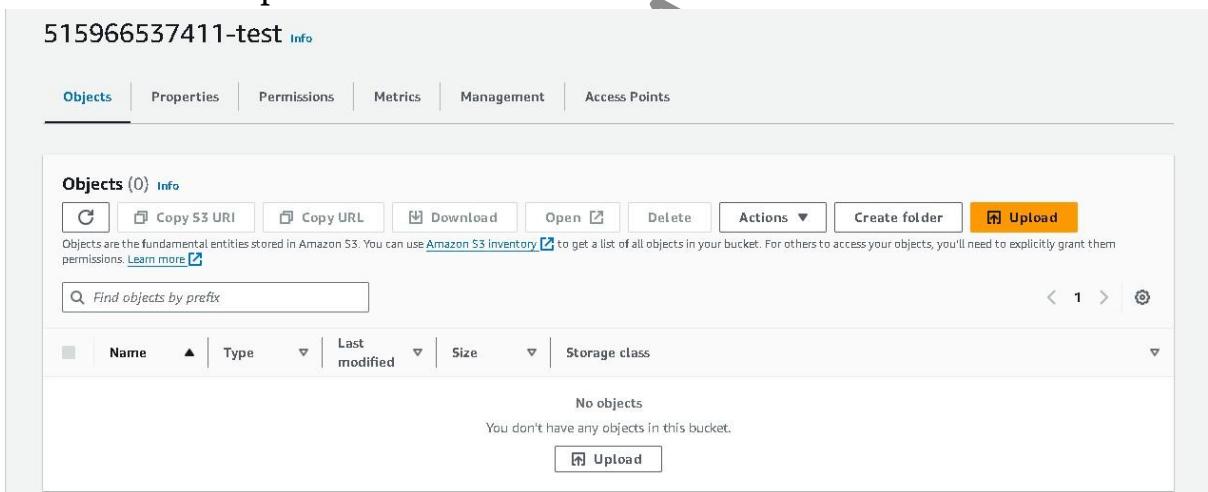
[Cancel](#) [Create bucket](#)

Step 3: Upload a File to the S3 bucket

- Once the bucket is created, click on your bucket's name to open it.



- Click upload on the bucket's dashboard.



- In the Upload window, either drag and drop a file from your local machine or click add files to browse for the file you want to upload.
- After adding the file, you can configure permission:

- Manage public permissions: Decide if the file should be public or private.

Multicloud Devops by veera nareshit

▼ Permissions
Grant public access and access to other AWS accounts.

Access control list (ACL)
Grant basic read/write permissions to other AWS accounts. [Learn more](#)

AWS recommends using S3 bucket policies or IAM policies for access control. [Learn more](#)

Access control list (ACL)
 Choose from predefined ACLs
 Specify individual ACL permissions
Predefined ACLs
 Private (recommended)
Only the object owner will have read and write access.
 Grant public-read access
Anyone in the world will be able to access the specified objects. The object owner will have read and write access. [Learn more](#)

Granting public-read access is not recommended
Anyone in the world will be able to access the specified objects. [Learn more](#)

I understand the risk of granting public-read access to the specified objects.

Upload [Info](#)

Add the files and folders you want to upload to S3. To upload a file larger than 160GB, use the AWS CLI, AWS SDK or Amazon S3 REST API. [Learn more](#)

Drag and drop files and folders you want to upload here, or choose **Add files** or **Add folder**.

| Files and folders (1 Total, 159.1 KB) | | Remove | Add files | Add folder |
|---|------------|------------------------|---------------------------|----------------------------|
| All files and folders in this table will be uploaded. | | | | |
| <input type="text"/> Find by name | | < 1 > | | |
| <input checked="" type="checkbox"/> | Name | Folder | Type | |
| <input checked="" type="checkbox"/> | RESUME.pdf | - | application/pdf | |

Destination [Info](#)

5. You can leave storage class settings as default (S3 Standard), or choose other storage classes like Glacier for archival.
6. Once you've set the file and permissions, click Upload.

Multicloud Devops by veera nareshit

The screenshot shows the AWS S3 console with a green success message at the top: "Upload succeeded" and "View details below." Below this is a "Summary" section with two rows: "Destination s3://515966537411-test" with "Succeeded" status and "1 file, 159.1 KB (100.00%)"; and "Failed" with "0 files, 0 B (0%)". At the bottom, there are tabs for "Files and folders" (selected) and "Configuration". Under "Files and folders", it says "(1 Total, 159.1 KB)" and lists a single file: "RESUME.pdf" (application/pdf, 159.1 KB, Succeeded).

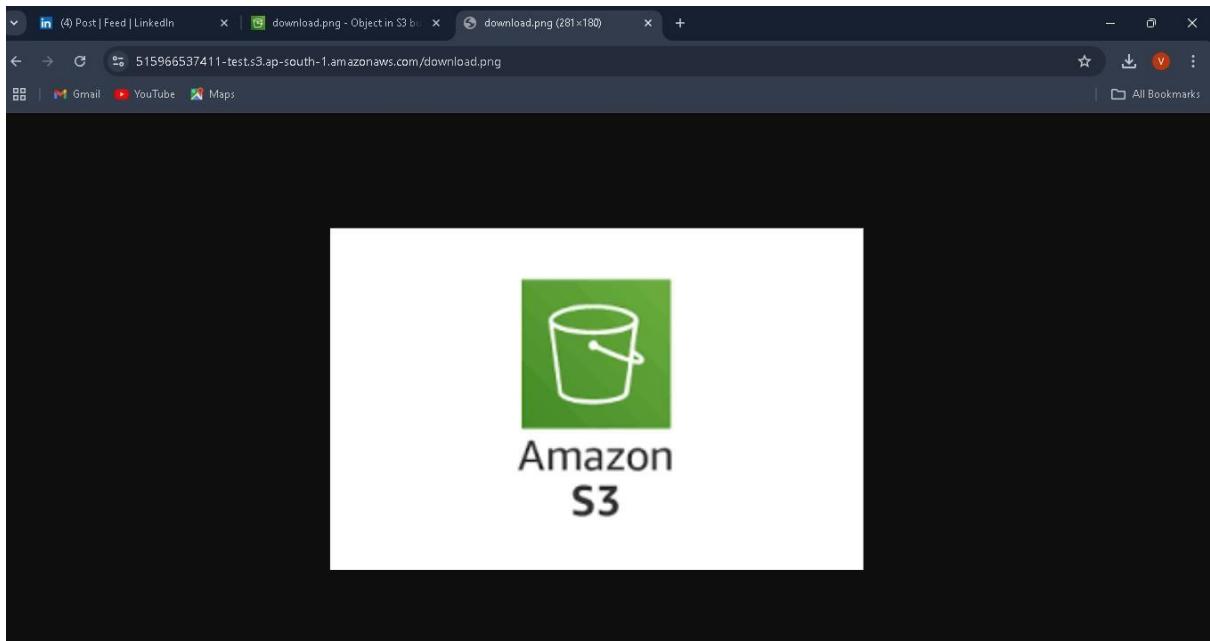
Step 4: Verify the Upload

- Once the file is uploaded, you will see it in the bucket. You can click on the file name to view its properties.

The screenshot shows the AWS S3 "Objects" page for the bucket "515966537411-test". It displays one object: "RESUME.pdf" (pdf, 159.1 KB, Standard storage class). The "Actions" dropdown menu is highlighted. The page includes standard S3 navigation and object management tools like Copy, Download, Open, Delete, and Create folder.

- The file has a URL that you can use to access it (based on permissions). The URL format will look like this:

`https://<bucket-name>.s3.<region>.amazonaws.com/<file-name>`



Step 5: Access and Manage Your File

- After uploading, you can manage your file's permissions, move it to different storage classes (like Glacier), or delete it.
- If you want to download the file from the bucket, You can either use the AWS CLI or the S3 console.

Additional Points:

- Public Access: If you want the file to be publicly accessible, make sure to set the object's permissions accordingly.
- Storage Classes: By default, objects are uploaded to the S3 Standard storage class. You can later change it to Glacier or Intelligent-Tiering, depending on your use case.

Reminder: Delete S3 Bucket After Completing Your Work

Once you've completed your tasks and no longer need the S3 bucket, make sure to delete it to avoid unnecessary storage costs. Deleting the bucket will permanently remove all data, so ensure you have saved any important files. Always clean up unused resources after finishing your work.

Multicloud Devops by veera nareshit

With these steps, you can successfully upload files from your local machine to Amazon S3!

VEERA

Multicloud Devops by veera nareshit

AWS IAM (Identity and Access Management Service)



AWS IAM is used to securely control individual and group access to AWS resources.

IAM makes it easy to provide multiple users secure access to AWS resources.

IAM can be used to manage:

- Users.
- Groups.
- Access policies.
- Roles.
- User credentials.
- User password policies.
- Multi-factor authentication (MFA).
- API keys for programmatic access (CLI).

Provides centralized control of your AWS account.

Enables shared access to your AWS account.

By default, new users are created with NO access to any AWS services – they can only login to the AWS console.

Permission must be explicitly granted to allow a user to access an AWS service.

IAM users are individuals who have been granted access to an AWS account.

Each IAM user has three main components:

- A username.
- A password.

- Permissions to access various resources.

You can apply granular permissions with IAM.

You can assign users individual security credentials such as access keys, passwords, and multi-factor authentication devices.

IAM is not used for application-level authentication.

Identity Federation (including AD, Facebook etc). can be configured allowing secure access to resources in an AWS account without creating an IAM user account.

Multi-factor authentication (MFA) can be enabled/enforced for the AWS account and for individual users under the account.

MFA uses an authentication device that continually generates random, six-digit, single-use authentication codes.

You can authenticate using an MFA device in the following three ways:

- Through the AWS Management Console – the user is prompted for a user name, password, and authentication code.
- Using the AWS API – restrictions are added to IAM policies and developers can request temporary security credentials and pass MFA parameters in their AWS STS API requests.
- Using the AWS CLI by obtaining temporary security credentials from STS (aws sts getsession-token).

Want to see how to setup MFA? In the brief AWS Hands-on Labs video tutorial below, you'll learn how to activate a virtual Multi-factor Authentication (MFA) for your AWS Root Account. In under 5 minutes, we cover: Deleting the Root Account Access Key and Activating MultiFactor Authentication.

It is a best practice to use MFA for all users and to use U2F or hardware MFA devices for all privileged users.

IAM is universal (global) and does not apply to regions.

IAM is eventually consistent.

IAM replicates data across multiple data centres around the world.

The “root account” is the account created when you setup the AWS account. It has complete Admin access and is the only account that has this access by default.

It is a best practice to not use the root account for anything other than billing.

Power user access allows all permissions except the management of groups and users in IAM.

Temporary security credentials consist of the AWS access key ID, secret access key, and security token.

IAM can assign temporary security credentials to provide users with temporary access to services/resources.

To sign-in you must provide your account ID or account alias in addition to a user name and password.

The sign-in URL includes the account ID or account alias, e.g.:

<https://My AWS Account ID.signin.aws.amazon.com/console/>.

Alternatively, you can sign-in at the following URL and enter your account ID or alias manually:

<https://console.aws.amazon.com/>.

IAM integrates with many different AWS services.

IAM supports PCI DSS compliance.

AWS recommend that you use the AWS SDKs to make programmatic API calls to IAM.

However, you can also use the IAM Query API to make direct calls to the IAM web service.

IAM Elements

Principals:

- An entity that can take an action on an AWS resource.
- Your administrative IAM user is your first principal.
- You can allow users and services to assume a role.
- IAM supports federated users.
- IAM supports programmatic access to allow an application to access your AWS account.
- IAM users, roles, federated users, and applications are all AWS principals.

Requests:

- Principals send requests via the Console, CLI, SDKs, or APIs.
- Requests are:
 - Actions (or operations) that the principal wants to perform.
 - Resources upon which the actions are performed.
 - Principal information including the environment from which the request was made.
- Request context – AWS gathers the request information:
 - Principal (requester).
 - Aggregate permissions associated with the principal.
 - Environment data, such as IP address, user agent, SSL status etc.
 - Resource data, or data that is related to the resource being requested.

Authentication:

- A principal sending a request must be authenticated to send a request to AWS.
- To authenticate from the console, you must sign in with your user name and password.
- To authenticate from the API or CLI, you must provide your access key and secret key.

Authorization:

- IAM uses values from the request context to check for matching policies and determines whether to allow or deny the request.
- IAM policies are stored in IAM as JSON documents and specify the permissions that are allowed or denied.
- IAM policies can be:
 - User (identity) based policies.
 - Resource-based policies.

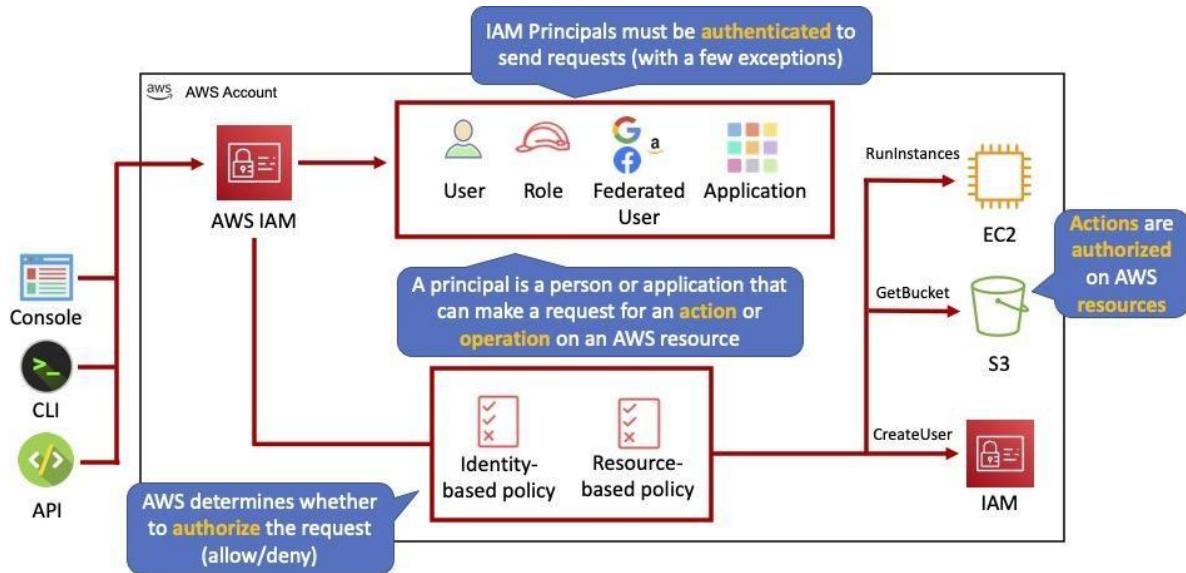
- IAM checks each policy that matches the context of your request.
- If a single policy has a deny action IAM denies the request and stops evaluating (explicit deny).
- Evaluation logic:
 - By default all requests are denied (implicit deny).
 - An explicit allow overrides the implicit deny.
 - An explicit deny overrides any explicit allows.
- Only the root user has access to all resources in the account by default.

Actions:

- Actions are defined by a service.
- Actions are the things you can do to a resource such as viewing, creating, editing, deleting.
- Any actions on resources that are not explicitly allowed are denied.
- To allow a principal to perform an action you must include the necessary actions in a policy that applies to the principal or the affected resource.

Resources:

- A resource is an entity that exists within a service.
- E.g. EC2 instances, S3 buckets, IAM users, and DynamoDB tables.
- Each AWS service defines a set of actions that can be performed on the resource.
- After AWS approves the actions in your request, those actions can be performed on the related resources within your account.



Authentication Methods

Console password:

- A password that the user can enter to sign into interactive sessions such as the AWS Management Console.
- You can allow users to change their own passwords.
- You can allow selected IAM users to change their passwords by disabling the option for all users and using an IAM policy to grant permissions for the selected users.

Access Keys:

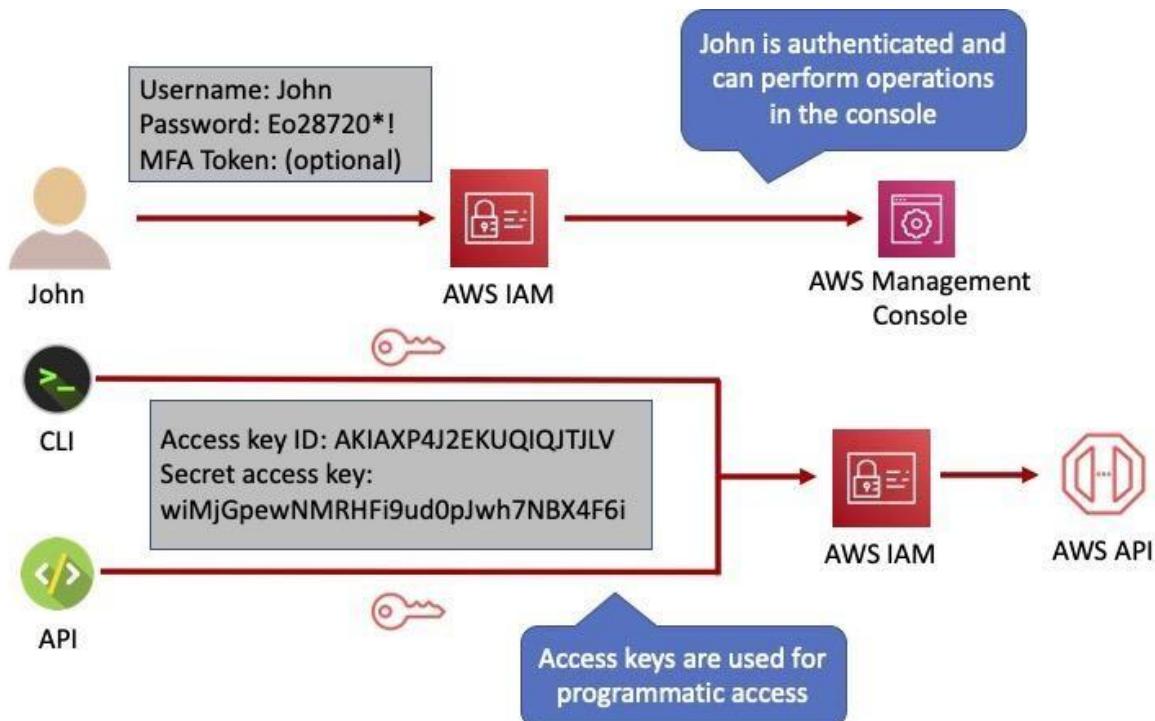
- A combination of an access key ID and a secret access key.
- You can assign two active access keys to a user at a time.
- These can be used to make programmatic calls to AWS when using the API in program code or at a command prompt when using the AWS CLI or the AWS PowerShell tools.
- You can create, modify, view, or rotate access keys.
- When created IAM returns the access key ID and secret access key.
- The secret access key is returned only at creation time and if lost a new key must be created.
- Ensure access keys and secret access keys are stored securely.
- Users can be given access to change their own keys through IAM policy (not from the console).

- You can disable a user's access key which prevents it from being used for API calls.

Server certificates:

- SSL/TLS certificates that you can use to authenticate with some AWS services.
- AWS recommends that you use the AWS Certificate Manager (ACM) to provision, manage and deploy your server certificates.
- Use IAM only when you must support HTTPS connections in a region that is not supported by ACM.

The following diagram shows the different methods of authentication available with IAM:



IAM Users

An IAM user is an entity that represents a person or service.

Can be assigned:

- An access key ID and secret access key for programmatic access to the AWS API, CLI, SDK, and other development tools.
- A password for access to the management console.

By default users cannot access anything in your account.

The account root user credentials are the email address used to create the account and a password.

The root account has full administrative permissions, and these cannot be restricted.

Best practice for root accounts:

- Don't use the root user credentials.
- Don't share the root user credentials.
- Create an IAM user and assign administrative permissions as required.
- Enable MFA.

IAM users can be created to represent applications, and these are known as "service accounts".

You can have up to 5000 users per AWS account.

Each user account has a friendly name and an ARN which uniquely identifies the user across AWS.

A unique ID is also created which is returned only when you create the user using the API, Tools for Windows PowerShell, or the AWS CLI.

You should create individual IAM accounts for users (best practice not to share accounts).

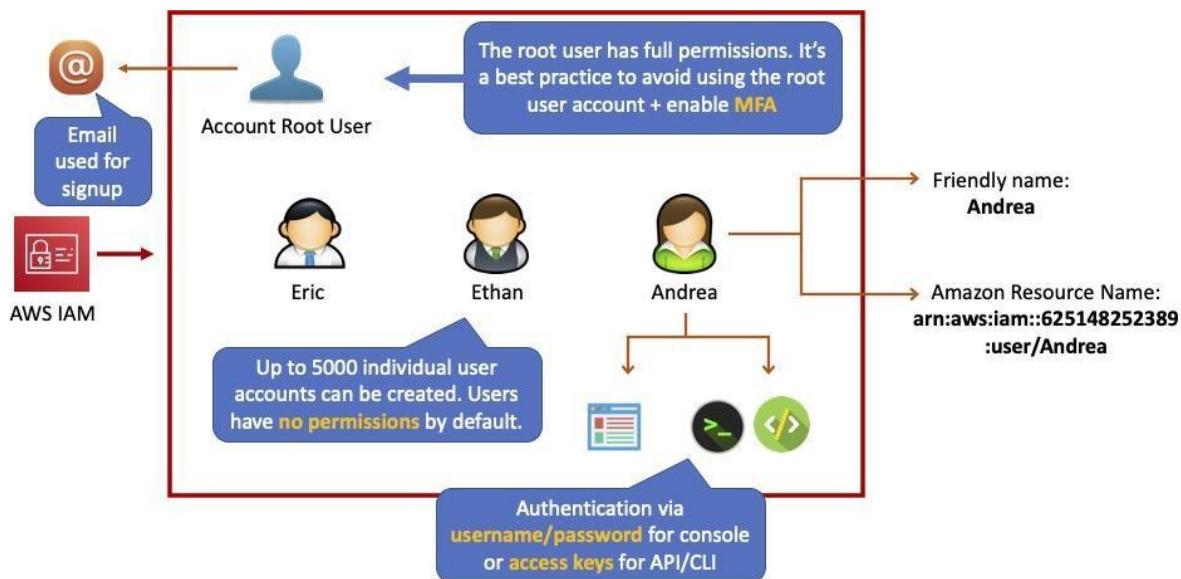
The Access Key ID and Secret Access Key are not the same as a password and cannot be used to login to the AWS console.

The Access Key ID and Secret Access Key can only be generated once and must be regenerated if lost.

A password policy can be defined for enforcing password length, complexity etc. (applies to all users).

You can allow or disallow the ability to change passwords using an IAM policy.

Access keys and passwords should be changed regularly.



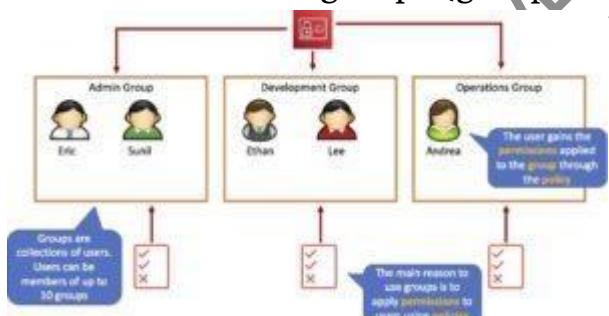
Groups

Groups are collections of users and have policies attached to them. A group is not an identity and cannot be identified as a principal in an IAM policy.

Use groups to assign permissions to users.

Use the principle of least privilege when assigning permissions.

You cannot nest groups (groups within groups).



Roles

Roles are created and then “assumed” by trusted entities and define a set of permissions for making AWS service requests.

With IAM Roles you can delegate permissions to resources for users and services without using permanent credentials (e.g. user name and password).

IAM users or AWS services can assume a role to obtain temporary security credentials that can be used to make AWS API calls.

You can delegate using roles.

There are no credentials associated with a role (password or access keys).

IAM users can temporarily assume a role to take on permissions for a specific task.

A role can be assigned to a federated user who signs in using an external identity provider. Temporary credentials are primarily used with IAM roles and automatically expire.

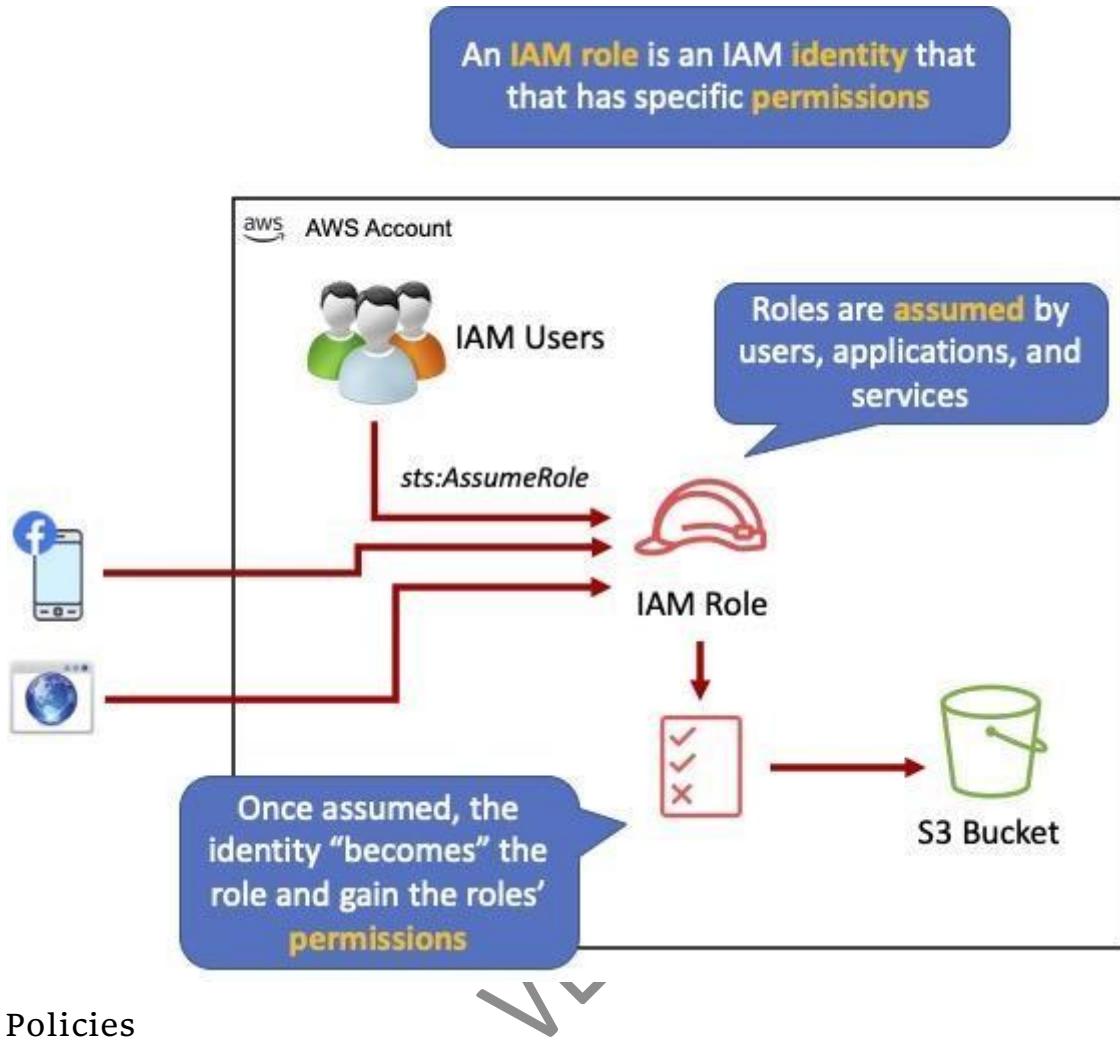
Roles can be assumed temporarily through the console or programmatically with the AWS CLI, Tools for Windows PowerShell, or API.

IAM roles with EC2 instances:

- IAM roles can be used for granting applications running on EC2 instances permissions to AWS API requests using instance profiles.
- Only one role can be assigned to an EC2 instance at a time.
- A role can be assigned at the EC2 instance creation time or at any time afterwards.
- When using the AWS CLI or API instance profiles must be created manually (it's automatic and transparent through the console).
- Applications retrieve temporary security credentials from the instance metadata.

Role Delegation:

- Create an IAM role with two policies:
 - Permissions policy – grants the user of the role the required permissions on a resource.
 - Trust policy – specifies the trusted accounts that are allowed to assume the role.
- Wildcards (*) cannot be specified as a principal.
- A permissions policy must also be attached to the user in the trusted account.



Policies are documents that define permissions and can be applied to users, groups, and roles.

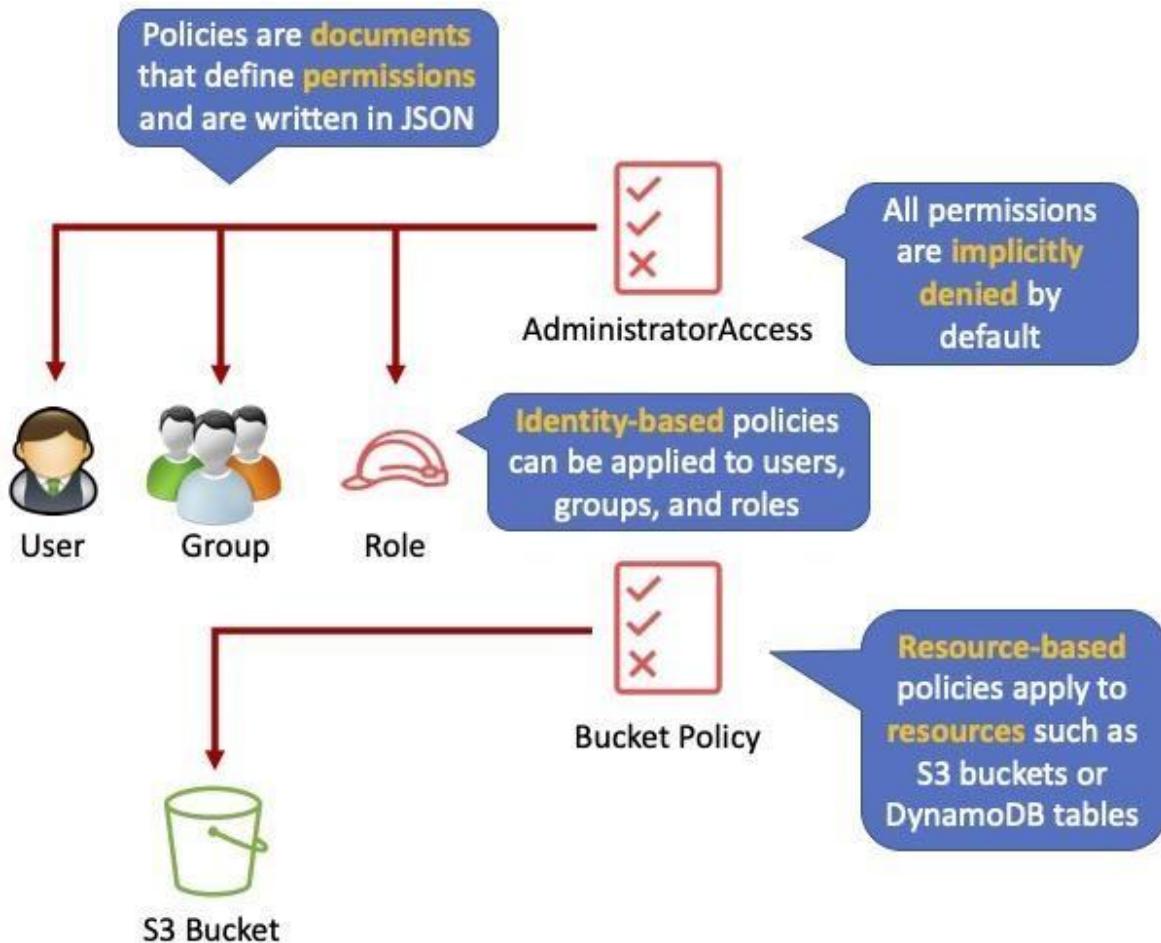
Policy documents are written in JSON (key value pair that consists of an attribute and a value).

All permissions are implicitly denied by default.

The most restrictive policy is applied.

The IAM policy simulator is a tool to help you understand, test, and validate the effects of access control policies.

The Condition element can be used to apply further conditional logic.



Inline Policies vs Managed Policies

There are 3 types of policies:

- Managed policies.
- Customer managed policies.
- Inline policies.

Managed Policy:

- Created and administered by AWS.
- Used for common use cases based on job function.
- Save you having to create policies yourself.
- Can be attached to multiple users, groups, or roles within and across AWS accounts.
- Cannot change the permissions assigned.

Customer Managed Policy:

- Standalone policy that you create and administer in your own AWS account.
- Can be attached to multiple users, groups, and roles – but only within your own account.
- Can be created by copying an existing managed policy and then customizing it.
- Recommended for use cases where the existing AWS Managed Policies don't meet the needs of your environment.

Inline Policy:

- Inline policies are embedded within the user, group, or role to which it is applied.
- Strict 1:1 relationship between the entity and the policy.
- When you delete the user, group, or role in which the inline policy is embedded, the policy will also be deleted.
- In most cases, AWS recommends using Managed Policies instead of inline policies.
- Inline policies are useful when you want to be sure that the permissions in a policy are not inadvertently assigned to any other user, group, or role.

AWS Managed and Customer Managed Policies

An AWS managed policy is a standalone policy that is created and administered by AWS.

Standalone policy means that the policy has its own Amazon Resource Name (ARN) that includes the policy name.

AWS managed policies are designed to provide permissions for many common use cases.

You cannot change the permissions defined in AWS managed policies. Some AWS managed policies are designed for specific job functions.

The job-specific AWS managed policies include:

- Administrator.
- Billing.

- Database Administrator.
- Data Scientist.
- Developer Power User.

VEERA

Network Administrator.

- Security Auditor.
- Support User.
- System Administrator.
- View-Only User.

You can create standalone policies that you administer in your own AWS account, which we refer to as customer managed policies.

You can then attach the policies to multiple principal entities in your AWS account.

When you attach a policy to a principal entity, you give the entity the permissions that are defined in the policy.

IAM Policy Evaluation Logic

By default, all requests are implicitly denied. (Alternatively, by default, the AWS account root user has full access).

An explicit allow in an identity-based or resource-based policy overrides this default.

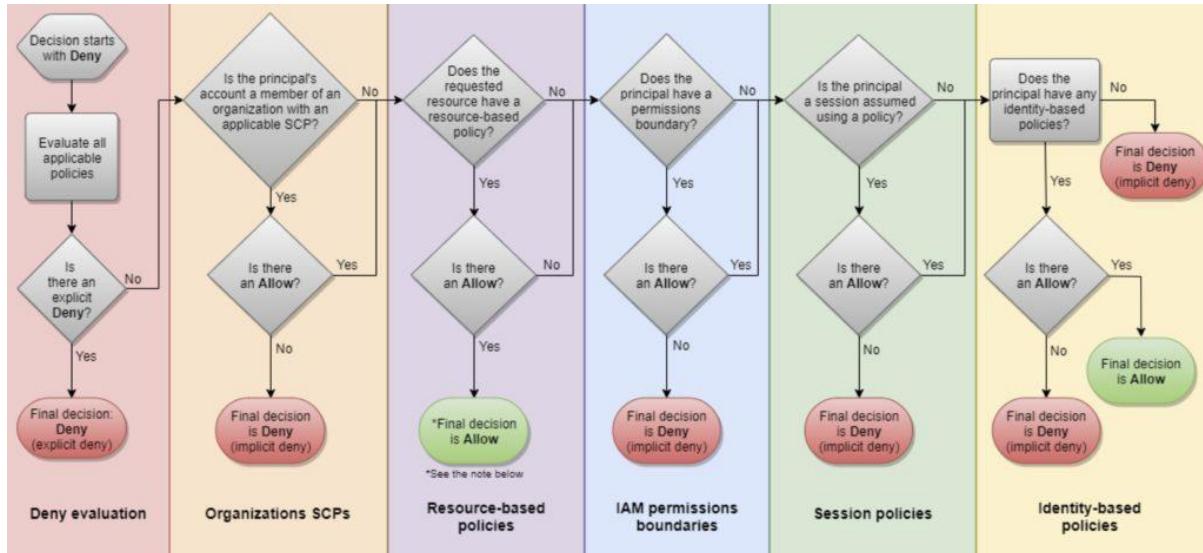
If a permissions boundary, Organizations SCP, or session policy is present, it might override the allow with an implicit deny.

An explicit deny in any policy overrides any allows.

A few concepts should be known to understand the logic:

- Identity-based policies – Identity-based policies are attached to an IAM identity (user, group of users, or role) and grant permissions to IAM entities (users and roles).
- Resource-based policies – Resource-based policies grant permissions to the principal (account, user, role, or federated user) specified as the principal.
- IAM permissions boundaries – Permissions boundaries are an advanced feature that sets the maximum permissions that an identity-based policy can grant to an IAM entity (user or role).
- AWS Organizations service control policies (SCPs) – Organizations SCPs specify the maximum permissions for an organization or organizational unit (OU). Session policies – Session policies are advanced policies that you pass as parameters when you programmatically create a temporary session for a role or federated user.

The following flowchart details the IAM policy evaluation logic:

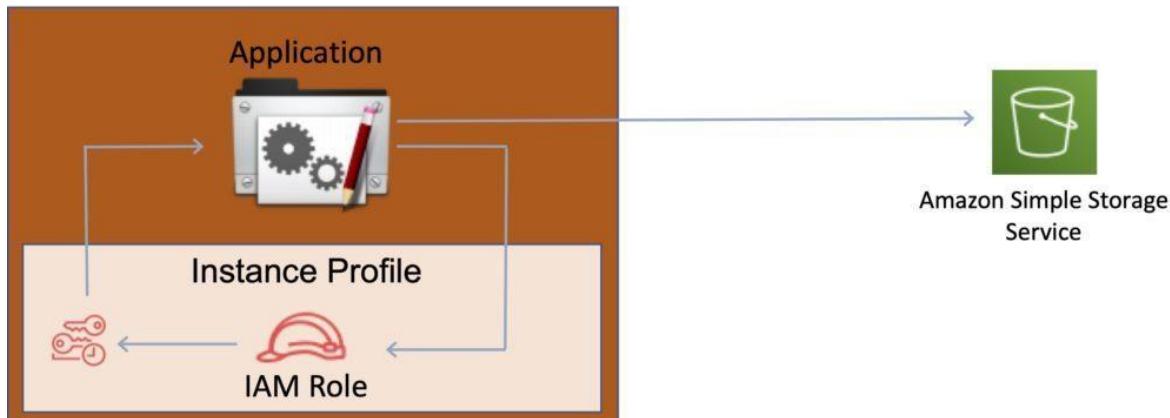


IAM Instance Profiles

An instance profile is a container for an IAM role that you can use to pass role information to an EC2 instance when the instance starts.

An instance profile can contain only one IAM role, although a role can be included in multiple instance profiles.

EC2 Instance



You can use the following AWS CLI commands to work with instance profiles in an AWS account:

- Create an instance profile: `aws iam create-instance-profile`
- Add a role to an instance profile: `aws iam add-role-to-instance-profile`
- List instance profiles: `aws iam list-instance-profiles`, `aws iam list-instance-profiles-forrole`
- Get information about an instance profile: `aws iam get-instance-profile`
- Remove a role from an instance profile: `aws iam remove-role-from-instance-profile`
- Delete an instance profile: `aws iam delete-instance-profile`

AWS Security Token Service

The AWS Security Token Service (STS) is a web service that enables you to request temporary, limited-privilege credentials for IAM users or for users that you authenticate (federated users).

By default, AWS STS is available as a global service, and all AWS STS requests go to a single endpoint at <https://sts.amazonaws.com>

You can optionally send your AWS STS requests to endpoints in any region (can reduce latency).

Credentials will always work globally.

STS supports AWS CloudTrail, which records AWS calls for your AWS account and delivers log files to an S3 bucket.

Temporary security credentials work almost identically to long-term access key credentials that IAM users can use, with the following differences:

- Temporary security credentials are short-term.
- They can be configured to last anywhere from a few minutes to several hours.
- After the credentials expire, AWS no longer recognizes them or allows any kind of access to API requests made with them.
- Temporary security credentials are not stored with the user but are generated dynamically and provided to the user when requested.
- When (or even before) the temporary security credentials expire, the user can request new credentials, if the user requesting them still has permission to do so.
- You do not have to distribute or embed long-term AWS security credentials with an application.
- You can provide access to your AWS resources to users without having to define an AWS identity for them (temporary security credentials are the basis for IAM Roles and ID Federation).
- The temporary security credentials have a limited lifetime, so you do not have to rotate them or explicitly revoke them when they're no longer needed.
- After temporary security credentials expire, they cannot be reused (you can specify how long the credentials are valid for, up to a maximum limit).

The AWS STS API action returns temporary security credentials that consist of:

- An access key which consists of an access key ID and a secret ID.
- A session token.
- Expiration or duration of validity.

- Users (or an application that the user runs) can use these credentials to access your resources.

With STS you can request a session token using one of the following APIs:

- AssumeRole – can only be used by IAM users (can be used for MFA).
- AssumeRoleWithSAML – can be used by any user who passes a SAML authentication response that indicates authentication from a known (trusted) identity provider.
- AssumeRoleWithWebIdentity – can be used by an user who passes a web identity token that indicates authentication from a known (trusted) identity provider.
- GetSessionToken – can be used by an IAM user or AWS account root user (can be used for MFA).
- GetFederationToken – can be used by an IAM user or AWS account root user.

AWS recommends using Cognito for identity federation with Internet identity providers.

Users can come from three sources.

Federation (typically AD):

- Uses SAML 2.0.
- Grants temporary access based on the users AD credentials.
- Does not need to be a user in IAM.
- Single sign-on allows users to login to the AWS console without assigning IAM credentials.

Federation with Mobile Apps:

- Use Facebook/Amazon/Google or other OpenID providers to login.

Cross Account Access:

- Lets users from one AWS account access resources in another.
- To make a request in a different account the resource in that account must have an attached resource-based policy with the permissions you need.

Or you must assume a role (identity-based policy) within that account with the permissions you need.

There are a couple of ways STS can be used.

Scenario 1:

1. Develop an Identity Broker to communicate with LDAP and AWS STS.
2. Identity Broker always authenticates with LDAP first, then with AWS STS.

- Application then gets temporary access to AWS resources.

Scenario 2:

- Develop an Identity Broker to communicate with LDAP and AWS STS.
- Identity Broker authenticates with LDAP first, then gets an IAM role associated with the user.
- Application then authenticates with STS and assumes that IAM role.
- Application uses that IAM role to interact with the service.

Cross Account Access

Useful for situations where an AWS customer has separate AWS account – for example for development and production resources.

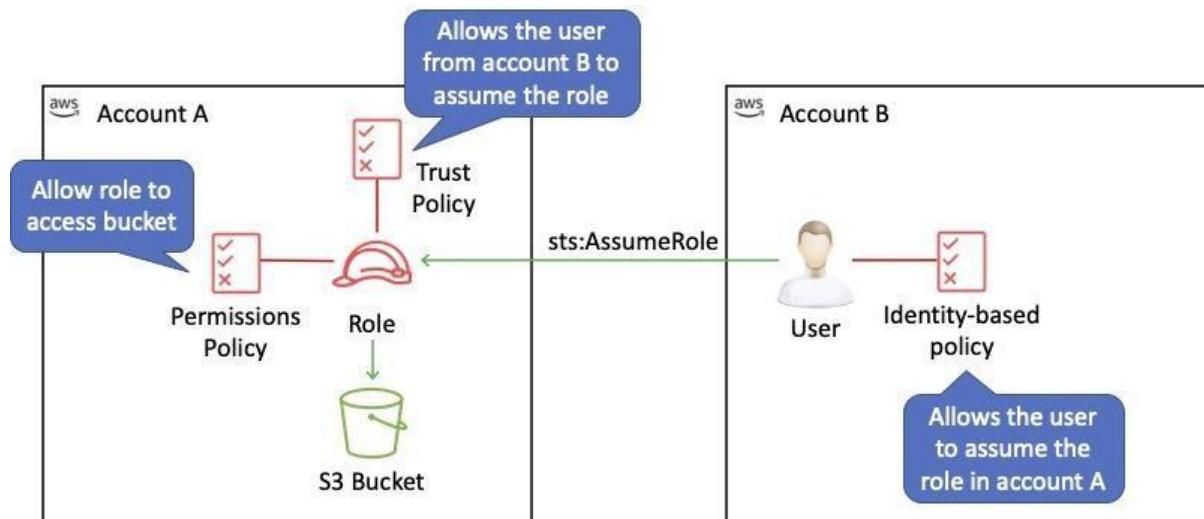
Cross Account Access makes it easier to work productively within a multi-account (or multirole) AWS environment by making it easy to switch roles within the AWS Management Console.

Can sign-in to the console using your IAM username and then switch the console to manage another account without having to enter another user name and password.

Let's users from one AWS account access resources in another.

To make a request in a different account the resource in that account must have an attached resource-based policy with the permissions you need.

Or you must assume a role (identity-based policy) within that account with the permissions you need.



IAM Best Practices

To secure AWS resources it is recommended that you follow these best practices:

- Lock away your AWS account root user access keys.

- Use roles to delegate permissions.
- Grant least privilege.
- Get started using permissions with AWS managed policies.
- Validate your policies.
- Use customer managed policies instead of inline policies.
- Use access levels to review IAM permissions.
- Configure a strong password policy for your users.
- Enable MFA.
- Use roles for applications that run on Amazon EC2 instances.
- Do not share access keys.
- Rotate credentials regularly.
- Remove unnecessary credentials.
- Use policy conditions for extra security.
- Monitor activity in your AWS account.

VEERA

RDS

Create RDS : Amazon Relational Database Service

Select create database

The screenshot shows the Amazon RDS Dashboard. On the left, there's a sidebar with links like Dashboard, Databases, Query Editor, etc. The main area has a sidebar with 'DB Instance (0)', 'Recent events (0)', and 'Event subscriptions (0/20)'. Below that is a 'Create database' section with a note about setting up RDS in the cloud, a 'Restore from S3' button, and a prominent orange 'Create database' button. A note at the bottom says 'Note: your DB instances will launch in the US East (N. Virginia) region'. At the bottom is a 'Service health' section with a 'View service health dashboard' button, showing 'Current status' for 'Amazon Relational Database Service (N. Virginia)' which is 'Service is operating normally'.

Select standard create and choose MySQL

EC2

Standard create
You set all of the configuration options, including ones for availability, security, backups, and maintenance.

Easy create
Use recommended best-practice configurations. Some configuration options can be changed after the database is created.

Engine options

Engine type [Info](#)

Aurora (MySQL Compatible)

Aurora (PostgreSQL Compatible)

MySQL

MariaDB

PostgreSQL

Oracle

VEERA

Choose free tier model

Templates
Choose a sample template to meet your use case.

- Production**
Use defaults for high availability and fast, consistent performance.
- Dev/Test**
This instance is intended for development use outside of a production environment.
- Free tier**
Use RDS Free Tier to develop new applications, test existing applications, or gain hands-on experience with Amazon RDS.
[Info](#)

Provide password

Amazon RDS

Dashboard
Databases
Query Editor
Performance Insights
Snapshots
Exports in Amazon S3
Automated backups
Reserved instances
Proxies

Subnet groups
Parameter groups
Option groups
Custom engine versions
Zero-ETL integrations [New](#)

Events
Event subscriptions

Databases (1)

| DB identifier | Status | Role | Engine | Region & AZ | Size | Recommendations | CPU | Current activity |
|---------------|----------|----------|-----------------|-------------|-------------|-----------------|-----|------------------|
| database-1 | Creating | Instance | MySQL Community | us-east-1a | db.t3.micro | - | - | - |

if you want to see your credentials, you can click on view credentials details

Creating database database-1
Your database might take a few minutes to launch. You can use settings from database-1 to simplify configuration of [suggested database add-ons](#) while we finish creating your DB for you.

[View credential details](#)

Connection details to your database database-1

X

This is the only time you can view this password. Copy and save the password for your reference. If you lose the password, you must modify your database to change it. You can use a SQL client application or utility to connect to your database.

[Learn about connecting to your database](#)

Master username

admin

Master password

admin123 [Copy](#)

Endpoint

database-1.cbcwa4kygaeq.us-east-1.rds.amazonaws.com [Copy](#)

[Close](#)

Master username [Info](#)

Type a login ID for the master user of your DB instance.

admin

1 to 16 alphanumeric characters. The first character must be a letter.

Credentials management

You can use AWS Secrets Manager or manage your master user credentials.

Managed in AWS Secrets Manager - *most secure*

RDS generates a password for you and manages it throughout its lifecycle using AWS Secrets Manager.

Self managed

Create your own password or have RDS create a password that you manage.

Auto generate password

Amazon RDS can generate a password for you, or you can specify your own password.

Master password [Info](#)

Minimum constraints: At least 8 printable ASCII characters. Can't contain any of the following symbols: / ' " @

Confirm master password [Info](#)

DB subnet group [Info](#)

Choose the DB subnet group. The DB subnet group defines which subnets and IP ranges the DB instance can use in the VPC that you selected.

default

Public access [Info](#)

Yes

RDS assigns a public IP address to the database. Amazon EC2 instances and other resources outside of the VPC can connect to your database. Resources inside the VPC can also connect to the database. Choose one or more VPC security groups that specify which resources can connect to the database.

No

RDS doesn't assign a public IP address to the database. Only Amazon EC2 instances and other resources inside the VPC can connect to your database. Choose one or more VPC security groups that specify which resources can connect to the database.

VPC security group (firewall) [Info](#)

Choose one or more VPC security groups to allow access to your database. Make sure that the security group rules allow the appropriate incoming traffic.

Click on create database button

The screenshot shows the AWS RDS Databases page. At the top, there is a green success message: "Successfully created database database-1" with a "View connection details" button and a close icon. Below this, there is a blue info message about Aurora I/O-Optimized. The main table lists a single database entry:

| DB identifier | Status | Role | Engine | Region & AZ | Size | Recommendations | CPU | Current a |
|---------------|-----------|----------|-----------------|-------------|-------------|-----------------|--------|-----------|
| database-1 | Available | Instance | MySQL Community | us-east-1a | db.t3.micro | | 18.40% | |

A red box highlights the "Available" status of the database row.

Database would take few mins to up and running.

Click on the Database-1 to capture the endpoint.

The screenshot shows the Amazon RDS console interface. At the top, a green banner indicates "Successfully created database database-1" and provides a link to "View connection details". Below this, a blue banner introduces Aurora I/O-Optimized, stating it offers predictable pricing and up to 40% cost savings for I/O-intensive applications. The main area is titled "Databases (1)" and shows a table with one row for "database-1". The table columns include DB identifier, Status, Role, Engine, Region & AZ, Size, Recommendations, CPU, and Current activity. A red box highlights the "DB identifier" column for "database-1". Below the table, the "Summary" tab is selected in the navigation bar, which includes tabs for Connectivity & security, Monitoring, Logs & events, Configuration, Zero-ETL integrations, Maintenance & backups, Tags, and Recommendations. The "Connectivity & security" section displays the Endpoint (database-1.cbcwa4kygaeq.us-east-1.rds.amazonaws.com), Port (3306), Networking (Availability Zone: us-east-1a, VPC: vpc-0b713227c5d8798dd), and Security (VPC security groups: default (sg-03f126d4a0dac83c), Active). A large red arrow points from the highlighted "DB identifier" in the first screenshot down to the "Endpoint" field in the second screenshot.

Copy the Endpoint

Endpoint: database-1.cbcwa4kygaeq.us-east-1.rds.amazonaws.com

This endpoint we are using it to connect from MySQL client.

MySQL :: MySQL Workbench download

<https://www.mysql.com/products/workbench/>

← → C https://www.mysql.com/products/workbench/ 🔍

IVYSQL MYSQL.COM DOWNLOADS DOCUMENTATION DEVELOPER ZONE

Products Cloud Services Partners Customers Why MySQL? News & Events How to Buy

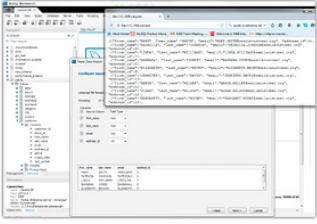
> MySQL HeatWave
▼ MySQL Enterprise Edition

- Datasheet (PDF)
- Technical Specification
- MySQL Database
- MySQL Document Store
- Oracle Enterprise Manager

MySQL Workbench

Enhanced Data Migration

[Download Now »](#)



MySQL Workbench is a unified visual tool for database architects, developers, and DBAs. MySQL Workbench prc

Click on Download Now

MySQL Workbench 8.0.36

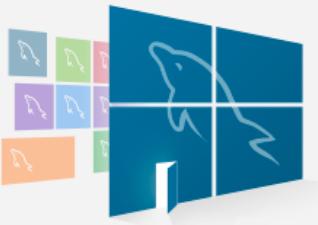
Select Operating System:

Microsoft Windows

Recommended Download:

MySQL Installer for Windows

All MySQL Products. For All Windows Platforms. In One Package.



Starting with MySQL 5.6 the MySQL Installer package replaces the standalone MSI packages.

[Windows \(x86, 32 & 64-bit\), MySQL Installer MSI](#) [Go to Download Page >](#)

Other Downloads:

| | | | |
|---|--------|-------|--------------------------|
| Windows (x86, 64-bit), MSI Installer (mysql-workbench-community-8.0.36-winx64.msi) | 8.0.36 | 42.0M | Download |
|---|--------|-------|--------------------------|

MD5: 2156fe0cb6f5ed83908e4636ba86390a | [Signature](#)

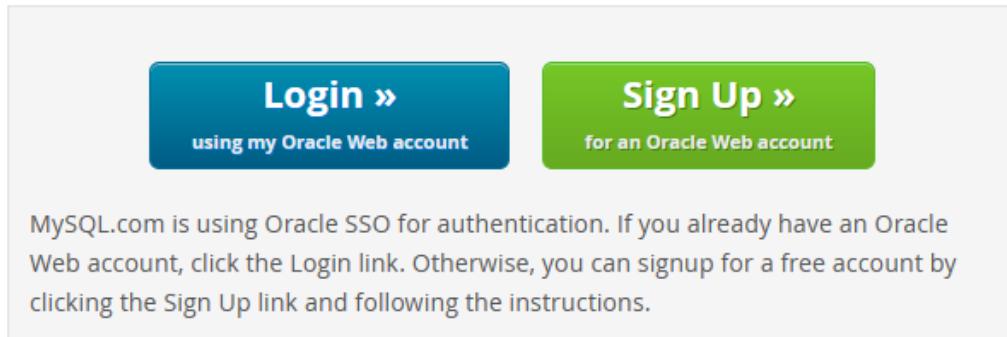
 We suggest that you use the [MD5 checksums and GnuPG signatures](#) to verify the integrity of the

Click on Download button

[Login Now](#) or [Sign Up](#) for a free account.

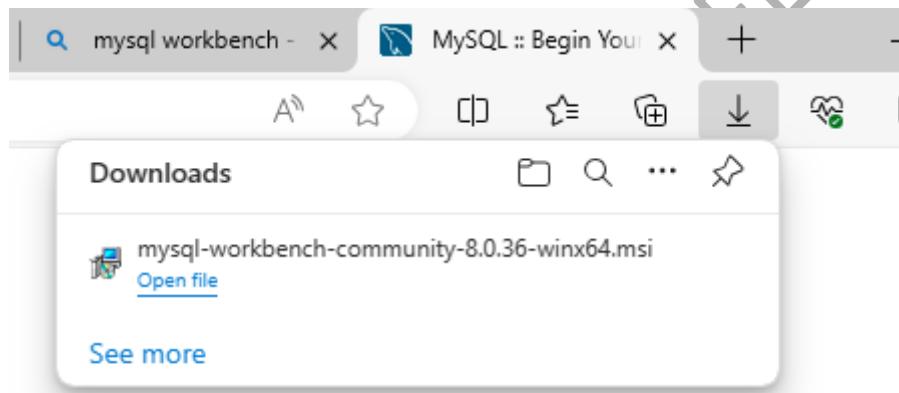
An Oracle Web Account provides you with the following advantages:

- Fast access to MySQL software downloads
- Download technical White Papers and Presentations
- Post messages in the MySQL Discussion Forums
- Report and track bugs in the MySQL bug system

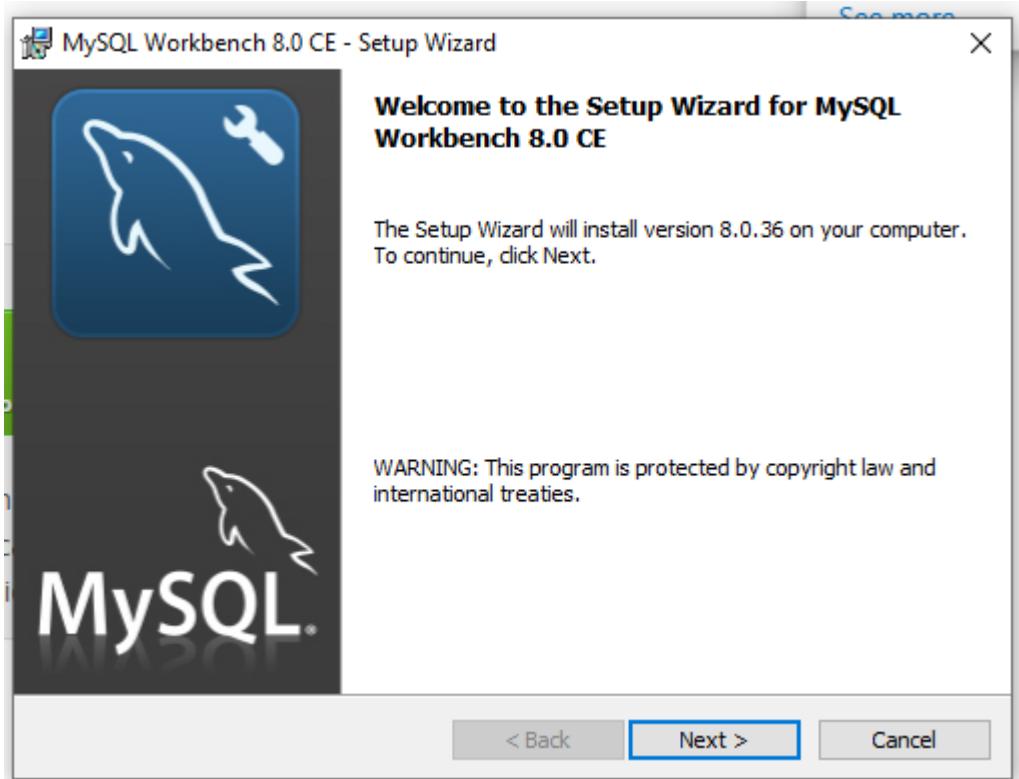


[No thanks, just start my download.](#)

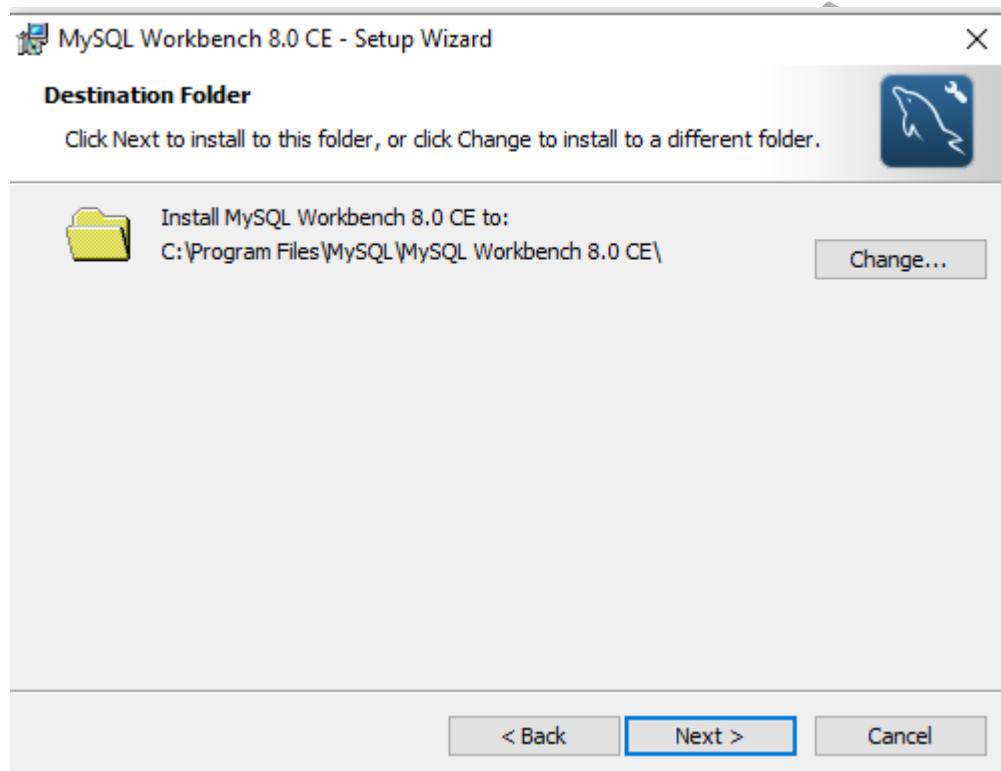
Click here to download with login : [No thanks, just start my download.](#)



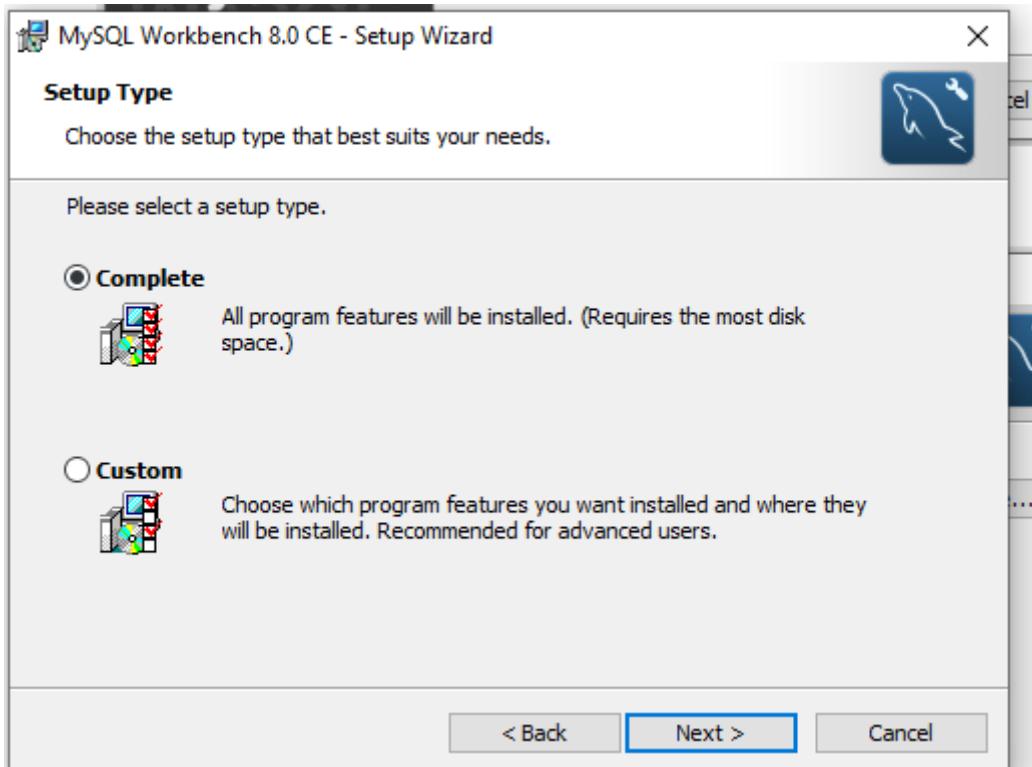
The download mysql client will be available in downloads, then double click on .msi file



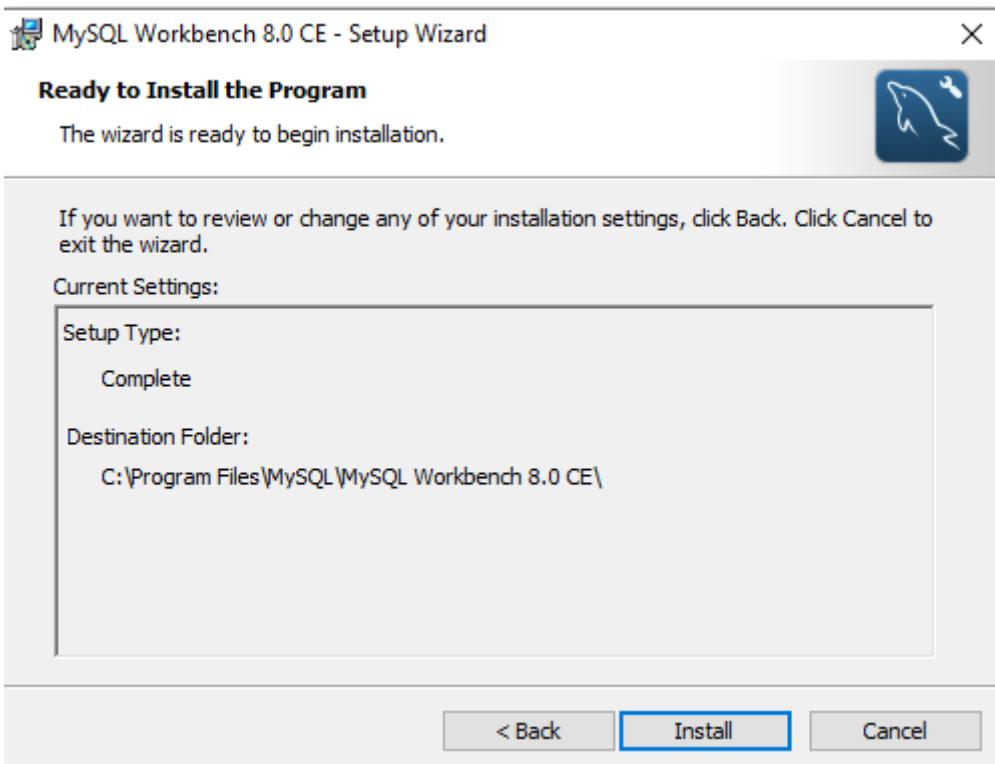
Click on Next



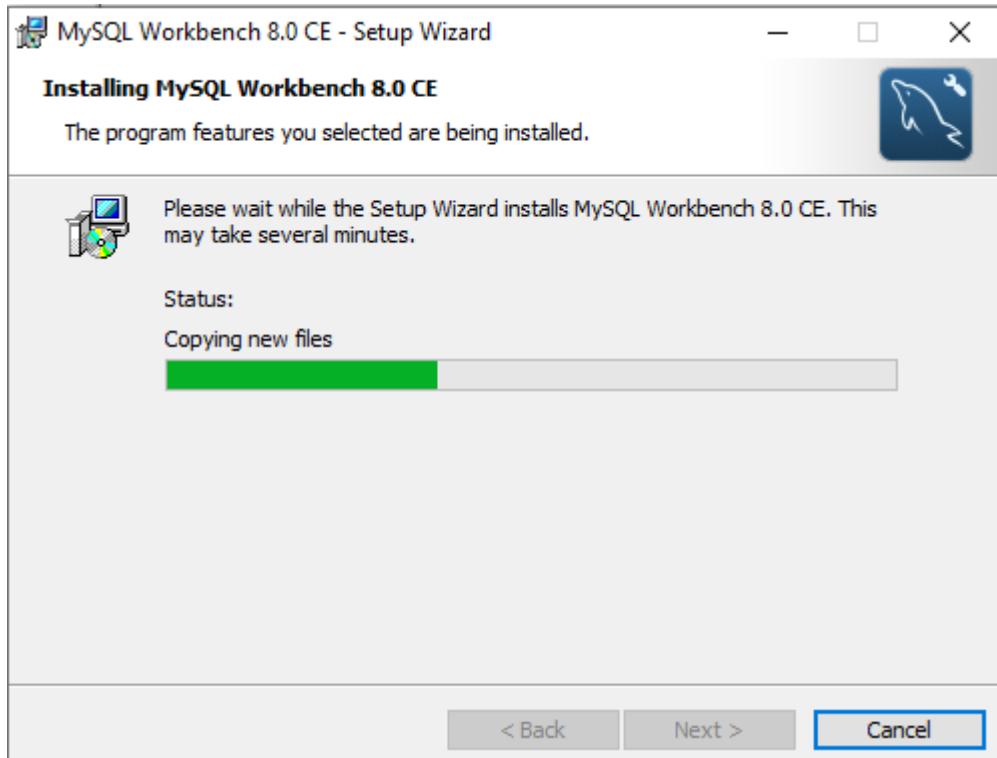
Click on Next



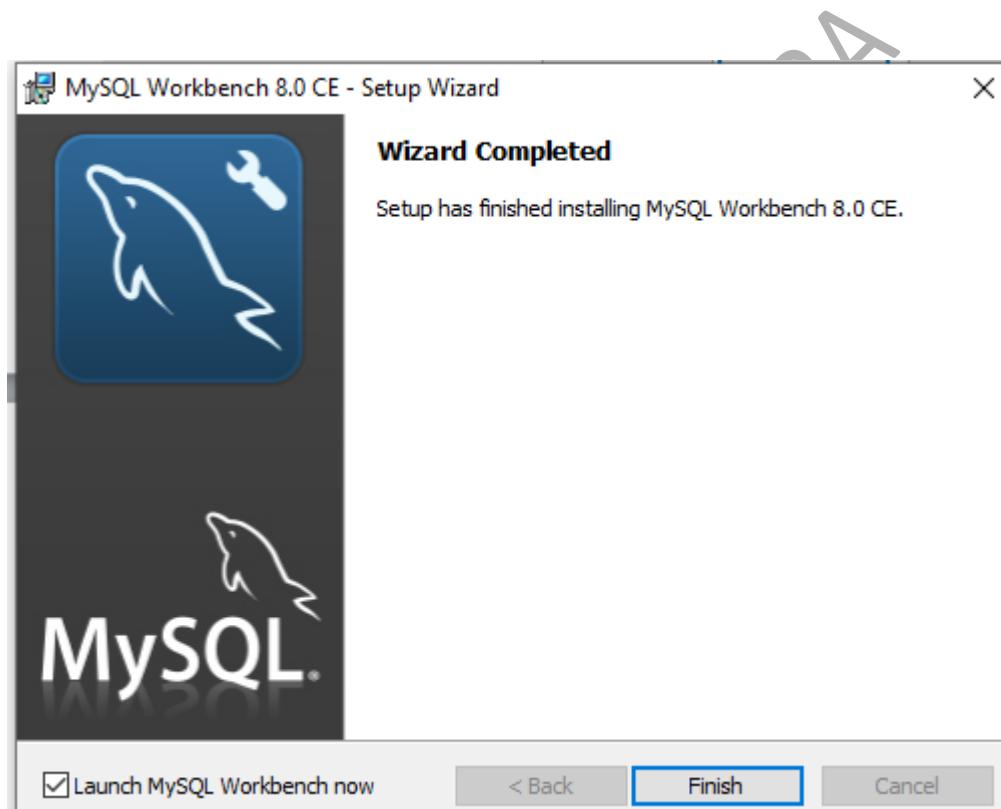
Click on Next



Click on Install

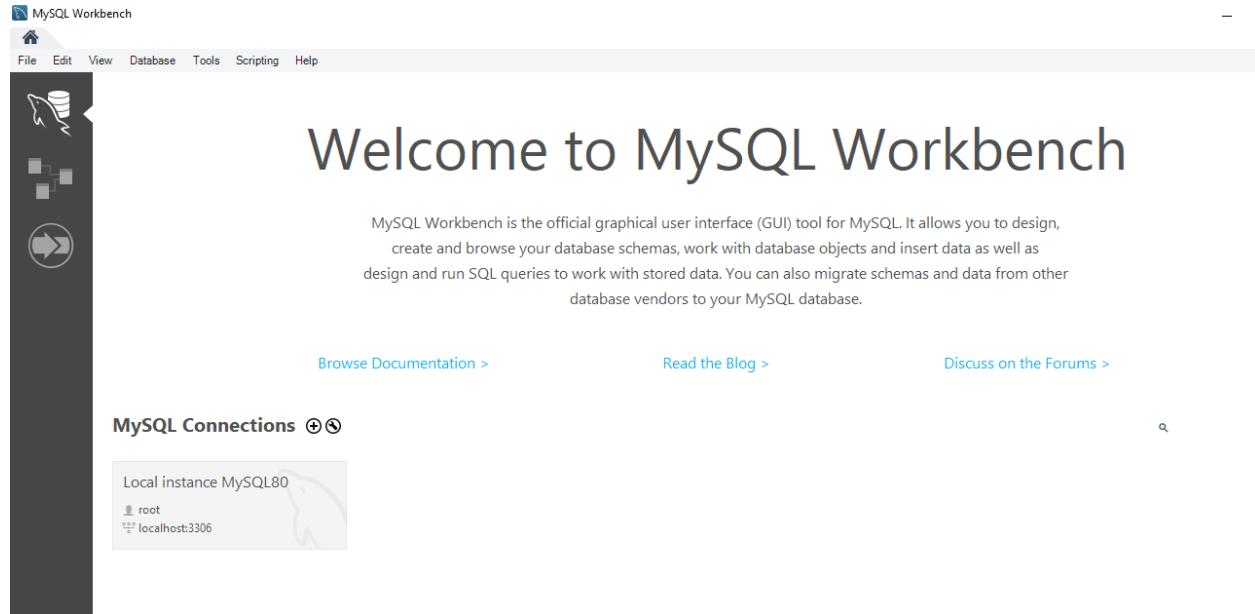


Wait till Finish button shows



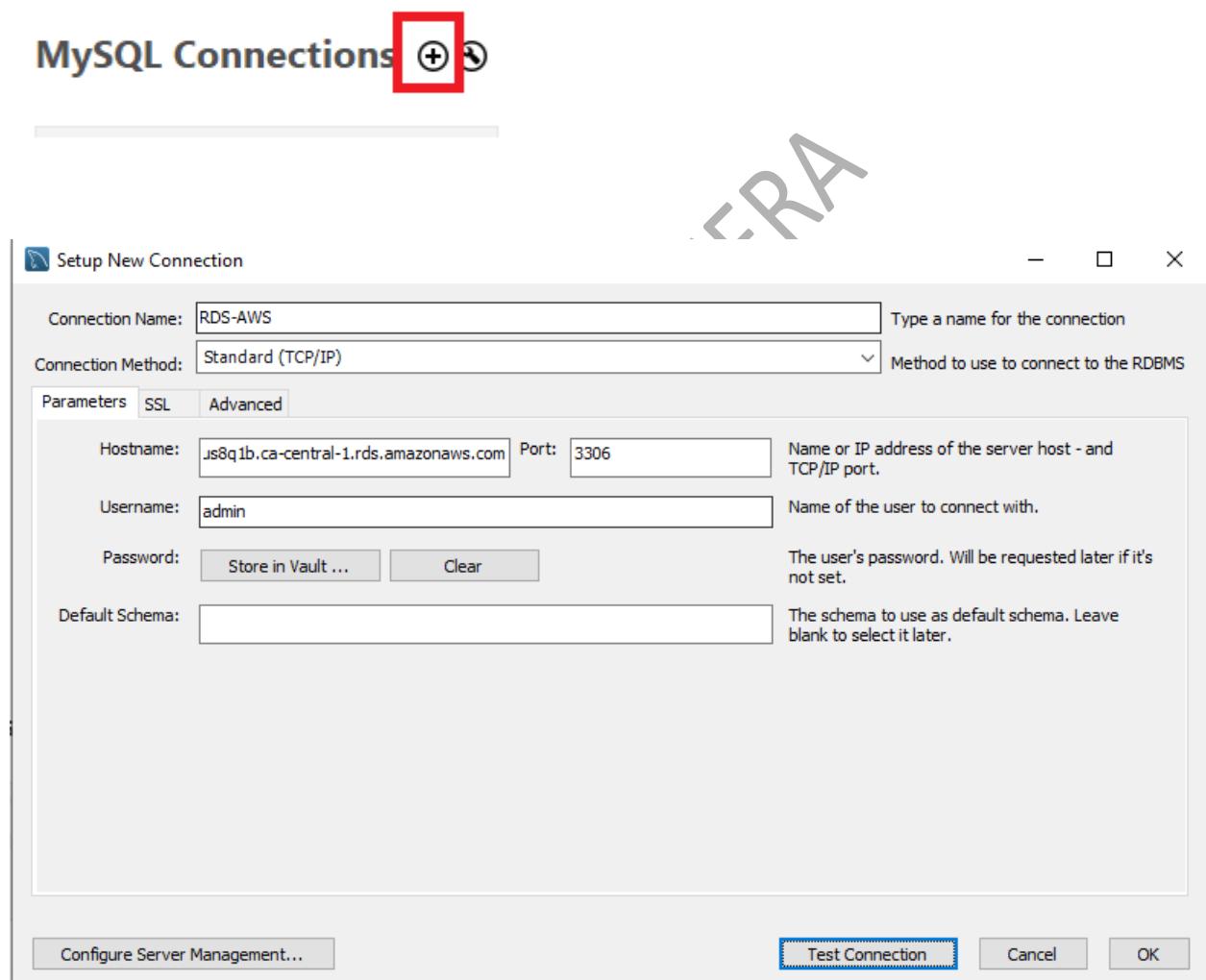
Click on Finish

Now launch the mysql workbench to connect AWS RDS



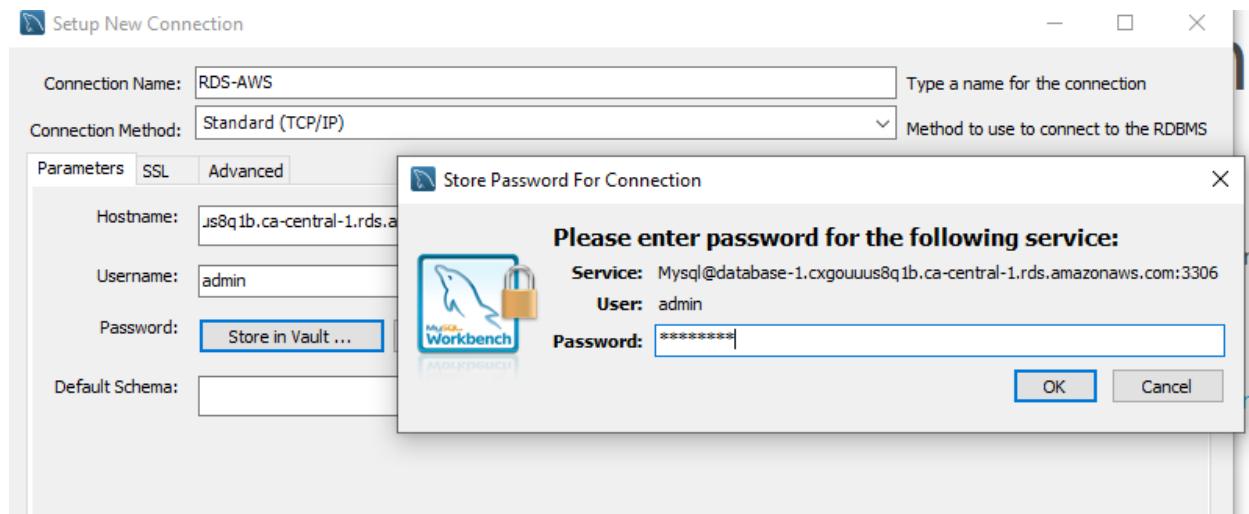
The screenshot shows the MySQL Workbench interface. At the top, there's a navigation bar with icons for Home, File, Edit, View, Database, Tools, Scripting, and Help. Below the navigation bar is a sidebar with three circular icons: a magnifying glass, a gear, and a double arrow. The main content area features a large title "Welcome to MySQL Workbench". Below the title is a brief description of what MySQL Workbench is used for. At the bottom of the main area are three links: "Browse Documentation >", "Read the Blog >", and "Discuss on the Forums >".

Click on MySQL connections + symbol to configure our RDS mysql db

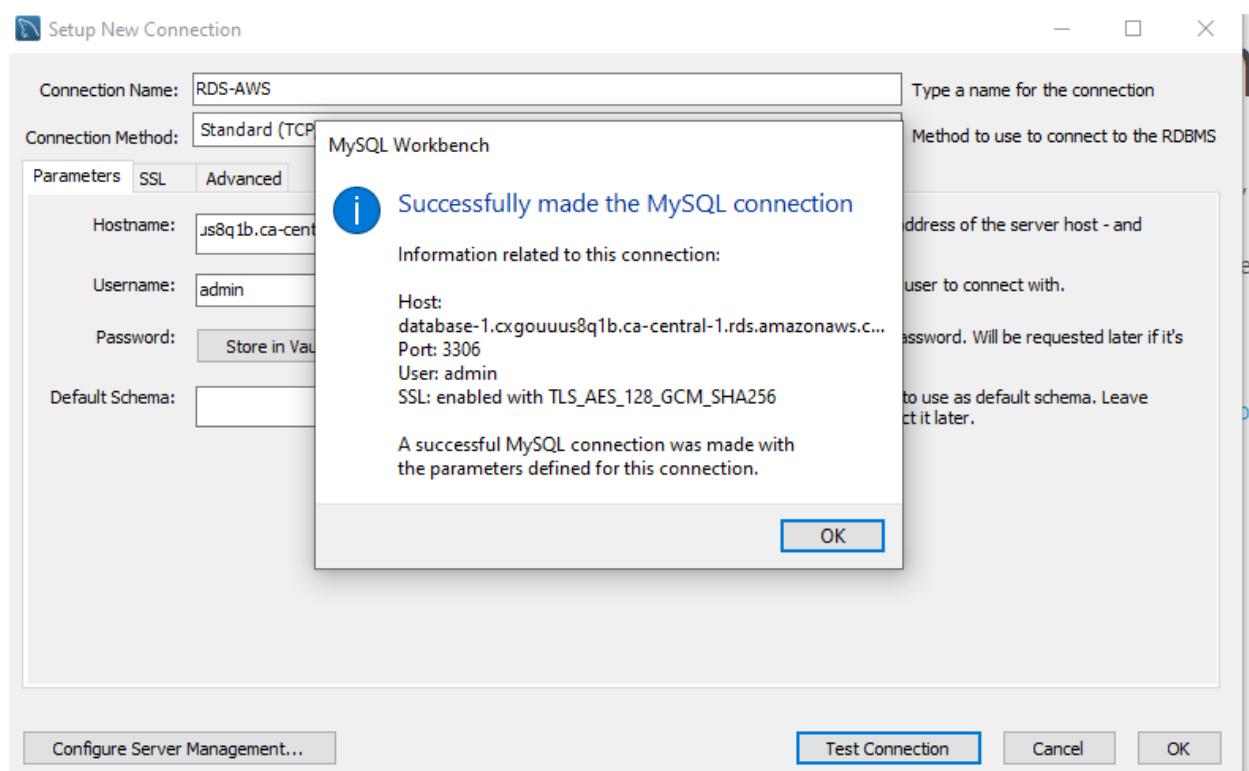


The screenshot shows the "Setup New Connection" dialog box. At the top, it says "MySQL Connections" followed by a red-bordered "+" button. The main area is titled "Setup New Connection". It has fields for "Connection Name" (set to "RDS-AWS") and "Connection Method" (set to "Standard (TCP/IP)"). Below these are tabs for "Parameters", "SSL", and "Advanced". Under "Parameters", there are fields for "Hostname" (set to "js8q1b.ca-central-1.rds.amazonaws.com"), "Port" (set to "3306"), "Username" (set to "admin"), and "Password". A note next to the password field says "The user's password. Will be requested later if it's not set." Under "Advanced", there is a field for "Default Schema". At the bottom of the dialog are buttons for "Configure Server Management...", "Test Connection" (which is highlighted with a blue border), "Cancel", and "OK".

Click on the Store in Vault to enter the password



Click on OK and then click on Test connection.



Connection got successful

MySQL Workbench

File Edit View Database Tools Scripting Help

Welcome to MySQL !

MySQL Workbench is the official graphical user interface (GUI) tool for MySQL. It allows you to design, create and browse your database schemas, work with database objects and insert data as well as design and run SQL queries to work with stored data. You can also migrate schemas and data from other database vendors to your MySQL database.

[Browse Documentation >](#) [Read the Blog >](#)

MySQL Connections [+ !\[\]\(879ae222611bb5a82937dbe6c3e11dba_img.jpg\)](#)

Local instance MySQL80

 root

 localhost:3306

RDS-AWS

 admin

 database-1.cxgouuuus8q1b.ca-central...

Click on our RDS-AWS db to create database and tables.

Welcome to MySQL Workben

MySQL Workbench is the official graphical user interface (GUI) tool for MySQL. It allows you to design, create and browse your database schemas, work with database objects and insert data as well as design and run SQL queries to work with stored data. You can also migrate schemas and data from other database vendors to your MySQL database.

Opening SQL Editor

An SQL editor instance for 'RDS-AWS' is opening and should be available in a moment.

[Read the Blog >](#)

Please stand by...

[Discuss on the Forum](#)

MySQL Connections [+](#)

[Browse Documentation >](#)

[Cancel](#)

Local instance MySQL80

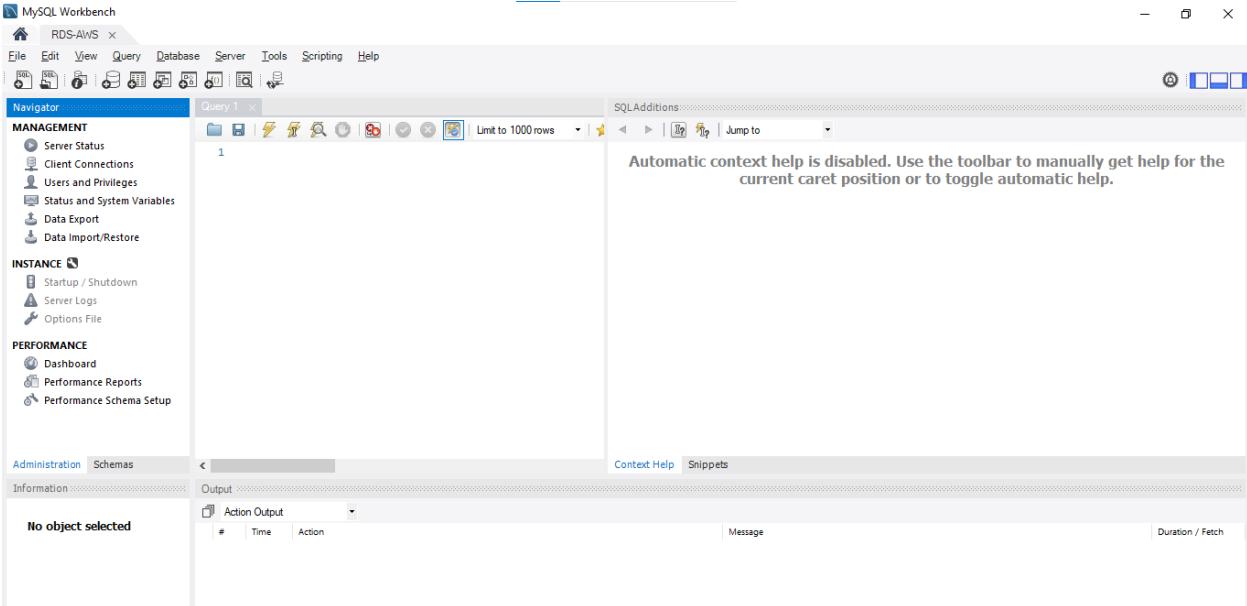
 root

 localhost:3306

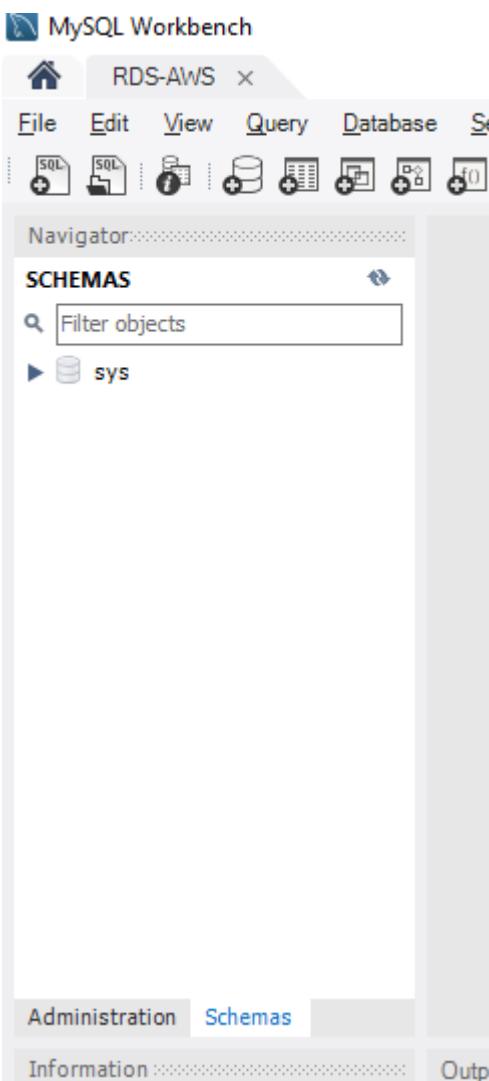
RDS-AWS

 admin

 database-1.cxgouuuus8q1b.ca-central...

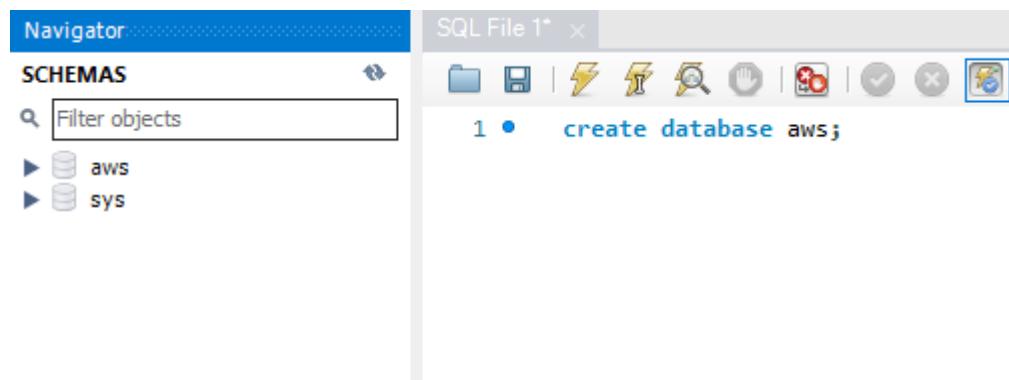
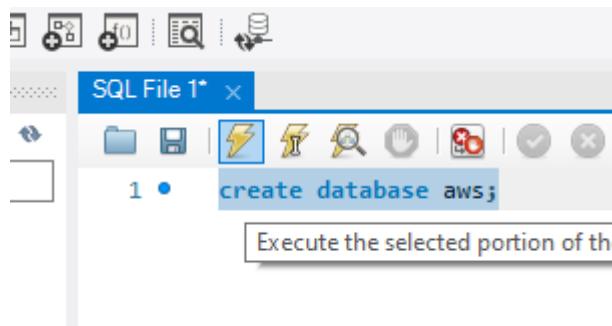


Click on Schemas to see the databases.



Create Database, Table and rows

Click on the execute query button



Aws db created

Run the below query to create a table in aws db

create table aws.employee (empid int auto_increment primary key,
first_name varchar(30),
last_name varchar(30),
email varchar(50),
department varchar(40),
salary decimal(10,2),
hire_date date
);

The screenshot shows the SQL Workbench interface. In the top panel, under the 'SCHEMAS' section, the 'aws' schema is selected. A SQL file is open with the following code:

```
1 • create database aws;
2
3 • create table aws.employee (empid int auto_increment primary key,
4                               first_name varchar(30),
5                               last_name varchar(30),
6                               email varchar(50),
7                               department varchar(40),
8                               salary decimal(10,2),
9                               hire_date date
10);
```

The screenshot shows the execution results in the 'Output' tab. It displays a single log entry:

| Action Output |
|--|
| # Time Action 1 00:15:00 create table aws.employee (empid int auto_increment primary key, first_name var... 0 row(s) affected |

The screenshot shows the Navigator pane with the 'SCHEMAS' section expanded. Under the 'aws' schema, the 'Tables' section is selected, and the 'employee' table is highlighted with a red box.

Employee table created

Run the below 2 query to insert the data and

```
insert into aws.employee (first_name, last_name, email, department,
salary, hire_date)
values ('David', 'John','david@12345.com','Engineer',50000.00,'2024-05-02');
```

```
insert into aws.employee (first_name, last_name, email, department,
salary, hire_date)
values ('Mark', 'Anthony','mark@45678.com','QA',60000.00,'2023-03-05');
```

Run the below query to view the data

```
select * from aws.employee;
```

The screenshot shows the MySQL Workbench interface. On the left, the 'SCHEMAS' tree shows 'aws' selected, with 'Tables' expanded to show 'employee'. The main pane displays the SQL code for creating the table and inserting two rows of data. The 'Result Grid' pane below shows the resulting data in a table format.

```
3 • ④ create table aws.employee (empid int auto_increment primary key,
4                               first_name varchar(30),
5                               last_name varchar(30),
6                               email varchar(50),
7                               department varchar(40),
8                               salary decimal(10,2),
9                               hire_date date
10                             );
11
12 •  insert into aws.employee (first_name, last_name, email, department, salary, hire_date)
13   values ('David', 'John', 'david@12345.com', 'Engineer', 50000.00, '2024-05-02');
14
15 •  insert into aws.employee (first_name, last_name, email, department, salary, hire_date)
16   values ('Mark', 'Anthony', 'mark@45678.com', 'QA', 60000.00, '2023-03-05');
17
18 •  select * from aws.employee;
```

| empid | first_name | last_name | email | department | salary | hire_date |
|-------|------------|-----------|-----------------|------------|----------|------------|
| 1 | David | John | david@12345.com | Engineer | 50000.00 | 2024-05-02 |
| 2 | Mark | Anthony | mark@45678.com | QA | 60000.00 | 2023-03-05 |
| * | HULL | HULL | HULL | HULL | HULL | HULL |

Now we able to see the data in employee table

Create Read Replica from Master Data Base:

Select database and click on read replica

The screenshot shows the AWS RDS console under the 'Databases' section. A context menu is open over the 'database-1' row, with 'Create read replica' highlighted. The main table lists the database details, including its status as 'Available'.

| DB identifier | Status | Role | Engine | Region & AZ | Size | Recommendations | CPU | Current ac |
|---------------|-----------|----------|-----------------|---------------|-------------|-----------------|-------|------------|
| database-1 | Available | Instance | MySQL Community | ca-central-1a | db.t3.micro | | 3.95% | 2 |

Create read replica

You are creating a replica DB instance from a source DB instance. This new DB instance will have the source DB instance's DB security groups and DB parameter groups.

Settings

Replica source

Source DB instance identifier

database-1
 Role: Instance

DB instance identifier

This is the unique key that identifies a DB instance. This parameter is stored as a lowercase string (for example, mydbinstance).

read-replica-1

Provide read-replica-1 and the click on create read replica button

The screenshot shows the AWS RDS Databases page with the following details:

- Databases (2)**: Shows two entries: "database-1" and "read-replica-1".
- Status**: "database-1" is "Modifying" and "read-replica-1" is "Creating".
- Role**: "database-1" is Primary and "read-replica-1" is Replica.
- Engine**: Both are MySQL Community.
- Region & AZ**: Both are in "ca-central-1a".
- Size**: Both are "db.t3.micro".
- Recommendations**: Both show low usage (e.g., 2.90% for database-1).
- CPU**: Both show low usage (e.g., 3.32% for database-1).
- Actions**: Includes "Group resources", "Modify", "Actions", "Restore from S3", and "Create database".
- Filter by databases**: A search bar at the top left.
- Page Navigation**: "1" of 1 page.

Lets wait for up and running

The screenshot shows the AWS RDS Databases page with the following details:

- Databases (2)**: Shows "database-1" and "read-replica-1".
- Status**: Both are now "Available".
- Role**: "database-1" is Primary and "read-replica-1" is Replica.
- Engine**: Both are MySQL Community.
- Region & AZ**: Both are in "ca-central-1a".
- Size**: Both are "db.t3.micro".
- Recommendations**: Both show low usage (e.g., 3.32% for database-1).
- CPU**: Both show low usage (e.g., 3.69% for read-replica-1).
- Actions**: Includes "Group resources", "Modify", "Actions", "Restore from S3", and "Create database".
- Filter by databases**: A search bar at the top left.
- Page Navigation**: "1" of 1 page.

Both replica and Primary db available

Open read-replica-1 to copy the endpoint to connect through mysql client.

Screenshot of the AWS RDS console showing the 'Connectivity & security' tab for a database instance.

Endpoint & port

- Endpoint copied: read-replica-1.cxgouuuus8q1b.ca-central-1.rds.amazonaws.com
- Port: 3306

Networking

- Availability Zone: ca-central-1d
- VPC: vpc-0ac628c6c2a97cc25
- Subnet group: default-vpc-0ac628c6c2a97cc25
- Subnets: subnet-0abb4c2761a8a3c61

Security

- VPC security groups: default (sg-0a3a4062d881ebd3a) - Active
- Publicly accessible: Yes
- Certificate authority: rds-ca-rsa2048-g1
- Certificate authority date: May 21, 2021, 16:02 UTC 07:00

Setup New Connection

Connection Name: AWS_RDS_Read_Replica-1 Type a name for the connection

Connection Method: Standard (TCP/IP) Method to use to connect to the RDBMS

Parameters Advanced

Hostname: us8q1b.ca-central-1.rds.amazonaws.com

Username: admin

Password:

Default Schema:

Store Password For Connection

Please enter password for the following service:

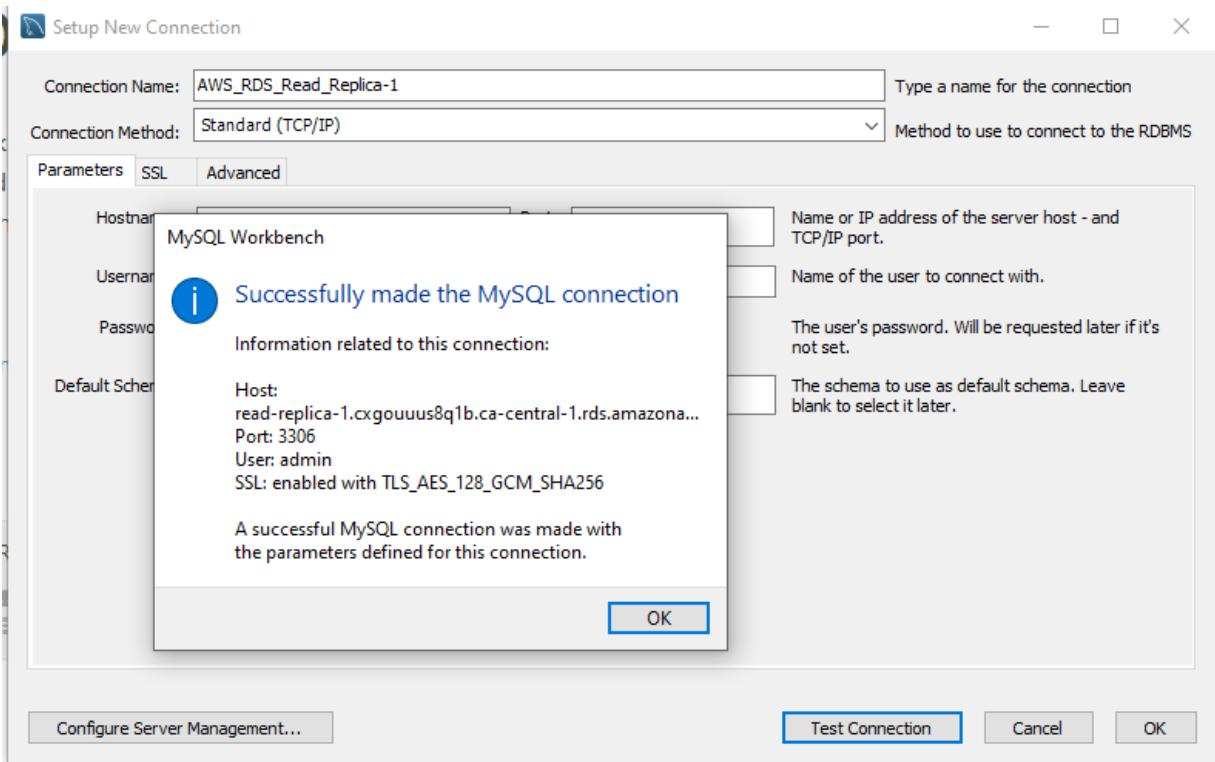
Service: Mysql@read-replica-1.cxgouuuus8q1b.ca-central-1.rds.amazonaws.com:3306

User: admin

Password: *****

OK Cancel

Configure Server Management... Test Connection Cancel OK



Welcome to MySQL Workbench

MySQL Workbench is the official graphical user interface (GUI) tool for MySQL. It allows you to design, create and browse your database schemas, work with database objects and insert data as well as design and run SQL queries to work with stored data. You can also migrate schemas and data from other

MySQL Connections + (0)

| Local instance MySQL80 | RDS-AWS | AWS_RDS_Read_Replica-1 |
|------------------------|---|---|
| root localhost:3306 | admin database-1.cxgouuuus8q1b.ca-central-1.rds.amazonaws.com:3306 | admin read-replica-1.cxgouuuus8q1b.ca-central-1.rds.amazonaws.com:3306 |

All table info available in read-replica-1 without creating any data.

Administration Schemas

Information

Table: employee

Columns:

| | empid | first_name | last_name | email | department | salary | hire_date |
|---|-------|------------|-----------|-----------------|------------|----------|------------|
| 1 | 1 | David | John | david@12345.com | Engineer | 50000.00 | 2024-05-02 |
| 2 | 2 | Mark | Anthony | mark@45678.com | QA | 60000.00 | 2023-03-05 |
| * | NULL | NULL | NULL | NULL | NULL | NULL | NULL |

Create 3rd record in Master db

```
insert into aws.employee (first_name, last_name, email, department, salary, hire_date)
values ('Honey', 'Mony','honey@88765.com','QA',66000.00,'2022-08-22');
```

Now try to create 4th record in read-replica-1, lets see..

The screenshot shows the MySQL Workbench interface. In the top-left, there's a 'Navigator' pane with 'SCHEMAS' expanded, showing 'aws' with 'Tables', 'Views', 'Stored Procedures', and 'Functions'. Below it is an 'Information' pane showing 'Table: employee' with its columns: empid (int AI PK), first_name (varchar(30)), last_name (varchar(30)), email (varchar(50)), department (varchar(40)), salary (decimal(10,2)), and hire_date (date). The main area is a 'Query 1' editor with the following SQL code:

```

1 • select * from aws.employee;
2
3 • insert into aws.employee (first_name, last_name, email, department, salary, hire_date)
4   values ('sony', 'clara', 'sony@88765.com', 'OA', 45000.00, '2022-07-18');
5
6
7
8
9

```

To the right of the editor is a 'SQLAdditions' toolbar with various icons. A message box is displayed: "Automatic context help is disabled. Use the toolbar to manually get help for the current caret position or to toggle automatic help." The bottom-right pane is an 'Output' window titled 'Action Output' showing three log entries:

| # | Time | Action | Message | Duration / Fetch |
|---|----------|--|---|-----------------------|
| 1 | 00:58:17 | select * from aws.employee LIMIT 0, 1000 | 2 row(s) returned | 0.078 sec / 0.016 sec |
| 2 | 01:04:55 | select * from aws.employee LIMIT 0, 1000 | 3 row(s) returned | 0.078 sec / 0.000 sec |
| 3 | 01:08:13 | insert into aws.employee (first_name, last_name, email, department, salary, hire_date) | Error Code: 1290. The MySQL server is running with the --read-only option so it cannot execute this statement | 0.078 sec |

A tooltip for the third entry states: "Error Code: 1290. The MySQL server is running with the --read-only option so it cannot execute this statement".

Error:

| | |
|---|---|
| 2 row(s) returned | 0.078 sec / 0.016 sec |
| 3 row(s) returned | 0.078 sec / 0.000 sec |
| Error Code: 1290. The MySQL server is running with the --read-only option so it cannot execute this statement | Error Code: 1290. The MySQL server is running with the --read-only option so it cannot execute this statement |



By seeing this error we have no access to execute insert, delete queries in read-replica-1

Try to delete a row from read-replica-1, lets see

The screenshot shows the MySQL Workbench interface. In the top-left, the connection is set to 'RDS-AWS' and the schema is 'AWS_RDS_Read_Replica-1'. The 'Query 1' tab contains the following SQL code:

```

1 • select * from aws.employee;
2
3 • insert into aws.employee (first_name, last_name, email, department, salary, hire_date)
4   values ('sony', 'clara', 'sony@88765.com', 'QA', 45000.00, '2022-07-18');
5
6
7 • delete from aws.employee where empid=1;
8
9
10
11

```

In the bottom-right pane, there is an error message: "Error Code: 1290. The MySQL server is running with the --read-only option so it cannot execute this statement".

Delete query also not having the permission in read-replica-1 db.

The screenshot shows the MySQL Workbench interface. In the top-left, the connection is set to 'RDS-AWS' and the schema is 'AWS_RDS_Read_Replica-1'. The 'Query 1' tab contains the following SQL code:

```

1 • select * from aws.employee;
2
3 • insert into aws.employee (first_name, last_name, email, department, salary, hire_date)
4   values ('sony', 'clara', 'sony@88765.com', 'QA', 45000.00, '2022-07-18');
5
6
7
8
9

```

In the bottom-right pane, there is an error message: "Error Code: 1290. The MySQL server is running with the --read-only option so it cannot execute this statement".

RDS-AWS x AWS_RDS_Read_Replica-1 x

File Edit View Query Database Server Tools Scripting Help

Navigator

SCHEMAS

Filter objects

aws

Tables

employee

Columns

Indexes

Foreign Keys

Triggers

Views

Stored Procedures

Functions

sys

SQL File 1*

Limit to 1000 rows

```
5     last_name varchar(30),
6     email varchar(50),
7     department varchar(40),
8     salary decimal(10,2),
9     hire_date date
10    );
11
12 • insert into aws.employee (first_name, last_name, email, department, salary, hire_date)
13   values ('David', 'John', 'david@12345.com', 'Engineer', 50000.00, '2024-05-02');
14
15 • insert into aws.employee (first_name, last_name, email, department, salary, hire_date)
16   values ('Mark', 'Anthony', 'mark@45678.com', 'QA', 60000.00, '2023-03-05');
17
18 • insert into aws.employee (first_name, last_name, email, department, salary, hire_date)
19   values ('Honey', 'Mony', 'honey@88765.com', 'QA', 66000.00, '2022-08-22');
20
21 • select * from aws.employee;
```

Administration Schemas

Information

Table: employee

Columns:

| empid | int AI PK |
|------------|-------------|
| first_name | varchar(30) |
| last_name | varchar(30) |
| email | varchar(50) |
| department | varchar(40) |

Result Grid | Filter Rows: | Edit: | Export/Import: | Wrap Cell Content: | Result Grid | Form Editor

| empid | first_name | last_name | email | department | salary | hire_date |
|-------|------------|-----------|-----------------|------------|----------|------------|
| 1 | David | John | david@12345.com | Engineer | 50000.00 | 2024-05-02 |
| 2 | Mark | Anthony | mark@45678.com | QA | 60000.00 | 2023-03-05 |
| 3 | Honey | Mony | honey@88765.com | QA | 66000.00 | 2022-08-22 |
| * | HULL | HULL | HULL | HULL | HULL | HULL |

The master data immediately reflecting into read-replica-1 db

The screenshot shows the SSMS interface with the following details:

- Toolbar:** Includes icons for New Query, Open, Save, Print, and various tools.
- File Menu:** File, Edit, View, Query, Database, Server, Tools, Scripting, Help.
- Navigator:** Shows the database schema with the **SCHEMAS** node expanded, revealing the **aws** schema which contains the **Tables** node and the **employee** table.
- Query Window:** Titled "Query 1", contains the SQL command: `select * from aws.employee;`. The results show the following data:

| | empid | first_name | last_name | email | department | salary | hire_date |
|---|-------|------------|-----------|-----------------|------------|----------|------------|
| 1 | 1 | David | John | david@12345.com | Engineer | 50000.00 | 2024-05-02 |
| 2 | 2 | Mark | Anthony | mark@45678.com | QA | 60000.00 | 2023-03-05 |
| 3 | 3 | Honey | Mony | honey@88765.com | QA | 66000.00 | 2022-08-22 |
| * | NULL | NULL | NULL | NULL | NULL | NULL | NULL |

Overview of Amazon RDS (Relational Database Service)

Amazon RDS (Relational Database Service) is a fully managed relational database service provided by AWS. It allows you to run a relational database in the cloud without needing to worry about infrastructure management tasks such as hardware provisioning, database setup, patching, and backups. RDS makes it easy to set up, operate, and scale relational databases for web applications and other workloads.

Key Features of Amazon RDS:

1. **Managed Service:** Amazon handles the administrative tasks such as backups, patch management, and failover, reducing the operational overhead.
2. **Multiple Database Engines:** Supports popular relational database engines, including:
 - o MySQL
 - o PostgreSQL
 - o MariaDB
 - o Oracle
 - o SQL Server

3. Scalability:

- You can easily scale the compute (CPU, memory) and storage (storage size and throughput) resources for your database as your workload grows.

4. High Availability and Durability:

- Multi-AZ Deployments: Ensures high availability by replicating data synchronously to a secondary instance in another Availability Zone.
- Automated Backups: Amazon RDS supports point-in-time backups, allowing you to restore your database to any specific time.

5. Security:

- Integration with AWS Identity and Access Management (IAM) for fine-grained access control.
- Encryption at rest and in transit for data protection.

6. Automated Backups: RDS automatically backs up your database and retains backups for up to 35 days.

7. Monitoring: RDS integrates with Amazon CloudWatch to provide metrics on database performance (e.g., CPU utilization, disk I/O, query performance).

Amazon RDS Use Cases:

- Web Applications: Running the backend relational databases for web and mobile applications.
- Analytics and Data Warehousing: Storing structured data and running complex queries.
- Business Applications: Running enterprise applications that rely on relational databases, such as ERP and CRM systems.

Overview of RDS Read Replicas

An RDS Read Replica is a read-only copy of an RDS primary database instance that asynchronously replicates data from the primary database. It is designed to offload read-heavy database traffic from the primary instance, improving read scalability and reducing latency for read requests.

Key Features of RDS Read Replicas:

1. Asynchronous Replication:

- Data is asynchronously copied from the primary database to the read replica, meaning there may be a slight delay (lag) between the primary and the replica.

- This lag is usually small (seconds to a few minutes) but may vary depending on traffic and database size.

2. Improved Read Performance:

- Offload read queries to the replica instead of the primary database, improving overall application performance.
- Multiple read replicas can be created to handle large numbers of read queries in parallel.

3. Geographic Distribution:

- Cross-Region Replication: You can create read replicas in a different AWS region to serve global users with lower latency.
- This is useful for improving application performance for users far away from the primary database's region.

4. Scalability:

- Adding more read replicas can horizontally scale the read capacity of your database, allowing you to handle higher read traffic without affecting the primary instance.

5. Failover and Disaster Recovery:

- Although read replicas are typically used for scaling read traffic, in the case of a failure in the primary instance, you can promote a read replica to become the new primary database instance.
- This promotes fault tolerance and enhances availability, especially when used in conjunction with Multi-AZ deployments.

6. Backup and Maintenance:

- Unlike the primary instance, read replicas don't have automatic backups or maintenance windows, but you can manually back them up if needed.

7. Write Restrictions:

- Read replicas can only serve read queries. You cannot write data directly to the replica. If you need to update data, it must be done on the primary instance, which will asynchronously replicate the changes to the replica.

How RDS Read Replicas Work

1. Replication Process:

- Asynchronous: The replication process is asynchronous, meaning the primary instance and the read replica are not always perfectly in sync. This ensures low latency for write operations on the primary instance.
- Replication Mechanism: For engines like MySQL and PostgreSQL, Amazon RDS uses database-level replication to synchronize data between the primary database and the replica. Changes made to the primary instance are propagated to the read replica.

2. Creating Read Replicas:

- You create read replicas from an existing RDS primary instance. AWS offers the option to create a read replica through the RDS management console, CLI, or API.
- When creating a replica, you can choose different instance sizes and even place them in a different availability zone or region.

3. Read Traffic Distribution:

- Applications can direct read queries to the read replica(s) by using the endpoint provided for the replica.
- The primary database instance continues to handle write operations, while the read replicas handle read queries.

4. Promotion of Read Replicas:

- If needed, you can promote a read replica to become the primary database instance. This is often done in the case of disaster recovery or if you want to change the primary instance for performance reasons.

5. Replication Lag:

- There is a small delay between when a change occurs on the primary instance and when it is reflected on the replica (known as replication lag).
- Typically, this lag is very small, but it can vary based on the database load and network conditions.

Benefits of Using RDS Read Replicas

- Offload Read Traffic: Read replicas can handle read-heavy workloads, which reduces the load on the primary instance and helps scale applications.
- Improved Performance: By distributing read queries across multiple replicas, you can significantly improve the performance of your application, particularly in read-intensive applications like reporting, analytics, or caching.

- **High Availability and Disaster Recovery:** Read replicas can serve as a backup in case of primary instance failure. In an emergency, you can promote a read replica to become the new primary instance, which reduces downtime.
 - **Global Reach:** Read replicas can be placed in different AWS regions, providing faster read access to users globally.
-

Limitations of RDS Read Replicas

- **Read-Only:** Read replicas are read-only, which means you cannot perform write operations directly on them.
- **Replication Lag:** Since replication is asynchronous, there is always a slight delay between changes made in the primary instance and when they are reflected in the replica.
- **Not a Full Backup:** Read replicas are not intended to be used as backups for disaster recovery, though you can promote them to become a standalone database in case of failover.

Use Cases for RDS Read Replicas

1. **Scalable Web Applications:** When your application has a high volume of read queries (e.g., e-commerce sites, content management systems), you can use read replicas to offload read traffic and enhance performance.
2. **Analytics:** Use read replicas to run analytics and reporting queries without affecting the performance of the production database.
3. **Disaster Recovery:** Promote a read replica to the primary in case of failure or for planned migrations.
4. **Global Applications:** Deploy read replicas in different regions to reduce latency and improve performance for geographically distributed users.

AWS CloudWatch

AWS CloudWatch

Amazon CloudWatch is a powerful monitoring and observability service from AWS that provides real-time insights into the performance, health, and operational status of your AWS resources and applications. It helps you collect and track metrics, collect and monitor log files, set alarms, and automatically react to changes in your AWS environment.

CloudWatch is primarily used for:

- Monitoring cloud resources and applications
- Collecting and storing log files
- Creating alarms based on specific metrics
- Visualizing data through dashboards
- Automating responses to events and thresholds

Core Features of AWS CloudWatch

1. CloudWatch Metrics:

- Metrics are the fundamental data points that CloudWatch uses to monitor the performance of AWS resources. AWS services (like EC2, S3, RDS, Lambda, etc.) automatically send metrics to CloudWatch at regular intervals (typically every minute).
- These metrics can include information like:
 - CPU utilization, memory usage, disk I/O, network traffic for EC2 instances
 - Request count, latency, error rates for AWS services like ELB, API Gateway
 - DB connections, read/write latency, free storage space for Amazon RDS
- You can also publish custom metrics from your own applications, giving you flexibility to monitor any data that is important to your specific use case.

2. CloudWatch Alarms:

- CloudWatch Alarms allow you to define thresholds for specific metrics. When these thresholds are breached (e.g., CPU utilization exceeds 80%), CloudWatch can trigger automated actions such as:
 - Sending notifications via Amazon SNS (Simple Notification Service)
 - Executing AWS Lambda functions
 - Auto-scaling actions (e.g., scaling EC2 instances up/down)
- You can configure alarms based on any metric, including custom metrics.

3. CloudWatch Logs:

- CloudWatch Logs helps you collect and monitor log data from AWS services, EC2 instances, Lambda functions, and other sources.
- Log data from your applications, system events, or security logs can be pushed to CloudWatch for centralized storage and real-time analysis.
- Features of CloudWatch Logs include:
 - Log Groups and Log Streams: Group and organize logs from different sources.
 - Log Retention: Automatically delete old logs after a specified retention period.
 - Metric Filters: Create CloudWatch metrics from specific patterns in your logs (e.g., error codes, application events).
 - Real-time Log Monitoring: Detect specific events or errors in your logs and set alarms for specific log patterns.

4. CloudWatch Dashboards:

- CloudWatch Dashboards provide a customizable, graphical interface for monitoring and visualizing your metrics and logs. You can create dashboards to display multiple metrics from different AWS services and even custom metrics.
- Dashboards can be shared with other team members, and you can view data in real time, track historical trends, and gain operational insights.
- You can visualize metrics like CPU utilization, application response time, or any custom metric across multiple AWS resources.

5. CloudWatch Events (Now part of Amazon EventBridge):

- CloudWatch Events enables you to respond to changes in your AWS environment in real time. It captures system events (like resource state changes, scheduled tasks, or application events) and routes them to targets (like Lambda functions, SNS topics, SQS queues).
- With EventBridge, CloudWatch Events now supports both AWS service events and custom events from your own applications. It can also connect to third-party services and integrate with other event-driven systems.

6. CloudWatch Synthetics:

- CloudWatch Synthetics allows you to monitor the availability and performance of your web applications by creating canaries, which are automated scripts that simulate user activity and test your endpoints (APIs, web applications) at regular intervals.
- This helps ensure that your application is functioning correctly and provides insights into the health of your resources from an end-user perspective.

7. CloudWatch Contributor Insights:

- Contributor Insights is a feature that analyzes and visualizes log data to identify and prioritize contributors to high load or performance issues.

- For example, you can use Contributor Insights to identify the top IP addresses causing high latencies or the services contributing to errors, helping you pinpoint issues in your environment quickly.

8. CloudWatch ServiceLens:

- CloudWatch ServiceLens provides a unified view of your application's performance, helping you monitor, visualize, and analyze traces across multiple AWS services.
- It integrates with AWS X-Ray for distributed tracing, allowing you to drill down into specific requests and track them through microservices, identifying bottlenecks and performance issues.

How AWS CloudWatch Works

1. Collecting Data:

- CloudWatch automatically collects metrics from a wide variety of AWS services. For instance, Amazon EC2 instances send metrics like CPU utilization, memory usage, and disk I/O to CloudWatch.
- For custom data, you can create your own metrics and publish them to CloudWatch using the AWS SDK or API. This allows you to monitor anything that's important to your application, such as request counts, business KPIs, or application health.

2. Setting Up Alarms:

- Once CloudWatch is collecting metrics, you can set up CloudWatch Alarms to monitor specific metrics. If a metric exceeds or falls below a predefined threshold, the alarm can send a notification, trigger a Lambda function, or start an auto-scaling action.
- For example, if the CPU utilization of an EC2 instance exceeds 80%, you can trigger an alarm to automatically scale the instance or send an alert via email to the DevOps team.

3. Storing and Analyzing Logs:

- AWS CloudWatch Logs can be used to centralize the log data from all your AWS resources and applications. Logs are automatically stored in CloudWatch and can be analyzed for specific patterns using metric filters.
- You can create custom metrics from logs (e.g., count occurrences of certain error messages or user actions) and use them to trigger alarms or visualize the trends on CloudWatch Dashboards.

4. Event-driven Automation:

- CloudWatch Events (or EventBridge) can capture events from AWS services and send them to various targets, such as Lambda functions, SQS queues, or Step Functions. This allows for event-driven automation.

- For example, when an EC2 instance starts or stops, CloudWatch Events can trigger a Lambda function to perform additional tasks, such as updating a database or notifying a user.
-

Use Cases for AWS CloudWatch

1. Application Monitoring:

- CloudWatch allows you to monitor application performance by tracking custom metrics (e.g., request count, error rates, response time) from your applications and microservices.
- You can set up alarms based on these metrics to notify you when something goes wrong, such as an increase in error rate or slow response times.

2. Infrastructure Monitoring:

- CloudWatch provides metrics for various AWS resources (EC2, RDS, ELB, Lambda, etc.) that allow you to track the health and performance of your infrastructure.
- You can monitor resource utilization (CPU, memory, storage) and set up auto-scaling triggers to maintain performance as demand changes.

3. Security and Compliance:

- CloudWatch can capture logs related to security events, such as login attempts, changes to IAM roles, and access to critical resources.
- CloudWatch Logs can be integrated with services like AWS CloudTrail to help track API calls and monitor for any unauthorized access or anomalous behavior.

4. Business Metrics Monitoring:

- CloudWatch can be used to monitor application-specific business metrics, such as transaction volume, sales data, and customer activity. This allows business teams to get real-time insights into application usage and performance.

5. Real-time Monitoring and Alerts:

- Using CloudWatch dashboards, you can visualize your metrics in real-time and get alerts when certain thresholds are crossed. This helps in keeping an eye on the operational health of your AWS environment.
-

Benefits of AWS CloudWatch

1. Centralized Monitoring:

- CloudWatch provides a single pane of glass to monitor and manage all your AWS resources, custom metrics, logs, and events in one place.

2. Real-time Insights:

- With CloudWatch, you can gain real-time insights into the health of your resources and applications, helping you quickly detect and respond to issues.

3. Automation:

- You can automate actions in response to events or metrics, improving operational efficiency and reducing manual intervention.

4. Scalability:

- CloudWatch scales automatically with your AWS resources. Whether you're running a few EC2 instances or thousands, CloudWatch can handle the volume of metrics and logs.

5. Cost-Effective:

- You only pay for what you use, whether it's for storing logs, collecting metrics, or using alarms and dashboards. This pay-as-you-go pricing model helps keep costs manageable.

6. Security and Compliance:

- CloudWatch integrates with AWS security and compliance tools, helping you monitor logs and events for regulatory and operational purposes.

Configure aws cloudwatch to ec2

Method 2

- Launch the ec2
- Attach the iam role {ec2- cloudwatch or admin acces}
- Connect to ec2
- Install cloud watch agent

yum install amazon-cloudwatch-agent -y

```
[root@ip-172-31-19-9 ec2-user]# sudo su -
Last login: Thu Oct 24 08:48:41 UTC 2024 on pts/1
[root@ip-172-31-19-9 ~]# yum install amazon-cloudwatch-agent -y
Last metadata expiration check: 1 day, 17:16:02 ago on Tue Oct 22 15:32:51 2024.
Dependencies resolved.
=====
Package           Architecture      Version
=====
Installing:
amazon-cloudwatch-agent          x86_64          1.300044.0-1.amzn2023
Transaction Summary
=====
Install 1 Package

Total download size: 135 M
Installed size: 445 M
Downloading Packages:
amazon-cloudwatch-agent-1.300044.0-1.amzn2023.x86_64.rpm
-----
Total
Running transaction check
Transaction check succeeded.
Running transaction test
Transaction test succeeded.
Running transaction
Preparing :
```

- Edit this file or create file using below commnd

vi /opt/aws/amazon-cloudwatch-agent/bin/config.json

```
[root@ip-172-31-19-9 ~]# vi /opt/aws/amazon-cloudwatch-agent/bin/config.json
[root@ip-172-31-19-9 ~]#
```

- Add below json config details
- Change the highlighted details based on your requirement

```
{
  "agent": {
    "metrics_collection_interval": 60,
    "run_as_user": "root"
  },
  "logs": {
    "logs_collected": {
      "files": {
        "collect_list": [
          {
            "file_path": "/var/log/*",
            "log_group_class": "STANDARD",
            "log_group_name": "aws",
            "log_stream_name": "devops",
            "retention_in_days": 1
          }
        ]
      }
    }
  },
  "metrics": {
    "aggregation_dimensions": [
      [
        "InstanceId"
      ]
    ],
    "append_dimensions": {
      "AutoScalingGroupName": "${aws:AutoScalingGroupName}",
      "ImageId": "${aws:ImageId}",
      "InstanceId": "${aws:InstanceId}",
      "InstanceType": "${aws:InstanceType}"
    },
    "metrics_collected": {

```

```
"cpu": {  
    "measurement": [  
        "cpu_usage_idle",  
        "cpu_usage_iowait",  
        "cpu_usage_user",  
        "cpu_usage_system"  
    ],  
    "metrics_collection_interval": 60,  
    "resources": [  
        "*"  
    ],  
    "totalcpu": false  
},  
"disk": {  
    "measurement": [  
        "used_percent",  
        "inodes_free"  
    ],  
    "metrics_collection_interval": 60,  
    "resources": [  
        "*"  
    ]  
},  
"diskio": {  
    "measurement": [  
        "io_time"  
    ],  
    "metrics_collection_interval": 60,  
    "resources": [  
        "*"  
    ]  
},  
"mem": {  
    "measurement": [  
        "mem_used_percent"  
    ],  
    "metrics_collection_interval": 60  
},  
"statsd": {  
    "metrics_aggregation_interval": 60,  
}
```

```

    "metrics_collection_interval": 10,
    "service_address": ":8125"
},
"swap": {
    "measurement": [
        "swap_used_percent"
    ],
    "metrics_collection_interval": 60
}
},
"traces": {
    "buffer_size_mb": 3,
    "concurrency": 8,
    "insecure": false,
    "region_override": "us-east-2",
    "traces_collected": {
        "xray": {
            "bind_address": "127.0.0.1:2000",
            "tcp_proxy": {
                "bind_address": "127.0.0.1:2000"
            }
        }
    }
}
}

```

- Enter below command for start the cloudwatch and creation of log group in cloud watch console based on config.json file

```
sudo /opt/aws/amazon-cloudwatch-agent/bin/amazon-cloudwatch-agent-ctl -a fetch-config -m ec2 -c file:/opt/aws/amazon-cloudwatch-agent/bin/config.json -s
```

```
[root@ip-172-31-19-9 ~]# sudo /opt/aws/amazon-cloudwatch-agent/bin/amazon-cloudwatch-agent-ctl -a fetch-config -m ec2 -c file:/opt/aws/amazon-cloudwatch-agent/bin/config.json -s
***** processing amazon-cloudwatch-agent *****
I! Trying to detect region from ec2 D! [EC2] Found active network interface I! imds retry client will retry 1 timesSuccessfully fetched the config and saved it to /opt/aws/amazon-cloudwatch-agent/etc/amazon-cloudwatch-agent.d/file_config.json.tmp
Start configuration validation...
2024/10/24 08:50:51 I! Reading json config file path: /opt/aws/amazon-cloudwatch-agent/etc/amazon-cloudwatch-agent.d/file_config.json.tmp ...
2024/10/24 08:50:51 I! Valid Json input schema.
2024/10/24 08:50:51 D! ec2tagger processor required because append_dimensions is set
2024/10/24 08:50:51 D! delta processor required because metrics with diskio or net are set
2024/10/24 08:50:51 D! ec2tagger processor required because append_dimensions is set
2024/10/24 08:50:51 Configuration validation first phase succeeded
I! Detecting run_as_user...
I! Trying to detect region from ec2
D! [EC2] Found active network interface
I! imds retry client will retry 1 times
/opt/aws/amazon-cloudwatch-agent/bin/amazon-cloudwatch-agent -schematest -config /opt/aws/amazon-cloudwatch-agent/etc/amazon-cloudwatch-agent.toml
Configuration validation second phase succeeded
Configuration validation succeeded
amazon-cloudwatch-agent has already been stopped
Created symlink /etc/systemd/system/multi-user.target.wants/amazon-cloudwatch-agent.service → /etc/systemd/system/amazon-cloudwatch-agent.service.
```

- Open your cloud watch and check

CloudWatch

- Favorites and recent
- Dashboards
- Alarms
- Logs
- Log groups**
- Log Anomalies
- Live Tail
- Logs Insights
- Contributor Insights
- Metrics
- X-Ray traces
- Events

Log groups (1/6)

By default, we only load up to 10000 log groups.

| <input type="checkbox"/> Log group | Log class | Anomaly d... | Data pr... | Sensiti... |
|------------------------------------|-----------|---------------------------|------------|------------|
| /aws/lambda/lambda-rds | Standard | Configure | - | - |
| /aws/lambda/rds | Standard | Configure | - | - |
| /aws/lambda/test | Standard | Configure | - | - |
| keke | Standard | Configure | - | - |
| aws-ec2logs | Standard | Configure | - | - |
| aws | Standard | Configure | - | - |

Log streams Tags Anomaly detection **Metric filters** Subscription filters Contributor Insights Data protection

Log streams (1/1)

| <input checked="" type="checkbox"/> Log stream | Last event time |
|--|---------------------------|
| devops | 2024-10-24 08:50:53 (UTC) |

CloudWatch

- Favorites and recent
- Dashboards
- Alarms
- Logs
- Log groups**
- Log Anomalies
- Live Tail
- Logs Insights
- Contributor Insights
- Metrics
- X-Ray traces
- Events
- Application Signals New

Filter events - press enter to search

| Timestamp | Message |
|--------------------------|--|
| 2024-10-24T08:50:53.254Z | There are older events to load. Load more . |
| 2024-10-24T08:50:53.254Z | 2024-10-10T21:23:36Z DEBUG Extra commands: [u'-c', u'/tmp/imagebuilder-UIMtPQ/dnfconf-sT0zB9', u'-d10', u'-e10', u'--ins... |
| 2024-10-24T08:50:53.254Z | 2024-10-10T21:23:36Z DEBUG Cleaning data: packages dbcache metadata |
| 2024-10-24T08:50:53.254Z | 2024-10-10T21:23:36Z INFO 0 files removed |
| 2024-10-24T08:50:53.254Z | 2024-10-10T21:23:36Z DEBUG Cleaning up. |
| 2024-10-24T08:50:53.254Z | 2024-10-22T15:32:41+0000 INFO --- logging initialized --- |
| 2024-10-24T08:50:53.254Z | 2024-10-22T15:32:41+0000 INFO timer: config: 25 ms |
| 2024-10-24T08:50:53.254Z | 2024-10-22T15:32:41+0000 DEBUG Loaded plugins: builddep, changelog, config-manager, copr, debug, debuginfo-install, downl... |
| 2024-10-24T08:50:53.254Z | 2024-10-22T15:32:41+0000 DEBUG DNF version: 4.14.0 |
| 2024-10-24T08:50:53.254Z | 2024-10-22T15:32:41+0000 DEBUG Command: dnf --debuglevel 2 updateinfo |
| 2024-10-24T08:50:53.255Z | 2024-10-22T15:32:41+0000 DEBUG InstallRoot: / |
| 2024-10-24T08:50:53.255Z | 2024-10-22T15:32:41+0000 DEBUG Releasever: 2023.6.20241010 |
| 2024-10-24T08:50:53.255Z | 2024-10-22T15:32:41+0000 DEBUG cachedir: /var/cache/dnf |

Certificate manager and route53

Overview of Amazon Route 53

Amazon Route 53 is a scalable and highly available Domain Name System (DNS) web service provided by AWS. It is designed to route end-user requests to various AWS resources, including EC2 instances, S3 buckets, CloudFront distributions, and other domain names, using DNS. Route 53 also integrates with AWS services, making it an ideal solution for managing DNS in the cloud.

The service is named after DNS (port 53), the protocol used for DNS queries.

Route 53 is fully managed and offers several features, including domain registration, DNS routing, health checking, and traffic management.

Key Features of Amazon Route 53

1. DNS Management

- DNS Resolution: Route 53 translates friendly domain names (like www.example.com) into IP addresses that computers use to communicate with each other (e.g., 192.0.2.1).
- Scalability: Route 53 can handle large volumes of DNS queries without requiring additional configuration or resources, scaling automatically based on demand.
- Low Latency: AWS has a global network of DNS servers, so Route 53 can respond to DNS queries from users with low latency.

2. Domain Registration

- Domain Purchase and Transfer: Route 53 allows you to register new domain names directly through AWS or transfer existing domains to AWS.
- WHOIS Privacy: For supported TLDs (Top-Level Domains), WHOIS privacy can be enabled to protect personal information associated with domain registration.
- Automatic Renewal: Domains can be set to renew automatically, avoiding potential disruptions in service due to expired domains.

3. Traffic Routing and Load Balancing

- Simple Routing: Route 53 can map domain names to one or more IP addresses, making it a basic way to resolve domain names to endpoints like web servers.
- Weighted Routing: Route 53 supports weighted routing policies that allow you to distribute traffic between different resources (e.g., EC2 instances) based on a set of weights. This is useful for gradual migrations, blue/green deployments, or testing purposes.
- Latency-Based Routing: Route 53 can route traffic to the region with the lowest latency for the end-user, improving response time.
- Geo-Location Routing: Route 53 can direct traffic based on the geographic location of users (e.g., route European users to a European server, and U.S. users to a U.S.-based server).
- Failover Routing: This allows Route 53 to detect health issues with a resource and route traffic to a backup (failover) resource automatically if the primary resource becomes unavailable.
- Multivalue Answer: Similar to simple routing, but Route 53 can return multiple IP addresses for a domain name. This allows for automatic load balancing across multiple endpoints.

4. Health Checking and Monitoring

- Health Checks: Route 53 can monitor the health of resources (e.g., web servers, load balancers) by periodically sending HTTP(S) requests. If the resource is found to be unhealthy, Route 53 will redirect traffic to healthy endpoints based on routing policies (e.g., failover, weighted routing).
- CloudWatch Integration: Route 53 integrates with Amazon CloudWatch to monitor DNS query activity and health check results, providing visibility into the health of your infrastructure.

5. DNS Failover

- DNS Failover: This feature helps ensure high availability by automatically rerouting traffic to a healthy resource when a failure is detected. It is typically used in conjunction with Route 53 health checks.

6. Alias Records

- Alias Records: Unlike traditional DNS records, Alias records allow you to map a domain name to AWS resources (e.g., S3 buckets, CloudFront distributions, Elastic Load Balancers) without needing an IP address. This is particularly useful for integrating Route 53 with AWS services.

- No Additional Cost for Alias Records: Alias records are treated as native DNS records, and AWS doesn't charge for query traffic to them (unlike standard CNAME records).

7. DNSSEC (Domain Name System Security Extensions)

- DNSSEC Support: Route 53 supports DNSSEC to add an additional layer of security, helping to protect against DNS spoofing and cache poisoning attacks. With DNSSEC, responses to DNS queries are signed to ensure that users are directed to the correct website.

8. Traffic Flow and Global Traffic Management

- Traffic Flow: Route 53's visual interface, Traffic Flow, allows you to create complex routing policies, including weighted, latency-based, geolocation, and failover routing, in an easy-to-use interface.
- Geoproximity Routing: Route 53 can route traffic based on the geographic location of users and resources, which is useful for managing global traffic.

9. Integration with AWS Services

- AWS CloudFront: Easily integrate Route 53 with Amazon CloudFront, AWS's Content Delivery Network (CDN), for faster and more secure delivery of web content globally.
- Elastic Load Balancing (ELB): Route 53 works seamlessly with AWS's Elastic Load Balancers to route traffic to your EC2 instances, providing load balancing and scaling.
- Amazon S3: You can use Route 53 to route traffic to Amazon S3 static websites or S3 buckets.
- AWS Elastic Beanstalk: Route 53 can be configured to manage traffic for applications deployed using Elastic Beanstalk.

Key Benefits of Amazon Route 53

1. Highly Available and Scalable: Route 53 is designed to be highly available with a global network of DNS servers and auto-scaling to handle millions of DNS queries with low latency.
2. Ease of Use: The Route 53 console offers a simple and easy-to-navigate interface for managing DNS records and configuring routing policies.
3. Integrated with AWS: Route 53 integrates tightly with other AWS services, enabling easy management of DNS for AWS resources (like EC2, S3, CloudFront, and more).

4. Cost-Effective: Route 53 offers pay-as-you-go pricing, meaning you only pay for the DNS queries made and the resources you use. This makes it an economical choice for businesses of all sizes.
 5. Security: Route 53 supports DNSSEC for domain name validation and integrates with AWS IAM for access control. Additionally, its health checks can help automatically reroute traffic during failures, ensuring continuity.
 6. Global Traffic Management: With features like geo-location routing, latency-based routing, and weighted routing, Route 53 enables intelligent traffic management to improve application performance and ensure a better user experience across the globe.
-

Common Use Cases for Amazon Route 53

1. Website Hosting:
 - Route 53 can route users to AWS-hosted websites, such as static sites hosted on S3 or dynamic applications hosted on EC2 instances.
 2. Disaster Recovery:
 - With DNS failover and health checks, Route 53 can reroute traffic to a healthy resource in the event of an outage or server failure, enabling high availability and business continuity.
 3. Global Traffic Management:
 - Businesses with a global customer base can use geolocation or latency-based routing to ensure that users are directed to the nearest or fastest server, improving load times and user experience.
 4. Multi-Region Applications:
 - Companies with applications deployed across multiple AWS regions can use Route 53's multi-region routing capabilities to improve redundancy and fault tolerance.
 5. Content Delivery Networks (CDNs):
 - Route 53 can be used in conjunction with Amazon CloudFront to manage DNS for global content delivery, caching, and content acceleration.
 6. Hybrid Cloud Environments:
 - Route 53 can be used to route traffic between on-premises systems and cloud resources, supporting hybrid cloud architectures.
-

Pricing for Amazon Route 53

Route 53 pricing is based on several factors:

1. Hosted Zones: The number of hosted zones you create.
2. DNS Queries: The number of DNS queries processed by Route 53.
3. Domain Registration: Costs associated with registering domain names.
4. Health Checks: Fees for performing health checks on resources.
5. Traffic Flow: There is a cost for creating and managing traffic flow policies.

The pay-as-you-go model makes Route 53 cost-effective for small websites to large-scale applications.

Practical

Add your domain name in route 53

| Record ... | Type | Routin... | Differ... | Alias | Value/Route to |
|-------------|------|-----------|-----------|-------|--|
| narni.co.in | NS | Simple | - | No | ns-1959.awsdns ns-674.awsdns ns-305.awsdns ns-1025.awsdns |
| narni.co.in | SOA | Simple | - | No | ns-1959.awsdns |

Then search for certificate manager in aws search open it

| Certificate ID | Domain name | Type | Status | In use | Renewal eligibility | Key algorithm |
|----------------|-------------|------|--------|--------|---------------------|---------------|
| | | | | | | |

Click on request

Then click on next

Request certificate

Certificate type Info

ACM certificates can be used to establish secure communications access across the internet or within an internal network. Choose the type of certificate for ACM to provide.

Request a public certificate

Request a public SSL/TLS certificate from Amazon. By default, public certificates are trusted by browsers and operating systems.

Request a private certificate

No private CAs available for issuance.

Requesting a private certificate requires the creation of a private certificate authority (CA). To create a private CA, visit [AWS Private Certificate Authority](#)

Cancel

Next

Enter your domain name here

Request public certificate

Domain names

Provide one or more domain names for your certificate.

Fully qualified domain name Info

narni.co.in

[Add another name to this certificate](#)

You can add additional names to this certificate. For example, if you're requesting a certificate for "www.example.com", you might want to add the name "example.com" so that customers can reach your site by either name.

Validation method Info

Select a method for validating domain ownership.

DNS validation - recommended

Choose this option if you are authorized to modify the DNS configuration for the domains in your certificate request.

Email validation

Choose this option if you do not have permission or cannot obtain permission to modify the DNS configuration for the domains in your certificate request.

Click on request

CRRA

Choose this option if you do not have permission or cannot obtain permission to modify the DNS configuration for the domains in your certificate request.

Key algorithm Info

Select an encryption algorithm. Some algorithms may not be supported by all AWS services.

RSA 2048

RSA is the most widely used key type.

ECDSA P 256

Equivalent in cryptographic strength to RSA 3072.

ECDSA P 384

Equivalent in cryptographic strength to RSA 7680.

Tags Info

To help you manage your certificates, you can optionally assign your own metadata to each resource in the form of tags.

No tags associated with this resource.

[Add tag](#)

You can add 50 more tag(s).

Cancel

Previous

Request

This is your certificate

Click on certificate id.

Certificates (1)

[C](#)

[Delete](#)

[Manage expiry events](#)

[Import](#)

[Request](#)

< 1 >

| <input type="checkbox"/> | Certificate ID | Domain name | Type | Status | In use | Renewal eligibility | Key algorithm |
|--------------------------|--------------------------------------|-------------|---------------|--------|--------|---------------------|---------------|
| <input type="checkbox"/> | 4c6f28d6-f6f1-4213-a4b4-f8f996c47ccb | narni.co.in | Amazon Issued | Issued | No | Ineligible | RSA 2048 |

Click on create records in route 53

The screenshot shows the AWS Certificate Manager interface. At the top, it displays the path: AWS Certificate Manager > Certificates > 4c6f28d6-f6f1-4213-a4b4-f8f996c47ccb. Below this, the certificate details are shown:

| | | | |
|------------|--|--------|--------|
| Identifier | 4c6f28d6-f6f1-4213-a4b4-f8f996c47ccb | Status | Issued |
| ARN | arn:aws:acm:ap-south-1:483216680875:certificate/4c6f28d6-f6f1-4213-a4b4-f8f996c47ccb | | |
| Type | Amazon Issued | | |

Below the details, there is a section titled "Domains (1)" with a table:

| Domain | Status | Renewal status | Type | CNAME name | CNAME value |
|--------------|--------------------|----------------|-------|---------------------------------|---|
| vardhan.live | Pending validation | - | CNAME | _5fae5fe08794f11532c58766d40bd0 | _e00cb41545e6ebe3bf0243... 3bf02430b443896 ce.mhbtbspdnt.acm -validations.aws. |

Buttons at the top right of the domain table include "Create records in Route 53" and "Export to CSV".

Click on create records

it will be added in to route 53 records list

The screenshot shows the "Create DNS records in Amazon Route 53" interface. It lists a single record:

| Domain | Validation status | Type | CNAME name | CNAME value | Is domain in Route 53? |
|--------------|--------------------|-------|---------------------------------|---|------------------------|
| vardhan.live | Pending validation | CNAME | _5fae5fe08794f11532c58766d40bd0 | _e00cb41545e6ebe3bf0243... 3bf02430b443896 ce.mhbtbspdnt.acm -validations.aws. | Yes |

At the bottom right, there are "Cancel" and "Create records" buttons.

The screenshot shows the "Records (3) Info" section of the Route 53 console. It lists three records:

| Record ... | Type | Routin... | Differ... | Alias | Value/Route traffic to | TTL (s...) | Health ... |
|--------------|-------|-----------|-----------|-------|---|------------|------------|
| vardhan.live | NS | Simple | - | No | ns-1560.awsdns-03.co.uk. ns-1189.awsdns-20.org. ns-858.awsdns-43.net. ns-68.awsdns-08.com. | 172800 | - |
| vardhan.live | SOA | Simple | - | No | ns-1560.awsdns-03.co.uk. a... | 900 | - |
| _5fae5fe0... | CNAME | Simple | - | No | _e00cb41545e6ebe3bf0243... | 300 | - |

A red box highlights the last record, which corresponds to the one created via the AWS Certificate Manager.

Then create a instance and add user data or manually connect to instance and install httpd server and check with public ip

| Instance summary for i-0f3b5da83536caa58 (ec2-lb) | | |
|--|---|---|
| Updated less than a minute ago | | |
| Instance ID i-0f3b5da83536caa58 (ec2-lb) | Public IPv4 address 52.66.180.51 open address | Private IPv4 addresses 10.0.2.28 |
| IPv6 address - | Instance state Running | Public IPv4 DNS - |
| Hostname type IP name: ip-10-0-2-28.ap-south-1.compute.internal | Private IP DNS name (IPv4 only) ip-10-0-2-28.ap-south-1.compute.internal | Elastic IP addresses - |
| Answer private resource DNS name - | Instance type t2.micro | AWS Compute Optimizer finding Opt-in to AWS Compute Optimizer for recommendations. Learn more |
| Auto-assigned IP address 52.66.180.51 [Public IP] | VPC ID vpc-0f745a839e3f1258d (my-cust-vpc) | Auto Scaling Group name |
| IAM Role | Subnet ID | |

Create a target group for above instance

Target group name

A maximum of 32 alphanumeric characters including hyphens are allowed, but the name must not begin or end with a hyphen.

Protocol : Port
Choose a protocol for your target group that corresponds to the Load Balancer type that will route traffic to it. Some protocols now include anomaly detection for the targets and you can set mitigation options once your target group is created. This choice cannot be changed after creation

HTTP 1-65535

IP address type
Only targets with the indicated IP address type can be registered to this target group.

IPv4
Each instance has a default network interface (eth0) that is assigned the primary private IPv4 address. The instance's primary private IPv4 address is the one that will be applied to the target.

IPv6
Each instance you register must have an assigned primary IPv6 address. This is configured on the instance's default network interface (eth0). [Learn more](#)

VPC
Select the VPC with the instances that you want to include in the target group. Only VPCs that support the IP address type selected above are available in this list.

my-cust-vpc
vpc-0f745a839e3f1258d
IPv4 VPC CIDR: 10.0.0.0/16

Mention your health check path as your application path

Then click on next

HTTP

Health check path
Use the default path of "/" to perform health checks on the root, or specify a custom path if preferred.

Up to 1024 characters allowed.

Advanced health check settings

Attributes

Certain default attributes will be applied to your target group. You can view and edit them after creating the target group.

Tags - optional
Consider adding tags to your target group. Tags enable you to categorize your AWS resources so you can more easily manage them.

Cancel **Next**

Select your instance and create a target group

| | |
|--|---|
| 80 | 1-65535 (separate multiple ports with commas) |
| <input type="button" value="Include as pending below"/> | |
| 1 selection is now pending below. Include more or register targets when ready. | |

Review targets

| Targets (1) | | | | | | | | <input type="button" value="Remove all pending"/> |
|--|--------|------|--|-----------------|-------------|----------------------|---------------|---|
| <input type="text" value="Filter targets"/> <input checked="checked" type="checkbox"/> Show only pending | | | | | | | | < 1 > <input type="button" value="Search"/> |
| Instance ID | Name | Port | State | Security groups | Zone | Private IPv4 address | Subnet ID | |
| i-0f3b5da83536caa58 | ec2-lb | 80 | Running | my-cust-sg | ap-south-1b | 10.0.2.28 | subnet-064357 | <input type="button" value="Edit"/> |

1 pending

Next create a load balancer

Create Application Load Balancer Info

The Application Load Balancer distributes incoming HTTP and HTTPS traffic across multiple targets such as Amazon EC2 instances, microservices, and containers, based on request attributes. When the load balancer receives a connection request, it evaluates the listener rules in priority order to determine which rule to apply, and if applicable, it selects a target from the target group for the rule action.

Basic configuration

Load balancer name
Name must be unique within your AWS account and can't be changed after the load balancer is created.

A maximum of 32 alphanumeric characters including hyphens are allowed, but the name must not begin or end with a hyphen.

Scheme Info
Scheme can't be changed after the load balancer is created.
 Internet-facing
An internet-facing load balancer routes requests from clients over the internet to targets. Requires a public subnet. Learn more

Internal
An internal load balancer routes requests from clients to targets using private IP addresses. Compatible with the IPv4 and Dualstack IP address types.

Select you subnets

VPC Info
Select the virtual private cloud (VPC) for your targets or you can [create a new VPC](#). Only VPCs with an internet gateway are enabled for selection. The selected VPC can't be deleted while the load balancer is created. To confirm the VPC for your targets, view your [target groups](#).

Mappings Info
Select at least two Availability Zones and one subnet per zone. The load balancer routes traffic to targets in these Availability Zones only. Availability Zones that are not supported by the load balancer or the VPC are not available for selection.

ap-south-1a (aps1-az1)
Subnet

IPv4 address
Assigned by AWS

ap-south-1b (aps1-az3)
Subnet

IPv4 address
Assigned by AWS

ap-south-1c (aps1-az2)

THIS IS MAJOR STEP FOR HTTPS

Select protocol as https

And select you target group

Listeners and routing [Info](#)

A listener is a process that checks for connection requests using the port and protocol you configure. The rules that you define for a listener determine how the load balancer routes requests to its registered targets.

▼ Listener HTTPS:443

[Remove](#)

Protocol Port Default action [Info](#)

HTTPS : 443
1-65535

Forward to **target-group**
Target type: Instance, IPv4

HTTP ▾

[Create target group](#)

Listener tags - optional

Consider adding tags to your listener. Tags enable you to categorize your AWS resources so you can more easily manage them.

[Add listener tag](#)

You can add up to 50 more tags.

Security policy as it is

Select you certificate here this certificate comes from certificate manager you previously created

Secure listener settings [Info](#)

These settings will apply to all of your secure listeners. Once created, you can manage these settings per listener.

Security policy [Info](#)

Your load balancer uses a Secure Socket Layer (SSL) negotiation configuration called a security policy to manage SSL connections with clients. [Compare security policies](#)

Security category

Policy name

All security policies

ELBSecurityPolicy-TLS13-1-2-2021-06 (recommended)

Default SSL/TLS server certificate

The certificate used if a client connects without SNI protocol, or if there are no matching certificates. You can source this certificate from AWS Certificate Manager (ACM), Amazon Identity and Access Management (IAM), or import a certificate. This certificate will automatically be added to your listener certificate list.

Certificate source

From ACM

From IAM

Import certificate

Certificate (from ACM)

The selected certificate will be applied as the default SSL/TLS server certificate for this load balancer's secure listeners.

narni.co.in
4c6f28d6-f6f1-4213-a4b4-f8f996c47ccb

[Request new ACM certificate](#)

Client certificate handling [Info](#)

Click on create load balancer

Service integrations [Edit](#)
AWS WAF: *None*
AWS Global Accelerator: *None*

Tags [Edit](#)
None

Attributes

Certain default attributes will be applied to your load balancer. You can view and edit them after creating the load balancer.

Creation workflow and status

▶ **Server-side tasks and status**
After completing and submitting the above steps, all server-side tasks and their statuses become available for monitoring.

Cancel **Create load balancer**

Then open your route 53

Click on create record

Records (4) DNSSEC signing Hosted zone tags (0)

Records (4) [Info](#) [C](#) Delete record Import zone file **Create record**

Automatic mode is the current search behavior optimized for best filter results. To change modes go to settings.

Filter records by property or value Type ▾ Routing policy ▾ Alias ▾ < 1 > ⚙

Select alias record

Type = alias to application and classic load balancer

Region= your load balancer region

Select your load balancer dns url

Click on create record

Create record [Info](#)

Quick create record [Switch to wizard](#)

▼ Record 1 [Delete](#)

Record name | [Info](#) subdomain narni.co.in
Keep blank to create a record for the root domain.

Alias

Route traffic to | [Info](#) Alias to Application and Classic Load Balancer
Asia Pacific (Mumbai)
Q dualstack.https-loadbalanerr-739461246.ap-south-1.elb.amazonaws.com

Alias hosted zone ID: ZP97RAFLXTNZK

Routing policy | [Info](#) Simple routing Evaluate target health Yes

[Add another record](#)

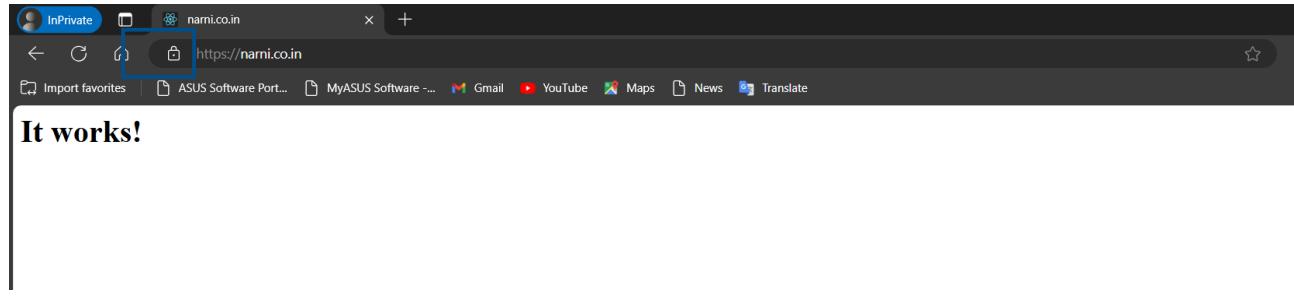
[Cancel](#) [Create records](#)

This is your record

| Records (4) Info | | DNSSEC signing | | Hosted zone tags (0) | |
|---|-------------|-------------------|-------------------------------|----------------------------------|--|
| | | C | Delete record | Import zone file | Create record |
| Filter records by property or value | | Type | Routing policy | Alias | |
| <input type="checkbox"/> | narni.co.in | A | Simple | - | Yes |
| | | | | | dualstack.https-loadbalanerr... |
| <input type="checkbox"/> | narni.co.in | NS | Simple | - | No |
| | | | | | ns-1959.awsdns-52.co.uk. ns-674.awsdns-20.net. ns-305.awsdns-38.com. ns-1025.awsdns-00.org. |
| <input type="checkbox"/> | narni.co.in | SOA | Simple | - | No |
| | | | | | ns-1959.awsdns-52.co.uk. a... 900 - |
| <input type="checkbox"/> | _73ed4ba... | CNAME | Simple | - | No |
| | | | | | _1ba7a9e81e796832b1edd6... 300 - |

Hit with your domain name

https working successfully



VEERA

Aws lambda

Overview of AWS Lambda

AWS Lambda is a serverless compute service provided by Amazon Web Services (AWS) that lets you run code in response to events without provisioning or managing servers. With AWS Lambda, you can run backend code for virtually any application or service, all without worrying about the infrastructure required to run it.

Lambda is one of the foundational services in the serverless computing model, where you don't need to manage servers or infrastructure yourself. Instead, AWS manages all the resources, scaling, and availability for you. You simply upload your code, define the event triggers, and Lambda takes care of execution, scaling, and monitoring.

Key Features of AWS Lambda:

1. Event-Driven Computing:

- AWS Lambda is event-driven, meaning that your Lambda functions are triggered by specific events or data. These events could come from many AWS services, such as S3, DynamoDB, SNS, API Gateway, CloudWatch, or even HTTP requests from your application.
- For example, you can trigger a Lambda function every time an object is uploaded to an S3 bucket or whenever an update happens in a DynamoDB table.

2. No Server Management:

- Lambda abstracts away the underlying infrastructure, so you don't have to worry about provisioning servers, scaling, or handling server maintenance.
- You upload your code (as a function) and specify the trigger (event). Lambda automatically takes care of scaling the resources to meet demand and managing fault tolerance.

3. Automatic Scaling:

- AWS Lambda scales automatically by running code in response to each event, and each event is handled by a separate instance of your function.
- If there are a high number of incoming events (e.g., many users accessing your application), Lambda will automatically create more instances of the function to handle the increased load, without any intervention from you.

4. Short-Running Tasks:

- Lambda functions are designed to handle short, stateless tasks. They typically run for a few seconds to a few minutes at most.
- AWS Lambda has a maximum execution timeout of 15 minutes per function invocation, but for most use cases, tasks are completed in less than a minute.

5. Supports Multiple Programming Languages:

- AWS Lambda supports several popular programming languages, including:
 - Node.js
 - Python
 - Java
 - Go
 - C# (.NET Core)
 - Ruby
- Custom Runtime: Lambda also allows you to create custom runtimes if you want to run code written in other languages.
- This flexibility lets developers choose the best programming language for their specific needs.

6. Pay-Per-Use Pricing:

- Lambda follows a pay-per-use pricing model. You only pay for the compute time your functions consume when they are triggered, and there are no charges when your code is idle.
- You are billed based on the number of requests and the duration of code execution, which is calculated in milliseconds.
- There is also a free tier that provides 1 million free requests and 400,000 GB-seconds of compute time per month.

7. Integrated with AWS Services:

- AWS Lambda integrates seamlessly with other AWS services like Amazon S3, DynamoDB, API Gateway, SNS, CloudWatch, SQS, and more.
- This enables developers to build complex workflows and event-driven applications without managing server infrastructure.

8. Stateless and Immutable:

- Lambda functions are stateless by design, meaning that each invocation of the function is independent, and there's no stored state between function calls.
- However, you can use external storage services like Amazon S3, DynamoDB, or ElastiCache to store data if needed.

9. Versioning and Aliases:

- AWS Lambda supports function versioning, allowing you to publish and manage different versions of a Lambda function. You can then create aliases to point to specific versions of a function, providing flexibility for testing, deployment, and version management.

10. Monitoring and Logging:

- AWS Lambda integrates with Amazon CloudWatch for monitoring and logging. Lambda automatically creates logs for every function invocation.
 - You can monitor the function's performance, success/failure rates, and error logs via CloudWatch, making it easy to track function behavior and diagnose issues.
-

How AWS Lambda Works:

1. Upload Your Code:

- You write your function and upload it to AWS Lambda, or you can store it in a zip file or a container image.
- You also specify the runtime environment (e.g., Node.js, Python) and any required permissions (via IAM roles) for your function to access other AWS services.

2. Define Event Triggers:

- You configure your Lambda function to be triggered by specific events (such as changes in S3 buckets, DynamoDB streams, SNS notifications, etc.).
- For example, an image uploaded to an S3 bucket can trigger a Lambda function to process the image or generate a thumbnail.

3. Execution:

- When the specified event occurs, AWS Lambda automatically runs the function.
- Lambda runs the function in a stateless container and allocates resources such as CPU and memory based on the configuration you provide.

4. Scaling:

- If multiple events are triggered at once (e.g., many users accessing your app), Lambda automatically scales to handle the load by running multiple instances of your function in parallel.
- If the event volume decreases, Lambda reduces the number of active function instances accordingly, without you needing to manage the infrastructure.

5. Pay for Usage:

- You are billed for the total compute time your function uses, based on the number of invocations and how long each function takes to execute.
-

Use Cases for AWS Lambda:

1. Real-Time File Processing:

- Lambda is often used for processing files in real-time. For example, if you upload files to an S3 bucket, Lambda can automatically trigger a function to process, transform, or analyze the files.
- Example: Image resizing, video transcoding, log file analysis.

2. Web Application Backends:

- You can use AWS Lambda to build scalable web application backends. By integrating with API Gateway, Lambda can respond to HTTP requests (GET, POST, etc.) without managing web servers.
- Example: User authentication, data processing, and serving API responses.

3. Data Pipelines:

- Lambda can help in building data pipelines where data is processed in real-time. For example, triggering Lambda functions to process data as it is inserted into DynamoDB or Kafka.
- Example: Real-time analytics, data transformation, or ETL jobs.

4. Serverless Microservices:

- Lambda is a perfect fit for implementing microservices architectures. Each Lambda function can perform a specific task (e.g., payment processing, user management), and together they form a serverless application.

5. IoT Applications:

- Lambda is commonly used in IoT applications where devices trigger events that Lambda processes.
- Example: A temperature sensor sends data to an S3 bucket or DynamoDB, and Lambda triggers processing or sends alerts.

6. Automating AWS Management Tasks:

- Lambda can be used to automate operational tasks like resource provisioning, backup operations, security monitoring, and scaling of AWS resources.
- Example: Automatically stop unused EC2 instances during off-hours.

7. Scheduled Jobs and Cron Jobs:

- You can run Lambda functions at scheduled intervals using CloudWatch Events. This is similar to running cron jobs, but with no need for provisioning servers.
 - Example: Periodic reports, system health checks, and backups.
-

Pricing for AWS Lambda:

- Requests: You are charged for the total number of requests to your Lambda function. The first 1 million requests per month are free.
 - Duration: Lambda is charged based on the duration of your function execution, measured in milliseconds. You specify the amount of memory allocated to the function, and the cost is based on the resources consumed during execution.
 - Free Tier: AWS offers a free tier that includes 1 million requests and 400,000 GB-seconds of compute time per month for free.
-

Advantages of AWS Lambda:

1. Serverless Architecture: You don't need to manage infrastructure, making it easier and cheaper to scale applications.
 2. Automatic Scaling: Lambda automatically scales based on incoming events, without needing you to configure or manage servers.
 3. Cost-Effective: You only pay for what you use, which makes Lambda cost-effective, especially for intermittent workloads.
 4. Speed of Development: You can focus on writing code without worrying about managing servers or infrastructure.
 5. High Availability: AWS Lambda is highly available by default, as AWS manages the infrastructure across multiple availability zones.
-

Limitations of AWS Lambda:

1. Execution Time Limit: Lambda functions have a maximum execution time of 15 minutes. If your function requires more time, you may need to use an alternative like EC2 or Fargate.
2. Statelessness: Lambda functions are stateless by design, meaning they don't retain data between invocations. You'll need external storage for any persistent data (e.g., S3, DynamoDB).

3. Cold Starts: While Lambda can scale automatically, the first time a function is invoked after a period of inactivity (called a cold start), it can experience some latency.

VEERA