# Container Design and Deployment

Deploying a Microservice Application over an AWS EC2 Instance and

implementing Auto Scaling Group Policy and load balancing.

**Introduction**

This research project focuses on the deployment of a Microservice on an Amazon Web Services (AWS) Elastic Compute Cloud (EC2) Instance whilst implementing specific policies and features available within the AWS EC2 Instance, a publicly available cloud-computing platform. The research uses a .jar application called a "Product Service" as a Microservice application to be deployed within the AWS EC2 Instance. As this is an introductory research project, all steps and required software applications are outlined within this research to demonstrate the student'sunderstanding of deploying a Microservice application. Brief terminologies are defined, and key aspects of the demonstration are provided as screenshots within this document. All demonstrations during this exercise are performed using a Free-tier AWS account. The research concludes with an overview of the AWS Auto Scaling group policy, and the challenges and limitations of its usage.

In the dynamic world of cloud computing, managing the scalability and availability of your applications is a critical endeavor. AWS Auto Scaling is a service provided by Amazon Web Services (AWS) that empowers businesses to automatically adjust the capacity of their resources based on specific performance metrics. While this service offers a powerful solution for ensuring optimal resource utilization and application performance, it comes with certain regional limitations.

As of the last knowledge update in September 2021, AWS Auto Scaling is inherently region-specific. This means that Auto Scaling policies are created and managed within a single AWS region and are applied only to resources located within that region. This regional confinement can pose a challenge for organizations with multi-region applications, as it necessitates the setup of Auto Scaling policies separately in each AWS region. This introduction explores the implications of these regional constraints on managing resources in a multi-region context while highlighting potential solutions to streamline the process. However, it's crucial to stay informed about any updates or enhancements that AWS might have introduced to Auto Scaling since that time

**Amazon Web Services (AWS):** A publicly available cloud computing platform owned by Amazon.

**AWS Elastic Compute Cloud (EC2):** A web service that you can use to launch and manage Linux/UNIX and Windows Server instances in Amazon data centres.

**Amazon EC2 Auto Scaling:** A web service that launches or terminates instances automatically based on user-defined policies, schedules, and health checks.

**Amazon Machine Image (AMI):** An encrypted machine image stored in Amazon Elastic Block Store (Amazon EBS) or Amazon Simple Storage Service. AMIs function like a template of a computer's root drive. They contain the operating system and can also include software and layers of your application, such as database servers, middleware, and web servers.

**Microservice**: Microservices are an architectural and organizational approach to software development created to speed up deployment cycles, foster innovation, and ownership, improve maintainability and scalability of software applications, and scale organizations delivering software and services by using an agile approach that helps teams work independently.

**S3 Bucket** is a storage for the internet within the AWS EC2 You can use it to store and retrieve any amount of data at any time, from anywhere on the web

**Postman**: is a complete API development platform that helps you manage your APIs in every stage of development, from designing and testing, to publishing API documentation and monitoring.

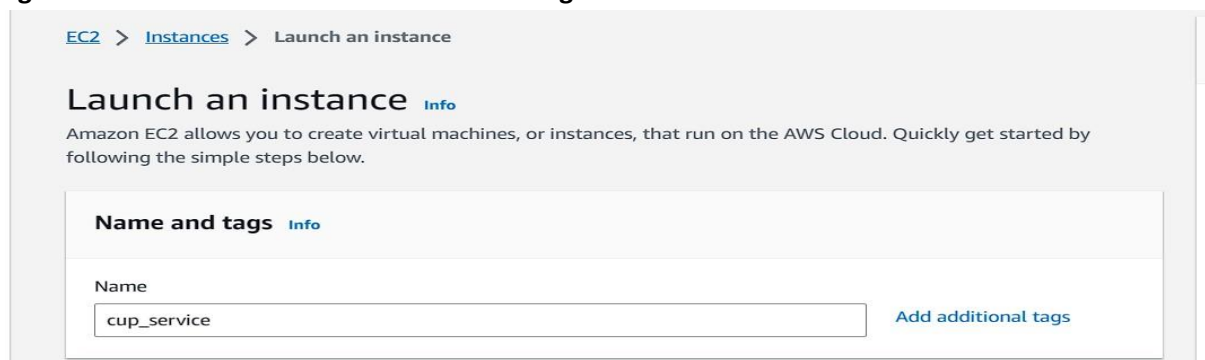**Steps for deploying a Microservice on an AWS EC2 Instance.**

The implementation considers that there exists a Free-tier AWS EC2 Instance Account. 1. Logging into your EC2 Instance using an IAM user account. It is important to note that granting users root access to an EC2 Instance is not the best security practice for accessing AWS resources, therefore the root user must set up another user account for multiple access to the resources. This section provides a detailed step for setting up an IAM user account.

Steps for creating a user account from the root user account

a. Go to the IAM Console (search IAM) select users and click Add user to create a new user

b. Create a username and set a password

c. Set user account policy or assigned user account to a group. For this demo, we set the user policy to Administrator which grants the user all access except billing.

 O Open the Amazon EC2 console at https://console.aws.amazon.com/ec2/  write the name of the instance and additional tags if you wish.

**Figure 1 shows how to launch an instance using AWS**



 O In the navigation pane, NETWORK & SECURITY, choose Key Pairs.

Note: The navigation pane is on the left side of the Amazon EC2 console. If you do not see the pane, it might be minimized; choose the arrow to expand the pane.  You can  allow Http traffic to expose port 80

**Figure 2 : shows creation of security groups and allow HTTP traffic for port 80**



**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-14'** 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

⚠ Rules with source of 0.0.0.0/0 allow all IP addresses to access your instance. We recommend setting security group rules to allow access from known IP addresses only.

○ Choose Create Key Pair.

○ Enter a name for the new key pair in the Key pair name field of the Create Key Pair dialog box, and then choose Create.

○ The private key file is automatically downloaded by your browser. The base file name is the name you specified as the name of your key pair, and the file name extension is .pem. Save the private key file in a safe place.

○ Important: This is the only chance for you to save the private key file. You'll need to provide the name of your key pair when you launch an instance and the corresponding private key each time you connect to the instance.

○ If you will use an SSH client on a Mac or Linux computer to connect to your Linux instance, use the following command to set the permissions of your private key file so that only you can read it.

**Figure 3: shows Instance launched form AWS**



Successfully started i-0fcbf48a1df437eaf    ✕

**Instances (1/1)** Info    ⟳    Connect    Instance state ▼    Actions ▼    **Launch instances** ▼

🔍 Find Instance by attribute or tag (case-sensitive)

Instance state = running    ✕    Clear filters    ‹ 1 › ⚙

| ☑ | Name ✎ | ▽ | Instance ID | Instance state | ▽ | Instance type | ▽ | Status check | Alarm status | Availability Zone | ▽ | Public IP |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ☑ | cup_service | | i-0fcbf48a1df437eaf | ⊘ Running ⊕ ⊖ | | t3.micro | | ⊘ Initializing | No alarms + | eu-north-1a | | ec2-13-5 |

**Instance: i-0fcbf48a1df437eaf (cup_service)**    ⚙ ✕

Launch Session Manager and install Java and MariaDb on the EC2 Instance. Create the necessary database and tables for the Product Service Application. Check both installations by checking the application versions. To do so, enter the following command in the terminal.

➢ sudo su – (to change to root directory of the EC2 Instance from the session manager)

➢ mariadb –version (to check if mariadb server is installed) else enter the next command

➢ yum install -y mariadb-server (to install mariadb server on your EC2) spercify the version if mariadb is not found  using  yum search mariadb

➢ systemctl enable mariadb (to enable the mariadb server for running)

➢ systemctl start mariadb (to start the mariadb server)

➢ mysql_secure_installation (to setup mysql root user access. You may choose to set a password that will be used to log on to your mysql. Follow the prompts and allow default settings for all other prompts)

➢ mysql -uroot -p (to securely login to mysql and enter your password if you have setup a login password from the previous command and your mysql environment should be set and ready for your sql commands) Enter the following sql commands to create your database and related tables.

➢ create database servicedb; (to create a database called serviced)

➢ java –version (to check if the Java application is installed) else enter the next command  ➢ yum install -y java (to install java on your EC2)

## Creating s3 bucket

Create the S3 Bucket: Use the aws s3api create-bucket command to create a new S3 bucket. Replace your-bucket-name with your desired bucket name and your-region with the AWS region where you want to create the bucket. You can also add additional options like versioning, logging, and ACLs by including additional parameters in the command.

**Figure 4:  shows creation of S3 bucket in AWS**



Amazon S3 > Buckets > Create bucket

## Create bucket Info
Buckets are containers for data stored in S3. Learn more ↗

**General configuration**

Bucket name

ken97awsbucket

Bucket name must be unique within the global namespace and follow the bucket naming rules. See rules for bucket naming ↗

AWS Region

Europe (Stockholm) eu-north-1 ▼

Copy settings from existing bucket - *optional*
Only the bucket settings in the following configuration are copied.

**Choose bucket**

Remember that S3 bucket names must be globally unique within AWS, and they have naming rules (e.g., no uppercase letters, no underscores). Also, keep in mind that there may be some limitations on the number of S3 buckets you can create within your AWS account.

Once your S3 bucket is created, you can start uploading and managing objects (files) within it using various methods, such as the AWS Management Console, AWS CLI, or SDKs.

**Figure 5: shows  properties of S3 AWS bucket**



## ken97awsbucket Info
Publicly accessible

Objects | Properties | Permissions | Metrics | Management | Access Points

**Objects** (4)
Objects are the fundamental entities stored in Amazon S3. You can use Amazon S3 inventory ↗ to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. Learn more ↗

C | Copy S3 URI | Copy URL | Download | Open ↗ | Delete | Actions ▼ | Create folder

Upload

Q Find objects by prefix          < 1 > ⚙

| | Name ▲ | Type ▽ | Last modified ▽ | Size ▽ | Storage class ▽ |
|---|---|---|---|---|---|
| ☐ | 📄 couponservice-0.0.1-SNAPSHOT.jar | jar | October 5, 2023, 17:56:02 (UTC+01:00) | 37.1 MB | Standard |
| ☐ | 📁 images/ | Folder | - | - | - |
| ☐ | 📄 index.html | html | October 5, 2023, 16:20:56 (UTC+01:00) | 1.7 KB | Standard |
| ☐ | 📁 styles/ | Folder | - | - | - |

Copy the  microservice couponservice.jar to the s3 bucket and make it accessible public. In the EC2 console update the rc.local to start the jar application automatically. The rc.local is a script file that is called to run set commands once all the services in a Linux shell have been started. **Figure 6: Run the java file in AWS console**



```
couponservice-0.0.1-SNAPSHOT.jar                    100%[==============================================

2023-10-21 22:37:34 (62.7 MB/s) - 'couponservice-0.0.1-SNAPSHOT.jar' saved [38880993/38880993]

[root@ip-172-31-25-142 ec2-user]# ls
couponservice-0.0.1-SNAPSHOT.jar
[root@ip-172-31-25-142 ec2-user]# java -jar couponservice-0.0.1-SNAPSHOT.jar
```

**Figure 7: Installation of Rc.local and make it active**



```
[root@ip-172-31-25-142 ~]# systemctl daemon-reload
[root@ip-172-31-25-142 ~]# systemctl stop rc-local.service
[root@ip-172-31-25-142 ~]# systemctl enable rc-local.service
[root@ip-172-31-25-142 ~]# systemctl start rc-local.service
[root@ip-172-31-25-142 ~]# systemctl status rc-local.service
• rc-local.service - /etc/rc.local Compatibility
    Loaded: loaded (/etc/systemd/system/rc-local.service; enabled; prese
    Active: active (running) since Sun 2023-10-22 12:00:05 UTC; 16s ago
```

Install amazon-linux extras and install stress package (package to simulate load)

**Figure 8: Install amazon-linux extras and install stress package (package to simulate load)**



```
  _/m/'
ast login: Sun Oct 22 17:41:39 2023 from 13.48.4.202
ec2-user@ip-172-31-23-33 ~]$ sudo su
[root@ip-172-31-23-33 ec2-user]# sudo dnf install stress -y
ast metadata expiration check: 8:33:24 ago on Sun Oct 22 09:18:30 2023.
Dependencies resolved.
===============================================================================
 Package                          Architecture                    Version
===============================================================================
```

## Creating an Amazon Machine Image (AMI)

is a process that involves capturing a snapshot of an Amazon EC2 instance. An AMI allows you to launch new EC2 instances with the same configuration and data as the original instance. Here are the steps to create an AMI using the AWS Management Console:

Using the AWS Management Console:

Sign into the AWS Management Console:

Launch an EC2 Instance: If you don't already have an EC2 instance running, you'll need to launch one.

Select the EC2 Instance: sIn the EC2 Dashboard, select the EC2 instance for which you want to create an AMI.

**Create an AMI**:

  Right-click on the instance or use the "Actions" dropdown menu, and then select "Create Image." Alternatively, you can select the instance and click the "Create Image" button from the top menu.

**Figure 9: creation for AMI image from an instance**



Configure the AMI: In the "Create Image" dialog box, you can provide the following information.

Name: Give a descriptive name to your AMI.

Description: Add an optional description for your AMI.

**Figure 10: AMI image name**



create the AMI   After configuring the AMI settings, click the "Create Image" button. This will initiate the AMI creation process.

Monitor the AMI Creation: You can monitor the status of the AMI creation in the "AMI" section under the "Images" category in the EC2 Dashboard.

Use the AMI once the AMI is created successfully, you can use it to launch new EC2 instances with the same configuration as the original instance. While creating an AMI, the instance may be temporarily stopped, and the process might take some time, especially if the instance has a large amount of data. Additionally, the storage costs for the AMI snapshots will be incurred. After creating the AMI, you can launch new instances from it, and they will have the same state as the original instance at the time of the AMI creation.

**Figure 11: Properties**



**Figure 12 testing the Coup Image using postman**



## Auto Scalling Group  Create launch template

- Under Launch template name and description, enter a name and description for the new launch template.

- Under Key pair (login), for Key pair name, choose the name of the previously created key pair to use when connecting to instances, for example, using SSH.

- Under Network settings, for Security groups, choose one or more previously created security groups.

**Figure 13:  Launch template details**



- Under Configure storage, update the storage configuration. The default storage configuration is determined by the AMI and the instance type.
- choose Create Auto Scaling group.

**Figure 14: image of the launch template**



## Create an Auto Scaling group

- On the Choose launch template or configuration page, enter a name for the Auto Scaling group. I.e couponserviceasg for group and couponconfig for template

**Figure 15:  create an auto scaling group**



- The launch template that you created is already selected for you.
  - For Launch template version, choose whether the Auto Scaling group uses the default, the latest, or a specific version of the launch template when scaling out.
- Choose Next to continue to the next step.
- On the Choose instance launch options page, if you're not using multiple instance types, you can skip the Instance type requirements section to use the EC2 instance type that is specified in the launch template.
- Under Network, for VPC, choose a VPC. The Auto Scaling group must be created in the same VPC as the security group you specified in your launch template.

**Figure 16: configuring security groups**



○ For Availability Zones and subnets, choose one or more subnets in the specified VPC.

**Figure 17: configuring subnets**



To go to the Configure group size and scaling policies page.

**Figure 18: scaling policies configuration**



○ Under Group size, define the Desired capacity initial number of instances to launch immediately after the Auto Scaling group is created, Minimum

capacity minimum number of instances, and Maximum capacity maximum number of instances.

**Figure 19: Group size configuration**



- Choose Skip to review.
- On the Review page, choose Create Auto Scaling group.

**Figure 20:  Auto scaling instance**



Testing your auto scaling group policy: to test the auto scaling, you must remote into an instance of the EC2 that has been instantiated by the auto scaling and running and execute a stress command. Instances created because of an auto scaling policy are unnamed or rather identify using a dash symbol. A demonstration of an auto scaling group policy is shown in the steps below

- click on an instantiated EC2 and copy the IPv4 address to securely remote into the instance
- **Figure 21: IP address of Autoscaling instance**

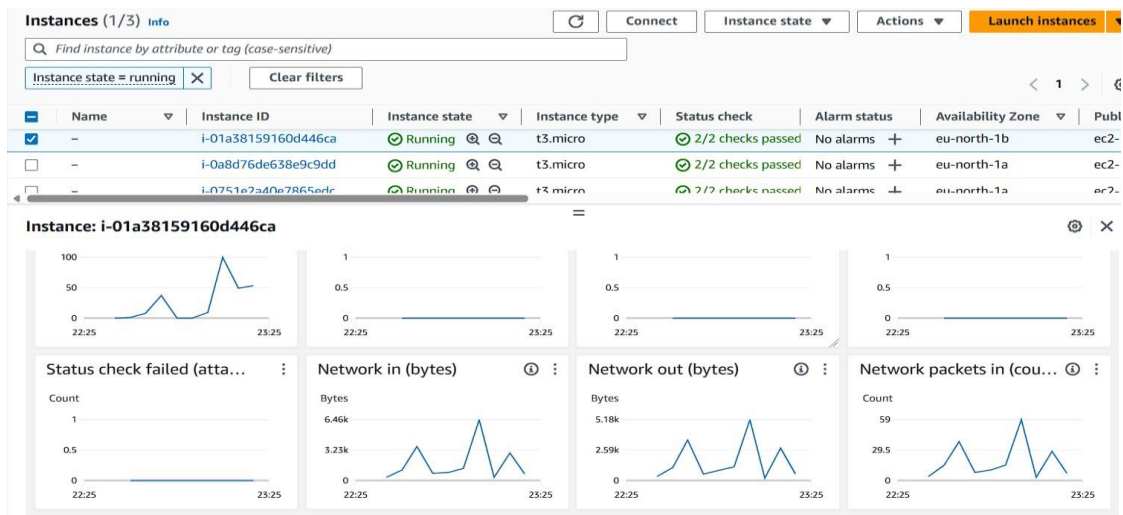**Figure 22: testing the Autoscaling instance through postman**



- ⭕ using a remote client application login into the EC2 instance execute a stress command to stress the system cpu and triggered the instantiation of instances as per the group policy in place. stress - - cpu 1
- ⭕ **Figure 23 using Linux command to stress instance in the AWS console**



allow the number of seconds you setup in your configuration, and check to see that the number of instances running have increased due to the stress command executed

**Figure 24: Monitoring stress and the new instances scaling**



Test that the jar application is also running on the new instance. This is done using "postman", a software application for testing REST API connections. Use the IPv4 address of the EC2 instance and test the application. After the GET command is send a RESPONSE of 200 is received, it indicates that the application is running perfectly.

Stop the stress command and view your instance. After the stress command has been issued, the instances that were created because of the auto scaling group policy gets terminated. When an instance is terminated, it is deleted from the ESB and it leaves the list of instances after a while. To terminate the stress command, press CTRL + C

**CREATING A LOAD BALANCER**

- On the navigation pane, under Load Balancing, choose Load Balancers.
- Choose Create Load Balancer.
- Expand the Classic Load Balancer section, then choose Create.
- Basic configuration
    - For Load balancer name, type a name for your load balancer.

    The name of your Classic Load Balancer must be unique within your set of Classic Load Balancers for the Region, can have a maximum of 32 characters, can contain only alphanumeric characters and hyphens, and must not begin or end with a hyphen.

    - For Scheme, select Internet-facing



**Figure 25: create a load balancer**

- Network mapping
    - For VPC, select the same VPC that you selected for your instances.
    - For Mappings, first select an Availability Zone, then choose a public subnet from its available subnets. You can only select one subnet per Availability Zone. To improve the availability of your load balancer, select more than one Availability Zone and subnet.
- Security groups

    For Security groups, select an existing security group that is configured to allow the required HTTP traffic on port 80.

- Listeners and routing
  - For Listener, ensure the protocol is HTTP and the port is 80.
  - For Instance, ensure the protocol is HTTP and the port is 9091 for the couponservice.

**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 ▼ | C

| default ✕ | launch-wizard-12 ✕ |
| sg-0c9ac2243352aa2d4   VPC: vpc-05fd7f68171bf37ef | sg-0cc288fd52d79d715   VPC: vpc-05fd7f68171bf37ef |

**Listeners and routing** Info

A listener is a process that checks for connection requests using the protocol and port you configure. The settings you define for a listener determine how the load balancer routes requests to its registered targets.

▼ Listener **HTTP:80**    Remove
   Instance HTTP:9091

| Listener protocol | Listener port | | Instance protocol | Instance port |
| HTTP ▼ | : | 80 | HTTP ▼ | : | 9091 ⬍ |
| | 1-65535 | | | 1-65535 |

**Figure 26: security group configuration**

- Health checks
  - For Ping Protocol, ensure the protocol is HTTP.
  - For Ping Port, ensure the port is 80.
  - For Ping Path, ensure the path is /couponapi/coupons/SUMMERSALE
  - For Advanced health check settings, use the values you desire but I recommend default

**2. Figure 27 Advance health check settings configuration**

**Ping target**

The health check ping is sent using the protocol and port you specify. If using HTTP/HTTPS protocol, you m path.

| Ping protocol | Ping port | Ping path |
| HTTP ▼ | : | 9091 | /couponapi/coupons/SUMMERSALE |
| | 1-65535 | |

▼ **Advanced health check settings**

**Response timeout**
Time to wait for EC2 instances to respond to health checks.

5    seconds

2-60 seconds. Must be less than the health check interval.

**Interval**
Amount of time between health checks sent to EC2 instances.

30    seconds

5-300 seconds. Must be greater than the health check response timeout.

**Unhealthy threshold**
Number of consecutive health check failures before declaring an EC2 instance unhealthy.

2 ▼

2-10

**Healthy threshold**
Number of consecutive health check successes before declaring an EC2 instance healthy.

10 ▼

2-10

- Instances
    - Select Add instances, to bring up the instance selection screen.
    - Under Available instances, you can select from the current instances that are available to the load balancer, based on the current network settings. In my case is coup service and cup _image.
    - select Confirm to add the instances to be registered to the load balancer

**Figure28: Choosing Instances from the coupon load balancer**



- Attributes
    - For Enable cross-zone load balancing, enable connection draining, and Timeout   keep the default values if you want.
- Load balancer tags (optional)
    - To add another tag, select Add new tag then input your values into the Key field, and optionally the Value field.
    - To remove an existing tag, select Remove next to the tag you want to remove.
- Summary and creation
    - If you need to change any settings, select Edit next to the setting needing to be changed.
    - After you're satisfied with all the settings shown in the summary, select Create load balancer to begin creation of your load balancer.
    - On the final creation page, select View load balancer to view your load balancer in the Amazon EC2 console.

**Figure 29:  load balancer summary**



- Verify
  - Select your new load balancer.
  - On the Target instances tab, check the health status column. After at least one of your EC2 instances is In-service, you can test your load balancer.
  - In the Details section, copy the load balancers DNS name, which would look like coupload-347727001.eu-north-1.elb.amazonaws.com
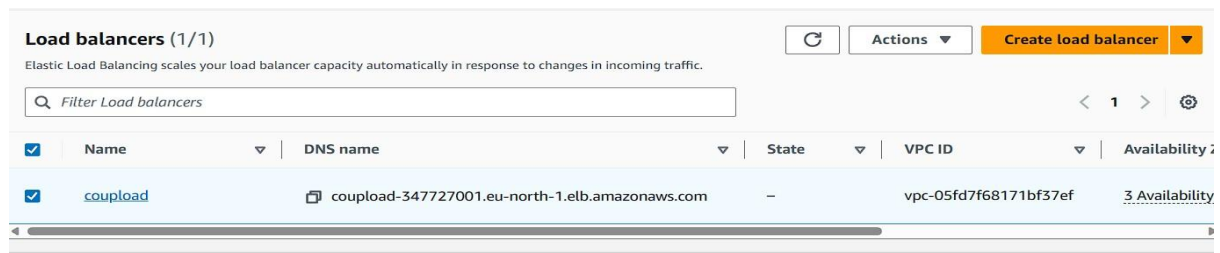
**Figure 30 : load balancer details**

**Figure 31:The load balancer**



To see the difference in the load balancer we need to change the database(mariadb) from one instance in the AWS.

**Figure 31: updating database of one instance**



Paste your load balancers DNS name into the address field of a public internet connected web browser.

If your load balancer is functioning correctly, you will see the default page of your server.

Or copy the DNS name to postman to get the 200ms ok.

Attach the load balancer to the Auto scaling group  his helps you keep up with the demand from your users while keeping costs down. Amazon's ELB ties in well with auto scaling groups automatically load balancing servers that have been brought up as the result of an auto scaling action. **Figure 32: attaching load balancer to the autoscaling group**

by attaching the load balancer to the auto scaling group next time when the EC2 are stress the auto scaling instance will launch from the load balancer instance see on the figure bellow

**Figure 33: shows auto scaling instance launching in the load balancer**

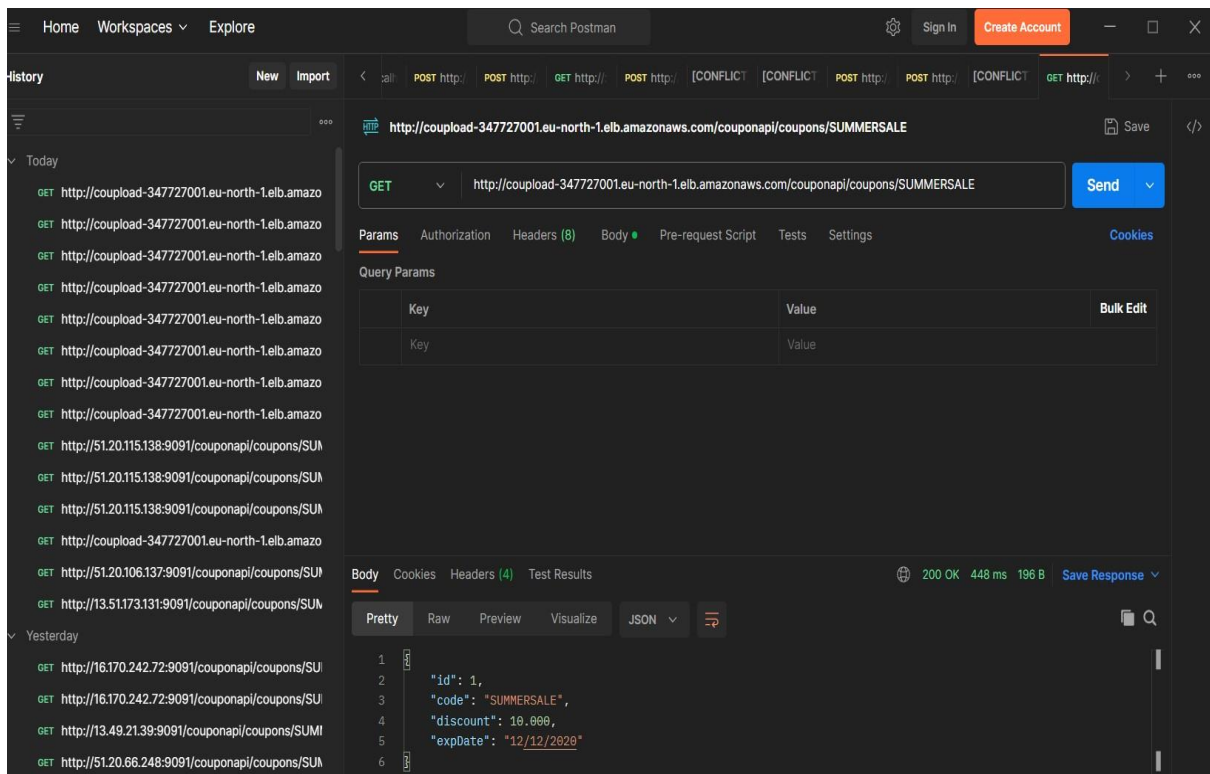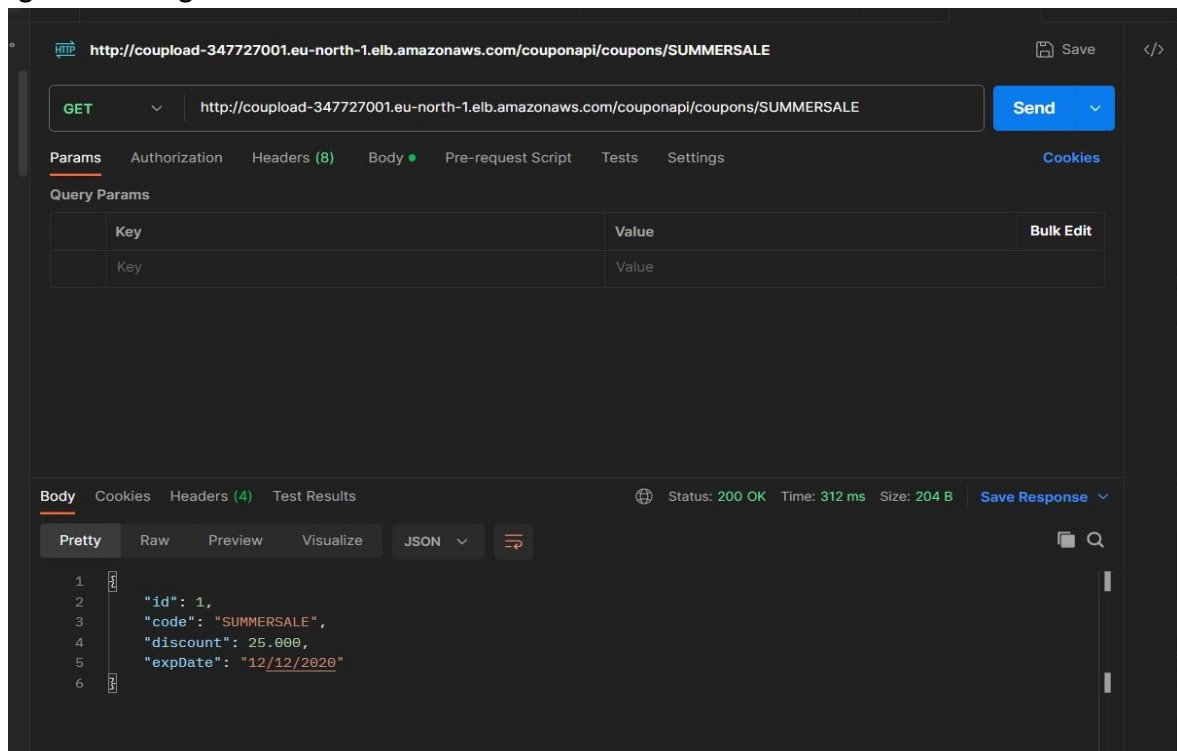**Figure 32: testing the load balancer from one instance**



**Figure 3: testing the load balancer from another instance**

**Evaluation**

After I created an instance, connected to it. Installed mariadb, gave it the password created the databases from mariadb gave the tables (coupon and product). We later installed java in that instance.

Created a s3 bucket where we uploaded the java file made it accessible publicly and went back to the instance wget the java file then run the java file run it. Copied the EC2 public address and the couponapi/cooupon posted the summersales to postman app for test.

Installed rc.local to enable the jar file automatically when the instance is active.

Create an AMI from the instance and get the image . run the images. Open auto scaling group of the two instance. Apply a stress package from one the instance load The EC2 scaled horizontally and duplicate to the desired cpu parameter to help up with the load.

Create an load balancing template add the template to the auto scaling on the network group use the existing network group from the instance. After finishing up creating load ballancing group,go to one off the EC2 and change the mariadb

Copy the load balancer URL to the postman and change on the EC2 observe the tables to see the deference.

AWS Auto Scaling service can be considered as regionally limited, such that, it is only effective in one region that are selected during the setup of group policies and it's not possible to use auto scaling policies across resources in different regions. It may become a challenge setting up auto scaling policies for multi-regions applications, therefore AWS auto scaling policies will need to be set up separately for each region

The auto-scaling worked efficiently adjusted the number of instances to handle variable workloads, ensuring optimal performance and Auto scaling effectively maintained the desired performance levels by adding or removing instances as needed hence completing all the task required.

**key limitations**

Definining Scaling Policies: Configuring effective scaling policies can be challenging. You need to set thresholds that trigger scaling actions, and choosing the right thresholds can require careful monitoring and fine-tuning. I.e cpu percentage or bytes.

Testing Auto Scaling:To demonstrate auto scaling in operation, you can simulate increased traffic to your application. This can be done by using load testing tools, sending additional requests manually, or creating a custom script to generate traffic likes linux stress test or other script languages like python script.

Observe Auto Scaling in Action:As the load on your application increases and breaches the alarm thresholds, you should see the Auto Scaling Group automatically launch new instances based on your scaling policies time i.e 300 sec you will have to wait for that time to see the outcome and another instance can be added or existing one can be stopped.

Scaling Speed: AWS Auto Scaling may not always respond instantly to rapid changes in traffic. There can be a delay in launching new instances, which can lead to temporary performance issues during sudden traffic spikes.

cost Management: Without careful monitoring and adjustment of your scaling policies, AWS Auto Scaling can lead to increased costs, while using free tier especially if it scales up resources unnecessarily.

**References**

1. Amazon Glossary of Terms

Available at: https://docs.aws.amazon.com/general/latest/gr/glos-chap.html [cited 25 October 2023]

2. Developing Cloud Native Applications : Module 1 Glossary: Introduction to Cloud Native Available at: https://learning.edx.org/course/course-v1:IBM+CC0101EN+2T2021/home [cited 25October 2023]  3. Amazon Auto Scaling

 Available at: https://docs.aws.amazon.com/autoscaling/ec2/userguide/what-is-amazonec2-autoscaling.html [cited 25 October 2023}

4.Intergrating Load balancer with auto scaling groups

Available at: https://www.loadbalancer.org/blog/integrating-your-load-balancer-with-auto-scalinggroups-in-ec2[cited 25 October 2023}

5 creating an EC2 Instance with Lambda in AWS

Available at: https://www.pluralsight.com/cloud-guru/labs/aws/creating-an-ec2-instance-withlambda-in-AWS[cited 25 October 2023}

6 Amazon Load balancing

Available at: https://aws.amazon.com/elasticloadbalancing/[cited 25 October 2023}

7 Postman

Available at: https://linuxize.com/post/how-to-install-postman-on-ubuntu-20-04/[cited 25 October 2023}