

DevOps

LOAD-BALANCING AND AUTO-SCALING WEB APP IMPLEMENTATION



| | |
|-----------------|--|
| Project Title: | Load-Balanced Auto-Scaling Web Application |
| Student Name: | Stephen Power |
| Student Number: | 20093364 |
| Start Date: | November 17, 2022 |

| | |
|--------------|---|
| Description: | This document reports the implementation of the load-balanced auto-scaling web application. The objective of this assignment is to deploy the web application and to automate it’s management by using load-balancing and auto-scaling through AWS. |
|--------------|---|

Contents

Architecture Diagram

1. Initial Setup

- 1.1 Amazon Machine Image (AMI) 3
- 1.2 Virtual Private Cloud (VPC) 3
- 1.3 Load Balancer 4
- 1.4 Auto-Scaling 5
- 1.5 Scaling Policies 5

2. Deployment of Application

- 2.1 Web Application..... 6
- 2.2 Distributed Load..... 7
- 2.3 Custom Metrics 7

3. Additional Functionality

- 3.1 Secure Load Balancer 8
- 3.2 Simple Queue Service (SQS)..... 9

4. Conclusion

Architecture Diagram

This diagram represents the architecture of the assignment. A virtual private cloud (VPC) is set up as the host. An application load balancer is created within the VPC in order to balance the load over EC2 instances. An auto-scaling group is set up to configure the minimum, maximum and desired amount of instances. Web server applications are launched on the EC2 instances.

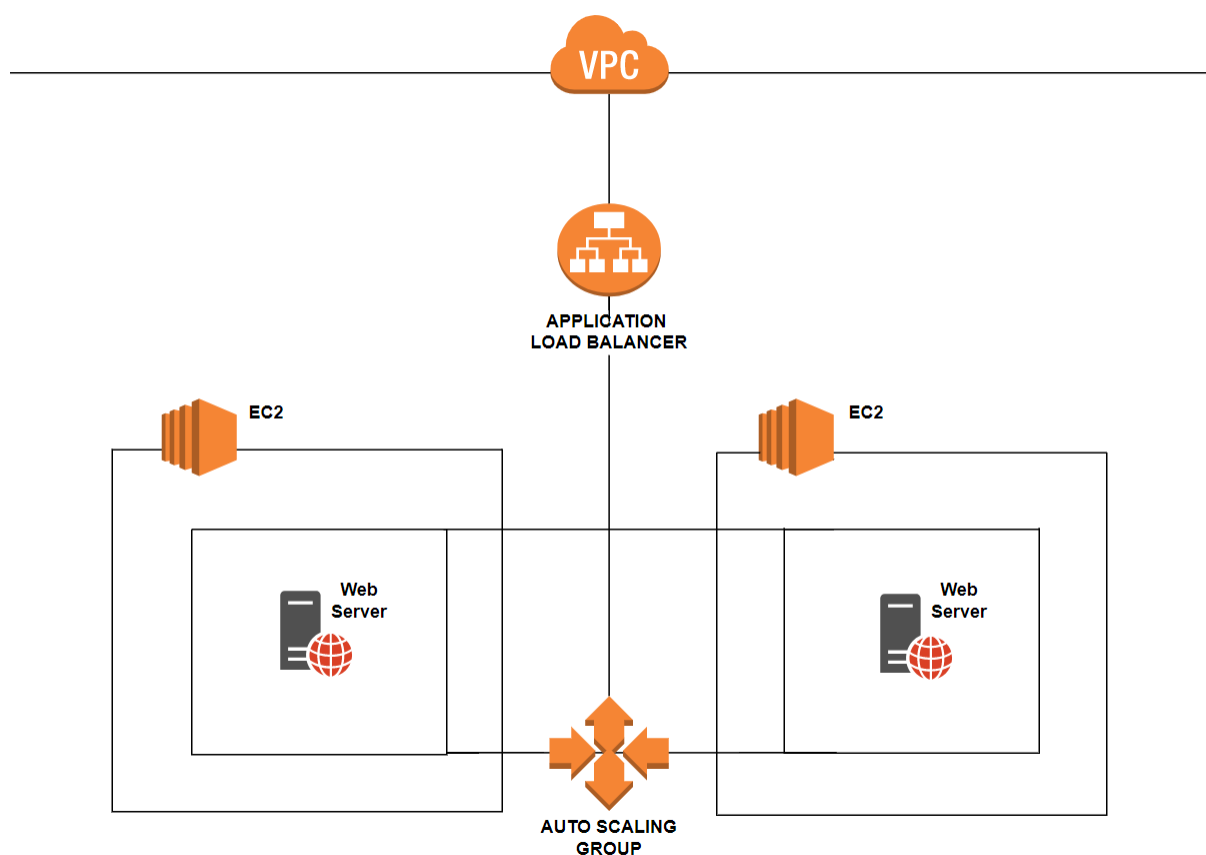


Figure 1: Architecture Diagram

1. Initial Setup

1.1 Amazon Machine Image (AMI)

A custom AMI was set up with the name 'Master web server'. The setting up of a custom AMI is required as an AMI is needed to launch an EC2 instance. This custom AMI is to be used for the master instance and auto-scaling. Once created, the image took a few minutes to become available.

| | | | |
|--|---|--|---------------------------------|
| AMI ID: ami-00fce4317adee2a41 | | | |
| | | | |
| Details | Permissions | Storage | Tags |
| | | | |
| AMI ID ami-00fce4317adee2a41 | Image type machine | Platform details Linux/UNIX | Root device type EBS |
| AMI name Master web server | Owner account ID 886086313441 | Architecture x86_64 | Usage operation RunInstances |
| Root device name /dev/xvda | Status Available | Source 886086313441/Master web server | Virtualization type hvm |
| Boot mode - | State reason - | Creation date Thu Nov 17 2022 14:31:09 GMT+0000 (Greenwich Mean Time) | Kernel ID - |
| Block devices /dev/xvda=snap-078afbdfb8f971533:8:true:gp2 | Description My own image | Product codes - | RAM disk ID - |
| Deprecation time - | Last launched time Thu Nov 17 2022 16:19:04 GMT+0000 (Greenwich Mean Time) | | |

Figure 2: Custom AMI

1.2 Virtual Private Cloud (VPC)

A custom VPC was created under the name 'vpc-stephenpower'. The VPC contained 4 subnets which would be used for deploying the application. The subnets were: private1-us-east-1a, private2-us-east-1b, public1-us-east-1a, public2-us-east-1b. Appropriate security groups were also set up in the VPC to allow SSH and HTTP traffic.

| | | | |
|--|---|---|---|
| vpc-04dd8d6c8b7403a7b / vpc-stephenpower-vpc | | | |
| Actions ▼ | | | |
| Details Info | | | |
| VPC ID vpc-04dd8d6c8b7403a7b | State Available | DNS hostnames Enabled | DNS resolution Enabled |
| Tenancy Default | DHCP option set dopt-06d34ae830d11ad27 | Main route table rtb-0620108523c05750a | Main network ACL acl-0fa08857d61ca6088 |
| Default VPC No | IPv4 CIDR 10.0.0.0/16 | IPv6 pool - | IPv6 CIDR (Network border group) - |

Figure 3: Custom VPC

| <input type="checkbox"/> | Name ▾ | Subnet ID ▾ | State ▾ | VPC ▲ |
|--------------------------|---|--------------------------|-------------|-------------------------------|
| <input type="checkbox"/> | vpc-stephenpower-subnet-private2-us-east-1b | subnet-03fcbca6eb73a17f5 | ✔ Available | vpc-04dd8d6c8b7403a7b vp... |
| <input type="checkbox"/> | vpc-stephenpower-subnet-private1-us-east-1a | subnet-0de6e9052b57752b4 | ✔ Available | vpc-04dd8d6c8b7403a7b vp... |
| <input type="checkbox"/> | vpc-stephenpower-subnet-public1-us-east-1a | subnet-019177157ec3cd488 | ✔ Available | vpc-04dd8d6c8b7403a7b vp... |
| <input type="checkbox"/> | vpc-stephenpower-subnet-public2-us-east-1b | subnet-04c4be72158657be3 | ✔ Available | vpc-04dd8d6c8b7403a7b vp... |

Figure 4: VPC - Subnets


| Security group ID ▾ | Security group name ▾ | VPC ID ▲ | Description ▾ |
|----------------------|-----------------------|-----------------------|----------------------------|
| sg-071be317e5973d669 | launch-wizard-2 | vpc-04dd8d6c8b7403a7b | launch-wizard-2 create... |
| sg-0ddf0d6c0ba26befc | default | vpc-04dd8d6c8b7403a7b | default VPC security gr... |
| sg-053bdc408657a1aff | WebServerSG | vpc-04dd8d6c8b7403a7b | For the web servers in ... |
| sg-0827bc18c00966472 | DBServerSG | vpc-04dd8d6c8b7403a7b | For the database serve... |

Figure 5: VPC - Security Groups

1.3 Load Balancer

An application load balancer was set up with the name 'load-balancer-sp'. The load balancer was mapped with the custom VPC so the public subnets could be used. A default target group was set up in the custom VPC under the name 'target-group-sp' and was selected as the designated target group for the load balancer.

load-balancer-sp

 Actions ▾

▼ Details

arn:aws:elasticloadbalancing:us-east-1:886086313441:loadbalancer/app/load-balancer-sp/e6789c5e91b1033f

| | | | |
|---|--|--|--|
| Load balancer type Application Load Balancer | DNS name load-balancer-sp-1715940905.us-east-1.elb.amazonaws.com (A Record) | Status Active | VPC vpc-04dd8d6c8b7403a7b |
| IP address type IPv4 | Scheme Internet-facing | Availability Zones subnet-019177157ec3cd488 us-east-1a (use1-az6) subnet-04c4be72158657be3 us-east-1b (use1-az1) | Hosted Zone Z35SXDOTRQ7X7K |

Figure 6: Load Balancer

target-group-sp

Actions ▾

Details

arn:aws:elasticloadbalancing:us-east-1:886086313441:targetgroup/target-group-sp/dae6e63c9dccc4a7

| | | | |
|-------------------------|---|---------------------------|--|
| Target type Instance | Protocol : Port HTTP: 80 | Protocol version HTTP1 | VPC vpc-04dd8d6c8b7403a7b |
| IP address type IPv4 | Load balancer load-balancer-sp | | |

Figure 7: Target Group

1.4 Auto-Scaling

A launch template was configured with the custom AMI. The launch template was then used to create an auto-scaling group under the name 'auto-scaling-sp'. The group size was configured to have a minimum capacity of 1 and a maximum of 2. The public subnets chosen for the load balancer were also configured for the auto-scaling group.

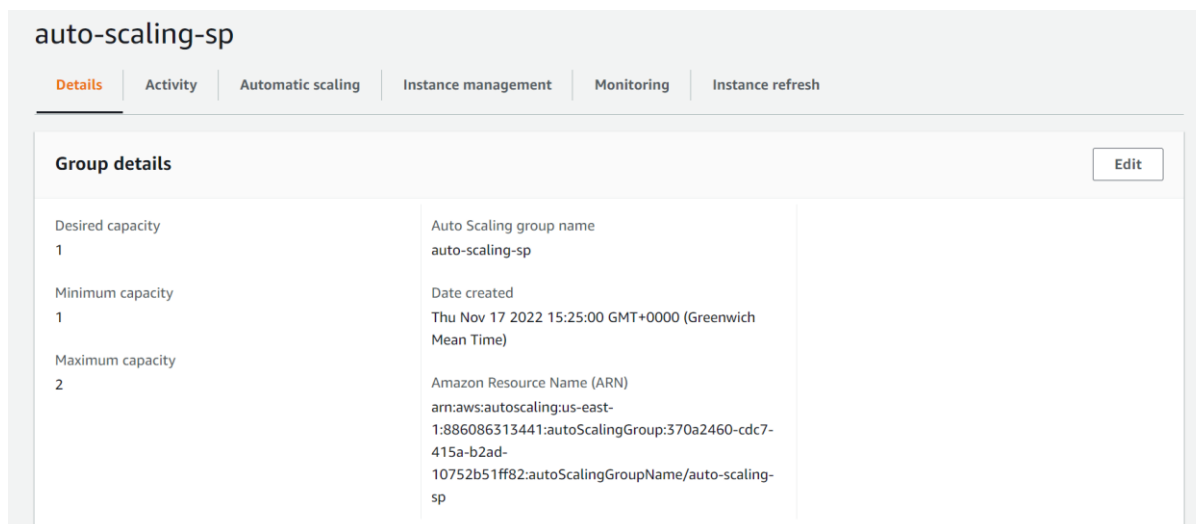


Figure 8: Auto Scaling Group

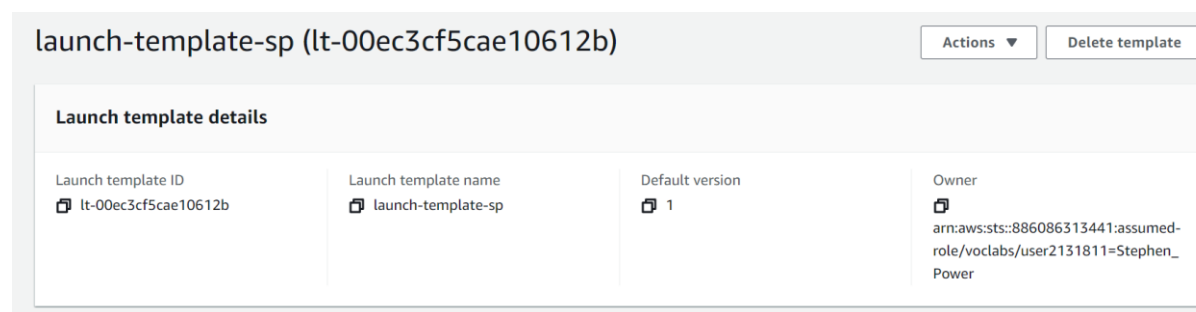


Figure 9: Launch Template

1.5 Scaling Policies

Two dynamic scaling policies were set up for the auto-scaling group:

1. Scale in on low CPU
2. Scale out on high CPU

Both policies used simple scaling based on CloudWatch alarms. The alarms were triggered using the CPU utilisation metric. If the CPU utilisation is greater than 50, the high CPU alarm will trigger and the 'Scale out on high CPU' policy will activate. Conversely, if the CPU utilisation is equal to or less than 50, the low CPU alarm will activate the 'Scale in on low CPU' policy. These scaling policies were chosen as when CPU utilisation is over 50, more resources are needed.

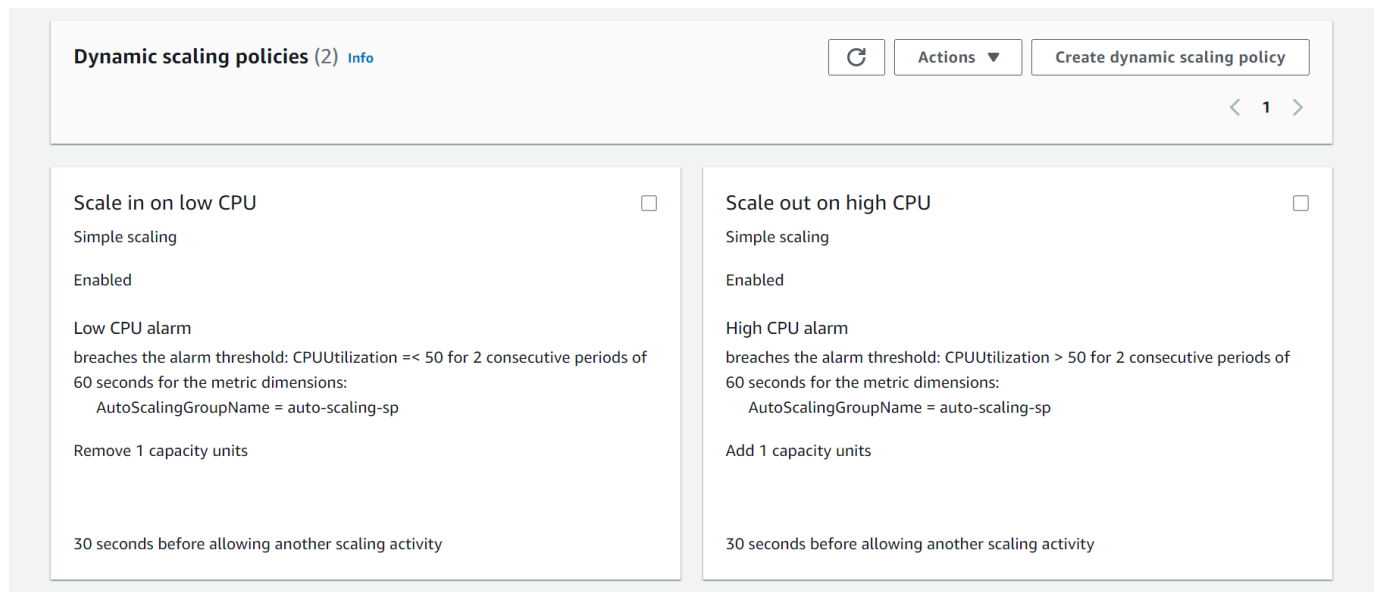


Figure 10: Scaling Policies

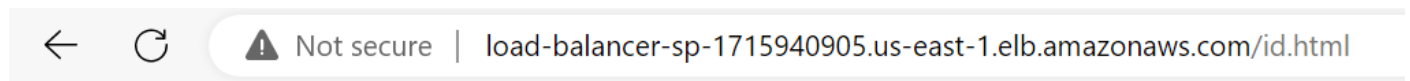
| | | | | | |
|--------------------------|----------------|----------|---------------------|--|-----------------|
| <input type="checkbox"/> | High CPU alarm | OK | 2022-11-23 11:04:29 | CPUUtilization > 50 for 1 datapoints within 2 minutes | Actions enabled |
| <input type="checkbox"/> | Low CPU alarm | In alarm | 2022-11-23 11:03:47 | CPUUtilization <= 50 for 1 datapoints within 2 minutes | Actions enabled |

Figure 11: CloudWatch Alarms

2. Deployment of Application

2.1 Web Application

A web application was launched on completion of the initial setup. The web application could be viewed by copying the DNS name of the load balancer and inputting it into a web browser search. The web application simply contained the ID of the EC2 instance. The application was dynamic as if the instance changed, the application would also change to display the current instance ID.



Instance ID: i-0672d2c821837aeb0

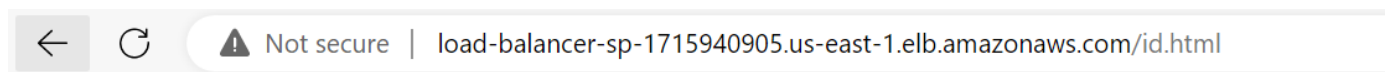
Figure 12: Web Application

2.2 Distributed Load

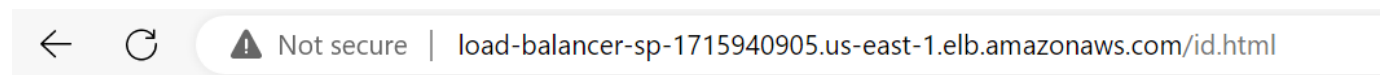
A simple infinite while loop (while true; do x=0; done) was placed on an instance through an SSH connection to show that the load is not just reliant on one server. Once the loop was initiated, the load balancer scaled out and another server was automatically set up. Both servers could be seen running using the 'top' command in the SSH connection, as well as displaying both instance id's on the web application.

| PID | USER | PR | NI | VIRT | RES | SHR | S | %CPU | %MEM | TIME+ | COMMAND |
|------|----------|----|----|--------|------|------|---|------|------|---------|---------|
| 3273 | ec2-user | 20 | 0 | 124740 | 3892 | 2972 | R | 99.9 | 0.8 | 2:45.39 | bash |
| 3354 | ec2-user | 20 | 0 | 168992 | 4480 | 3904 | R | 0.3 | 0.9 | 0:00.03 | top |

Figure 13: Multiple Instances



Instance ID: i-0c60a547096b3d9a7



Instance ID: i-0c4a19da2b5e610be

Figure 14: Web Page (Multiple Instances)

2.3 Custom Metrics

A bash script was used to monitor custom metrics on the EC2 instances. The script monitored the following metrics: percentage of used memory, the number of total and port 80 TCP connections, the percentage of I/O wait time, the number of HTTPD processes running and the total number of processes running. A cron job was created in order to push these custom metrics to CloudWatch.

| Metrics (6) Info | | Graph with SQL | Graph search |
|---------------------------------------|-----------------------|---|------------------------------|
| N. Virginia ▾ All > Custom > Instance | | 🔍 Search for any metric, dimension or resource id | |
| <input type="checkbox"/> | Instance (6) | ▲ | Metric name ▲ |
| <input type="checkbox"/> | i-031ec11a12c8b0573 ▾ | | IO_WAIT ▾ |
| <input type="checkbox"/> | i-031ec11a12c8b0573 ▾ | | TCP_connection_on_port_80 ▾ |
| <input type="checkbox"/> | i-031ec11a12c8b0573 ▾ | | memory-usage ▾ |
| <input type="checkbox"/> | i-031ec11a12c8b0573 ▾ | | Tcp_connections ▾ |
| <input type="checkbox"/> | i-031ec11a12c8b0573 ▾ | | processes ▾ |
| <input type="checkbox"/> | i-031ec11a12c8b0573 ▾ | | httpd-processes ▾ |

Figure 15: Custom Metrics


```
#!/bin/bash
INSTANCE_ID=$(curl -s http://169.254.169.254/latest/meta-data/instance-id)
USEDMEMORY=$(free -m | awk 'NR==2{printf "%.2f\t", $3*100/$2 }')
TCP_CONN=$(netstat -an | wc -l)
TCP_CONN_PORT_80=$(netstat -an | grep 80 | wc -l)
IO_WAIT=$(iostat | awk 'NR==4 {print $4}')
PROCESSES=$(expr $(ps -A | grep -c .) - 1)
HTTPD_PROCESSES=$(ps -A | grep -c httpd)

aws cloudwatch put-metric-data --metric-name memory-usage --dimensions Instance=$INSTANCE_ID --namespace "Custom" --value $USEDMEMORY
aws cloudwatch put-metric-data --metric-name Tcp_connections --dimensions Instance=$INSTANCE_ID --namespace "Custom" --value $TCP_CONN
aws cloudwatch put-metric-data --metric-name TCP_connection_on_port_80 --dimensions Instance=$INSTANCE_ID --namespace "Custom" --value $TCP_CONN_PORT_80
aws cloudwatch put-metric-data --metric-name IO_WAIT --dimensions Instance=$INSTANCE_ID --namespace "Custom" --value $IO_WAIT
aws cloudwatch put-metric-data --metric-name processes --dimensions Instance=$INSTANCE_ID --namespace "Custom" --value $PROCESSES
aws cloudwatch put-metric-data --metric-name httpd-processes --dimensions Instance=$INSTANCE_ID --namespace "Custom" --value $HTTPD_PROCESSES
```

Figure 16: Monitoring Script

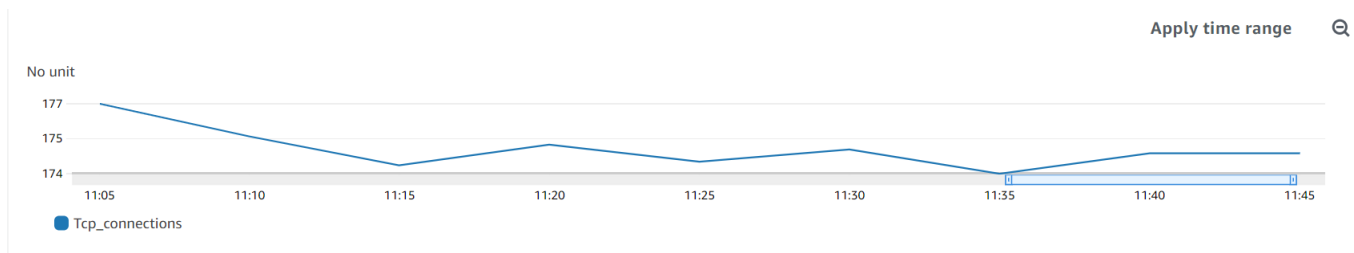


Figure 17: TCP Connections Monitoring (CloudWatch)

3. Additional Functionality

3.1 Secure Load Balancer

The load balancer was initially set up to accept HTTP traffic (port 80). Additionally, a self-signed certificate was created for the load balancer using the AWS Certificate Manager so it could now also accept HTTPS traffic (port 443). The self-signed certificate was generated using OpenSSL.

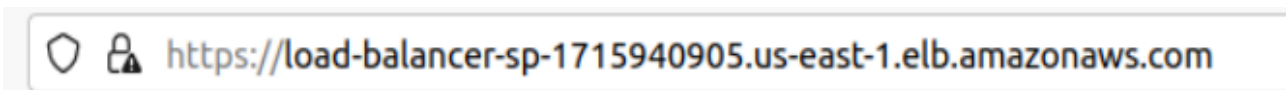


Figure 18: Secure Load Balancer

Certificate Viewer: load-balancer-sp-1715940905.us-east-1.elb.amazonaws.com

General Details

Issued To

| | |
|--------------------------|---|
| Common Name (CN) | load-balancer-sp-1715940905.us-east-1.elb.amazonaws.com |
| Organization (O) | Reliable web apps |
| Organizational Unit (OU) | <Not Part Of Certificate> |

Issued By

| | |
|--------------------------|---|
| Common Name (CN) | load-balancer-sp-1715940905.us-east-1.elb.amazonaws.com |
| Organization (O) | Reliable web apps |
| Organizational Unit (OU) | <Not Part Of Certificate> |

Validity Period

| | |
|------------|---|
| Issued On | Friday, November 25, 2022 at 3:37:56 PM |
| Expires On | Sunday, December 25, 2022 at 3:37:56 PM |

Fingerprints

| | |
|---------------------|--|
| SHA-256 Fingerprint | 46 3F 6F FC 72 E2 79 8A 83 96 38 06 D5 95 15 29 BA 68 7F 7A 95 71 8A 8A 60 BE A1 B7 15 56 20 78 |
| SHA-1 Fingerprint | 4E C0 29 DC E2 C2 D0 ED C9 EC D8 60 E2 17 30 87 FA 62 3D EF |

Figure 19: Self-Signed Certificate for the Load Balancer

3.2 Simple Queue Service (SQS)

A message queuing service was set up under the name 'SQS-SP'. Messages could be sent and received in a queue which would be used for auto-scaling. Two CloudWatch Alarms were set up:

1. High messages visible alarm
2. Low messages visible alarm

If the amount of messages in the queue was greater than 10, the high messages visible alarm would activate and if the amount of messages visible in the queue were equal to 10 or less, the low messages visible alarm would activate. Two dynamic scaling policies were set up using these alarms to carry out the auto-scaling. The auto-scaling group would scale out if the high messages visible alarm was active and it would scale in if the low messages visible alarm was active.

| | |
|--|--|
| <p>Scale out on high CPU <input type="checkbox"/></p> <p>Simple scaling</p> <p>Enabled</p> <p>High CPU alarm breaches the alarm threshold: CPUUtilization > 50 for 2 consecutive periods of 60 seconds for the metric dimensions: AutoScalingGroupName = auto-scaling-sp</p> <p>Add 1 capacity units</p> <p>30 seconds before allowing another scaling activity</p> | <p>Scale out on high messages <input type="checkbox"/></p> <p>Simple scaling</p> <p>Enabled</p> <p>High messages visible alarm breaches the alarm threshold: ApproximateNumberOfMessagesVisible > 10 for 2 consecutive periods of 60 seconds for the metric dimensions: QueueName = SQS-SP</p> <p>Add 1 capacity units</p> <p>30 seconds before allowing another scaling activity</p> |
|--|--|

Figure 20: SQS Scaling Policies

| <input type="checkbox"/> | Name | State | Last state update | Conditions |
|--------------------------|-----------------------------|------------|---------------------|--|
| <input type="checkbox"/> | Low messages visible alarm | 🟢 OK | 2022-11-25 20:05:21 | ApproximateNumberOfMessagesVisible <= 10 for 1 datapoints within 2 minutes |
| <input type="checkbox"/> | High messages visible alarm | 🔴 In alarm | 2022-11-25 20:04:19 | ApproximateNumberOfMessagesVisible > 10 for 1 datapoints within 2 minutes |

Figure 21: SQS CloudWatch Alarms


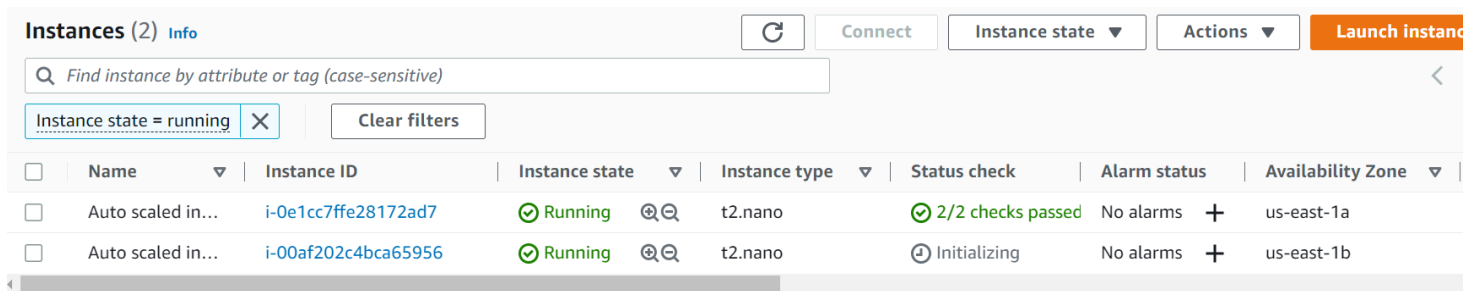
| Queues (1) | | | | |  | Edit |
|--|--------|----------|--------------------------|--------------------|---|------|
| <input type="text" value="Search queues by prefix"/> | | | | | | |
| | Name | Type | Created | Messages available | | |
| <input type="radio"/> | SQS-SP | Standard | 11/25/2022, 19:02:12 GMT | 11 | | |

Figure 22: SQS-SP (11 Messages Available)



| <input type="checkbox"/> | Name | Instance ID | Instance state | Instance type | Status check | Alarm status | Availability Zone |
|--------------------------|-------------------|---------------------|----------------|---------------|-------------------|--------------|-------------------|
| <input type="checkbox"/> | Auto scaled in... | i-0e1cc7ffe28172ad7 | Running | t2.nano | 2/2 checks passed | No alarms | us-east-1a |
| <input type="checkbox"/> | Auto scaled in... | i-00af202c4bca65956 | Running | t2.nano | Initializing | No alarms | us-east-1b |

Figure 23: SQS Scaling Out

4. Conclusion

It was clear to see the benefits of load balancing and auto-scaling through setting these up for a web application using AWS. AWS elastic load balancing distributes traffic across multiple EC2 instances, in multiple availability zones. Traffic is routed to instances based on their health status which is determined by health checks. More memory, disk space, and faster CPU speeds is a result of balancing an application on multiple instances. The goal of load balancing is to increase performance.

A group of EC2 instances can be grouped together in an auto-scaling group. The minimum, maximum, and desired amount of instances in the group can be specified. The auto-scaling group is used with the load balancer to scale in or scale out when needed. Scaling policies that use CloudWatch alarms can be set up to automatically determine when a scaling group should scale in or scale out. Auto-scaling policies were set up for EC2 and SQS metrics in this assignment which gave an insight into how these services work.

Working with different AWS services in this assignment such as VPC, EC2, CloudWatch, SQS, and Certificate Manager has resulted in a strong learning experience. Learning the theory about these services at first, but then actually working with these services practically for the assignment really helped in fully understanding how these services function and what they are used for.