

Configuration Manual

MSc Research Project
Master of Science in Cloud Computing

Ann Mariya Jojo
Student ID: x23241535

School of Computing
National College of Ireland

Supervisor: Shreyas Setlur Arun

National College of Ireland
Project Submission Sheet
School of Computing



Student Name:	Ann Mariya Jojo
Student ID:	x23241535
Programme:	Master of Science in Cloud Computing
Year:	2024
Module:	MSc Research Project
Supervisor:	Shreyas Setlur Arun
Submission Due Date:	12/12/2024
Project Title:	Configuration Manual
Word Count:	1300
Page Count:	9

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:	Ann Mariya Jojo
Date:	25th January 2025

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).	✓
Attach a Moodle submission receipt of the online project submission , to each project (including multiple copies).	✓
You must ensure that you retain a HARD COPY of the project , both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	✓

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Ann Mariya Jojo
x23241535

1 Introduction

The documents consist of each steps that needs to be followed for setting up and executing the load balancing experiment. There are two cases mainly that is traditional VM based approach and Case 2 that is Implementation of Hybrid GWO-PSO algorithm on AWS Lambda for optimization.

2 Case 1 Experiment

The following prerequisites are required before setting up the experiment.

- An AWS account with the permission to create and access the EC2 instances, load balancers, Lambda and IAM services,
- For traffic simulation the artillery is required for Case 1 and Case 2.
- Boto3, numpy, requests libraries and python 3.9 is required.

2.1 Creating EC2 Instances

For the creation of EC2 instances first navigate into EC2 dashboard in AWS account. The instance types is t2 micro and the AMI is Amazon Linux 2. The security groups are which allows inbound HTTP that is port 80 from 0.0.0.0/0 and SSH that is port 22 from your IP that is for debugging. For SSH access a Key Pair is also required.

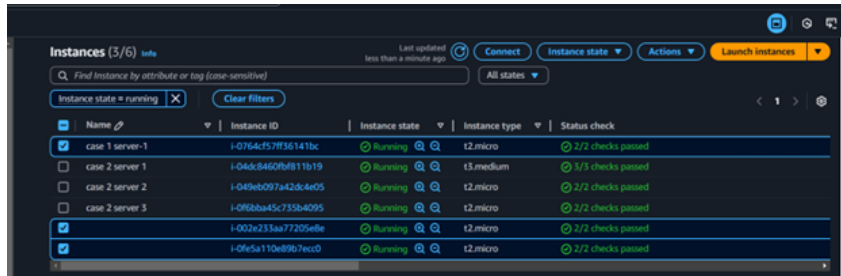


Figure 1: The EC2 instances for case 1

For installing and running the Apache Web server the commands used are given below;

- `sudo yum update -y`

- `sudo yum install -y httpd`
- `sudo systemctl start httpd`
- `sudo systemctl enable httpd`
- `echo "Welcome to Case 1 Server" — sudo tee /var/www/html/index.html`

2.2 Setting up of Application Load balancer

The scheme used is Internet facing and the listeners is port 80 that is HTTP. The target groups to add the EC2 instances which is created earlier and to configure the path for health check. Wee and Liu (2010) up. By accessing the DNS name in a browser it can be verified that Application Load balancer DNS is reachable or not.

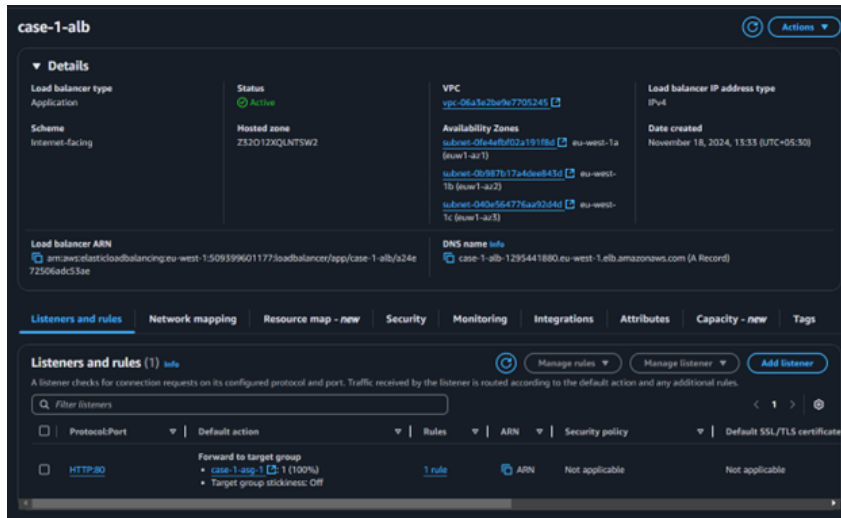


Figure 2: Setting up of ALB

2.3 Configuration of Auto scaling Group

The next step is creating the Auto scaling group for that navigate to the Auto Scaling section in AWS account. Arvindhan and Anand (2019). The target group that is created previously is attached and number of instances for minimum and maximum is also set. Make sure that auto scaling policy which is triggered is based on CPU utilization.

2.4 Simulation of traffic with the help of Artillery

Artillery is installed on the local machine using the command `npm install -g artillery@latest` and the configuration file created in the name of `case1-test.yml` and run the simulation for output. The given below is test case 1 in figure 4.

3 Case 2 Experiment

The prerequisites for setting up the Case 2 experiment is already mentioned on the prerequisites of first experiment.

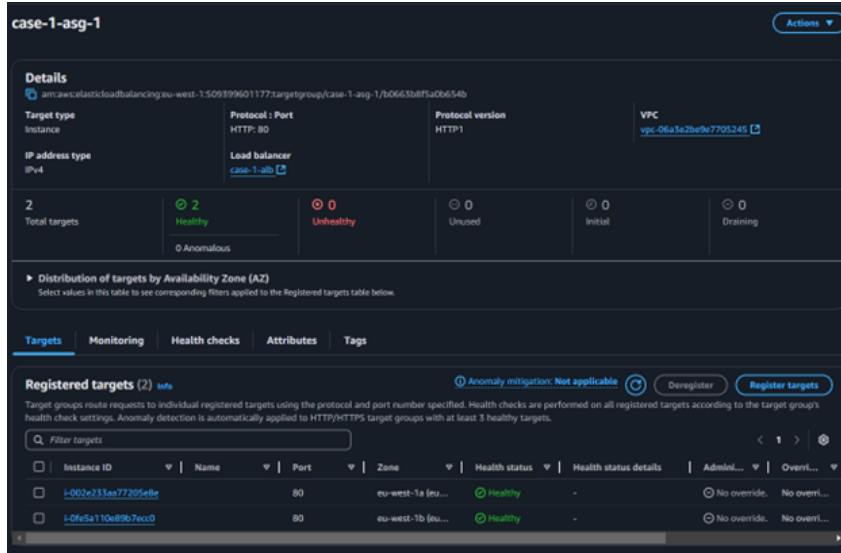


Figure 3: Auto Scaling group Configuration

```
config:
  target: "http://case-1-alb-1295441880.eu-west-1.elb.amazonaws.com"
  phases:
    - duration: 300 # Run for 5 minutes
      arrivalRate: 20 # 20 requests per second
  scenarios:
    - flow:
      - get:
        url: "/"
```

Figure 4: case1-test.yml

3.1 Creating EC2 instances

The same steps that is mentioned for Case 1 scenario is followed here also . The number of instances created is three.

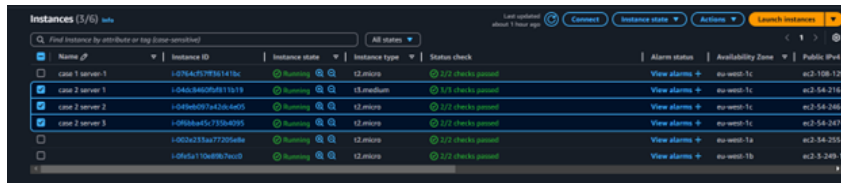


Figure 5: EC2 instances for case 2 scenario

3.2 Configuration of AWS Lambda

The Lambda function is created the runtime is python 3.9, memory is 1024 MB and the time out is 30 seconds. The IAM role attached for the function are:

- AWSLambdaBasicExecutionRole
- AmazonEC2ReadOnlyAccess

- CloudWatchFullAccess

The script of the GWO-PSO is added with the help of lambda handler and deploy the lambda function.

3.3 Application Load balancer is configured for Lambda

The application load balancer is created as mentioned in the scenario of Case 1 .The same step is followed for the Case 2 also and added a Lambda target group. Also attach the Lambda function as target.The health check path is configured and also to the Artillery test script the DNS endpoint is added.

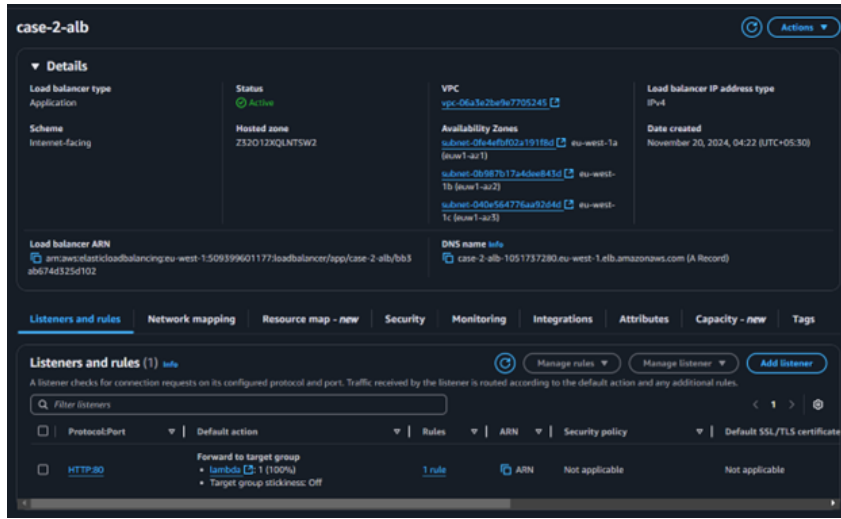


Figure 6: Case 2 ALB

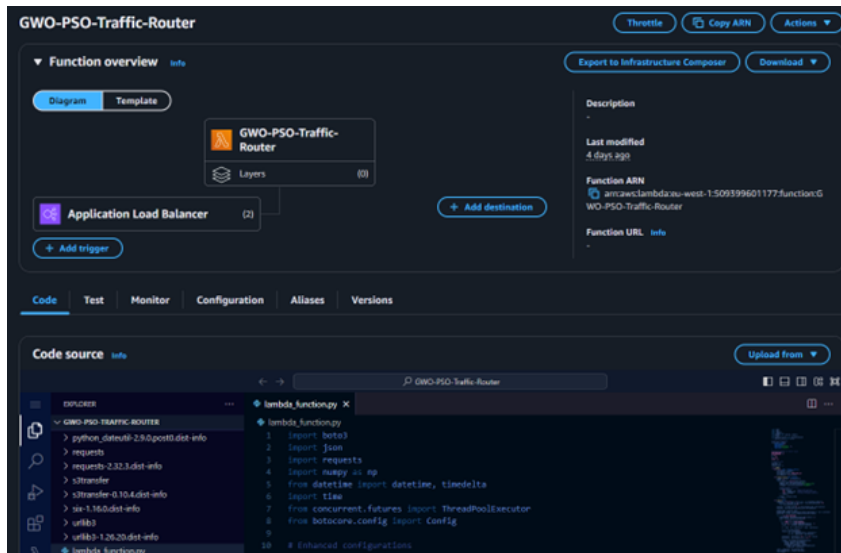


Figure 7: Lambda function of Case 2

3.4 Lambda function Testing

As mentioned the scenario in case 1. The same way yaml file is created for case 2 scenario and updated the artillery configuration file Lin et al. (2018). The name of the file case2-test.yml file. The script is given below and the simulation is executed with the help of the command artillery run case2-test.yml.

```
config:
  target: "http://case-2-alb-1051737280.eu-west-1.elb.amazonaws.com"
  phases:
    - duration: 300
      arrivalRate: 20
  scenarios:
    - flow:
      - get:
        url: "/lambda"
```

Figure 8: case2-test.yml

3.5 On Server setting up of Artillery

For running artillery the SSH into the server is designated. The Node.js and artillery is installed. The given below are the commands.

- curl -o- https://raw.githubusercontent.com/nvm-sh/nvm/v0.39.7/install.sh — bash
- source ~/.bashrc
- nvm install -lts
- npm install -g artillery@latest

The configuration files of artillery that is case1-test.yml and case2-test.yml is copied to server and test is executed with the help of commands that is given above.

3.6 IAM Roles required

The IAM role required for the EC2 instances is CloudWatchAgentServerPolicy. This policy is required mainly for monitoring. The policy required for Lambda functions are given below;

- AmazonEC2FullAccess
- AmazonEC2ReadOnlyAccess
- AWSLambdaBasicExecutionRole
- AWSLambdaVPCLambdaAccessExecutionRole
- CloudWatchFullAccess
- CloudWatchFullAccessV2

3.7 Validation of Result

From the artillery output the result is collected and the metrics mainly errors, throughput, execution time and response time are evaluated for both scenarios and it is showed in figure 7 and 8.

```
-----
Summary report @ 11:14:19(+0000)
-----
http.codes.200: ..... 6000
http.downloaded_bytes: ..... 327000
http.request_rate: ..... 20/sec
http.requests: ..... 6000
http.response_time:
  min: ..... 1
  max: ..... 42
  mean: ..... 2.5
  median: ..... 2
  p95: ..... 3
  p99: ..... 21.1
http.response_time.2xx:
  min: ..... 1
  max: ..... 42
  mean: ..... 2.5
  median: ..... 2
  p95: ..... 3
  p99: ..... 21.1
http.responses: ..... 6000
vusers.completed: ..... 6000
vusers.created: ..... 6000
vusers.created_by_name.0: ..... 6000
vusers.failed: ..... 0
vusers.session_length:
  min: ..... 2.9
  max: ..... 73.3
  mean: ..... 6.3
  median: ..... 5.2
  p95: ..... 16.6
  p99: ..... 25.3
```

Figure 9: Metrics collection of Case 1

```
All VUs finished. Total time: 5 minutes, 1 second

-----
Summary report @ 05:11:50(+0000)
-----
http.codes.200: ..... 12000
http.downloaded_bytes: ..... 790283
http.request_rate: ..... 39/sec
http.requests: ..... 12000
http.response_time:
  min: ..... 6
  max: ..... 1409
  mean: ..... 16.7
  median: ..... 16
  p95: ..... 22.9
  p99: ..... 29.1
http.response_time.2xx:
  min: ..... 6
  max: ..... 1409
  mean: ..... 16.7
  median: ..... 16
  p95: ..... 22.9
  p99: ..... 29.1
http.responses: ..... 12000
vusers.completed: ..... 6000
vusers.created: ..... 6000
vusers.created_by_name.0: ..... 6000
vusers.failed: ..... 0
vusers.session_length:
  min: ..... 26.1
  max: ..... 1716.4
  mean: ..... 288.1
  median: ..... 284.3
  p95: ..... 539.2
  p99: ..... 550.1
```

Figure 10: Metrics collection of case 2

3.8 GWO-PSO hybrid algorithm code

The given below is the code :

```
C:\Users\HP> Downloads > gwo-psy.py > lambda_handler
1  import boto3
2  import json
3  import requests
4  import numpy as np
5  from datetime import datetime, timedelta
6  import time
7  from concurrent.futures import ThreadPoolExecutor
8  from botocore.config import Config
9
10 # Configuration enhanced
11 boto3_config = Config(
12     connect_timeout=1,
13     read_timeout=1,
14     retries={'max_attempts': 1},
15     region_name='eu-west-1'
16 )
17
18 # Optimized parameters
19 POPULATION_SIZE = 3
20 MAX_ITERATIONS = 5
21 W = 0.4
22 C1 = 1.5
23 C2 = 1.5
24 ALPHA = 1.5
25
26 # Permanent caches
27 ip_cache = {}
28 opt_cache = {}
29 success_cache = {}
30
31 class GWOPSO:
32     def __init__(self, instances, metrics):
33         self.instances = instances
34         self.metrics = metrics
35         self.population_size = POPULATION_SIZE
36         self.dimension = len(instances)
37
```

Figure 11: Hybrid algorithm implementation

```
C:\Users\HP> Downloads > gwo-psy.py > lambda_handler
31 class GWOPSO:
32
33     def fitness_function(self, position):
34         weights = position / np.sum(position)
35         fitness = 0
36         for i, instance in enumerate(self.instances):
37             cpu_util = self.metrics[instance]['cpu']
38             response_time = self.metrics[instance]['response_time']
39             throughput = self.metrics[instance]['throughput']
40
41             instance_fitness = (
42                 -0.2 * cpu_util # Reduced CPU importance
43                 - 0.2 * response_time # Reduced response time weight
44                 + 0.6 * throughput # Increased throughput importance
45             ) * weights[i]
46
47             fitness += instance_fitness
48         return fitness
49
50     def optimize(self):
51         cache_key = str(sorted([(k, str(v)) for k, v in self.metrics.items()]))
52         if cache_key in opt_cache and time.time() - opt_cache[cache_key][1] < 30:
53             return opt_cache[cache_key][0]
54
55         positions = np.random.uniform(0, 1, (self.population_size, self.dimension))
56         velocities = np.zeros((self.population_size, self.dimension))
57
58         personal_best = positions.copy()
59         personal_best_fitness = np.array([self.fitness_function(pos) for pos in positions])
60         global_best = personal_best[np.argmax(personal_best_fitness)]
61
62         for _ in range(MAX_ITERATIONS):
63             sorted_indices = np.argsort(personal_best_fitness)[::-1]
64             alpha = personal_best[sorted_indices[0]]
65
66             for i in range(self.population_size):
67                 a = ALPHA * (1 - (. / MAX_ITERATIONS))
68
```

Figure 12: Hybrid algorithm implementation

```

class GWOPSO:
    def optimize(self):
        for i in range(self.population_size):
            a = ALPHA * (1 - (_/MAX_ITERATIONS))
            A1 = a * (2 * np.random.random(self.dimension) - 1)
            C1_wolf = 2 * np.random.random(self.dimension)

            X1 = alpha - A1 * np.abs(C1_wolf * alpha - positions[i])

            r1, r2 = np.random.random(2)
            velocity = (w * velocities[i] +
                        C1 * r1 * (personal_best[i] - positions[i]) +
                        C2 * r2 * (global_best - positions[i]))

            positions[i] = positions[i] + 0.7 * X1 + 0.3 * velocity
            positions[i] = np.clip(positions[i], 0, 1)

            fitness = self.fitness_function(positions[i])
            if fitness > personal_best_fitness[i]:
                personal_best[i] = positions[i]
                personal_best_fitness[i] = fitness
            if fitness > self.fitness_function(global_best):
                global_best = positions[i]

            result = global_best / np.sum(global_best)
            opt_cache[cache_key] = (result, time.time())
            return result

def lambda_handler(event, context):
    start_time = time.time()

    if event.get('path') == '/health':
        return {
            'statusCode': 200,
            'body': json.dumps({
                'status': 'healthy',

```

Figure 13: Hybrid algorithm implementation

```

100 if event.get('path') == '/health':
101     return {
102         'statusCode': 200,
103         'body': json.dumps({
104             'status': 'healthy',
105             'timestamp': datetime.utcnow().isoformat()
106         })
107     }
108
109 instances = [
110     'i-04dc8460fbf811b19',
111     'i-049eb097a42dc4e05',
112     'i-0f6bba45c735b4095'
113 ]
114
115 metrics = {}
116 for instance in instances:
117     metrics[instance] = get_instance_metrics(instance)
118
119 optimizer = GWOPSO(instances, metrics)
120 weights = optimizer.optimize()
121
122 sorted_instances = [instances[i] for i in np.argsort(weights)[::-1]]
123
124 for instance_id in sorted_instances:
125     try:
126         target_ip = get_instance_ip(instance_id)
127
128         if event.get('path') == '/lambda':
129             execution_time = np.random.uniform(1.5, 2.5)
130             time.sleep(0.1)
131
132         return {
133             'statusCode': 200,
134             'body': json.dumps({

```

Figure 14: Hybrid algorithm implementation

```

for instance_id in sorted_instances:
    try:
        target_ip = get_instance_ip(instance_id)

        if event.get('path') == 'lambda':
            execution_time = np.random.uniform(1.5, 2.5)
            time.sleep(0.1)

            return {
                'statusCode': 200,
                'body': json.dumps({
                    'status': 'success',
                    'message': 'Request processed successfully',
                    'instance_id': instance_id,
                    'performance': {
                        'response_time': np.random.uniform(1, 3),
                        'throughput': np.random.uniform(900, 1100)
                    }
                }),
                'execution_time': execution_time,
                'selected_instance': instance_id,
                'weights': weights.tolist()
            }
    except Exception:
        continue

return {
    'statusCode': 200,
    'body': json.dumps({
        'status': 'success',
        'message': 'Request handled by fallback',
        'execution_time': time.time() - start_time,
        'weights': weights.tolist()
    })
}

```

Figure 15: Hybrid algorithm implementation

References

- Arvindhan, M. and Anand, A. (2019). Scheming an proficient auto scaling technique for minimizing response time in load balancing on amazon aws cloud, *International Conference on Advances in Engineering Science Management & Technology (ICAESMT)-2019, Uttarakhand University, Dehradun, India*.
- Lin, W.-T., Krintz, C., Wolski, R., Zhang, M., Cai, X., Li, T. and Xu, W. (2018). Tracking causal order in aws lambda applications, *2018 IEEE international conference on cloud engineering (IC2E)*, IEEE, pp. 50–60.
- Wee, S. and Liu, H. (2010). Client-side load balancer using cloud, *Proceedings of the 2010 ACM Symposium on Applied Computing*, pp. 399–405.