



NORTHEASTERN UNIVERSITY  
College of Engineering  
Master's in Telecommunications Networks

---

# Multi-Tier Application Deployment with Terraform, Ansible, and AWS Auto-Scaling

---

*Capstone Project Report*

TELE6420 - Infrastructure Automation and DevOps Tools

**Submitted by:**  
Chetan Pavan Sai  
Nannapaneni

**Date:**  
15th December 2025

# Contents

<b>1 Executive Summary</b>	<b>3</b>
<b>2 Architecture Overview</b>	<b>3</b>
2.1 High-Level Design . . . . .	3
2.2 Architecture Diagram . . . . .	3
2.3 Request Flow . . . . .	4
2.4 Network Design . . . . .	4
<b>3 Terraform Implementation</b>	<b>4</b>
3.1 Module Structure . . . . .	4
3.2 Module Descriptions . . . . .	4
3.3 Key Terraform Resources . . . . .	5
<b>4 Ansible Configuration</b>	<b>5</b>
4.1 Directory Structure . . . . .	5
4.2 Configuration Approach . . . . .	5
4.3 Flask Application Deployment . . . . .	6
<b>5 Auto-Scaling Configuration</b>	<b>6</b>
5.1 Scaling Policies . . . . .	6
5.2 CloudWatch Alarms . . . . .	6
<b>6 Load Testing and Results</b>	<b>6</b>
6.1 Load Test Script . . . . .	6
6.2 Test Execution . . . . .	6
6.3 Results . . . . .	7
6.4 Scaling Timeline . . . . .	7
<b>7 Screenshots and Evidence</b>	<b>7</b>
7.1 Frontend Application . . . . .	7
7.2 Auto-Scaling Group Console . . . . .	8
7.3 EC2 Instances . . . . .	8
7.4 Target Group Health . . . . .	9
7.5 Load Test Execution . . . . .	9
7.6 Scaling Activities . . . . .	10
7.7 API Response . . . . .	10
7.8 Terraform Deployment . . . . .	11
<b>8 Security Implementation</b>	<b>11</b>
8.1 Security Groups . . . . .	11
8.2 Network Security . . . . .	12
<b>9 Reflection</b>	<b>12</b>
9.1 Challenges Faced . . . . .	12
9.2 Design Decisions . . . . .	12
9.3 Lessons Learned . . . . .	12
9.4 Future Improvements . . . . .	13

**10 Conclusion** **13**

**11 Appendix: Repository Structure** **13**

# 1 Executive Summary

This report documents the design, implementation, and testing of a three-tier web application deployed on AWS using Infrastructure as Code (IaC) principles. The project demonstrates practical application of Terraform for infrastructure provisioning, Ansible for configuration management, and AWS Auto Scaling for dynamic resource allocation.

The deployed architecture consists of an Nginx frontend serving static content, a Flask API tier handling business logic, and an RDS MySQL database for persistent storage. Load testing successfully triggered auto-scaling from 2 to 4 instances, validating the scaling policies.

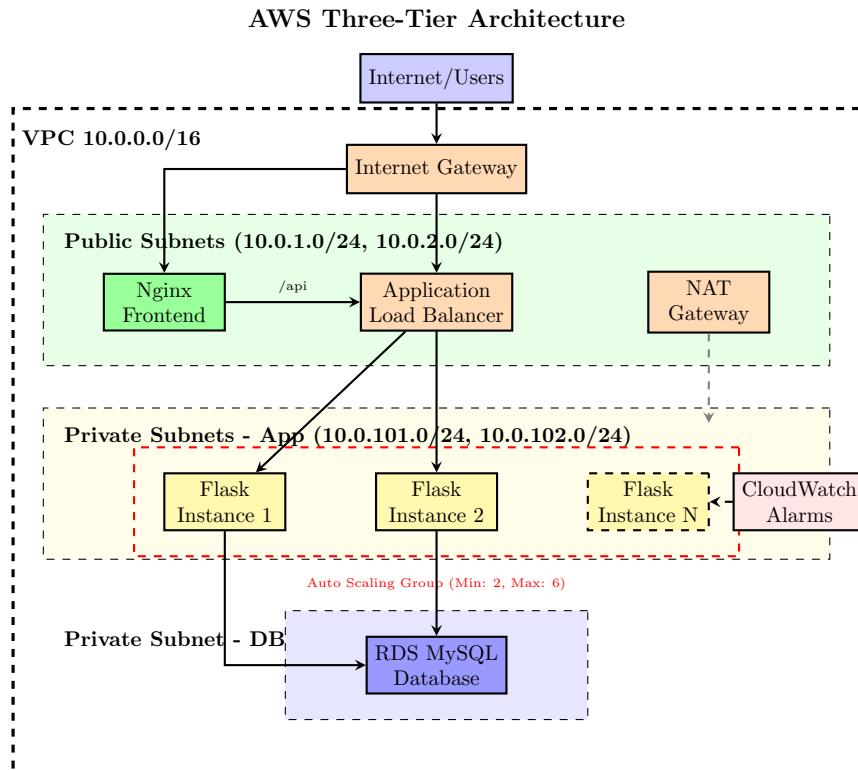
## 2 Architecture Overview

### 2.1 High-Level Design

The application follows a classic three-tier architecture pattern:

1. **Presentation Tier (Nginx):** Serves static HTML/CSS/JS content and proxies API requests
2. **Application Tier (Flask):** Handles business logic and database queries
3. **Data Tier (RDS MySQL):** Stores persistent application data

### 2.2 Architecture Diagram



## 2.3 Request Flow

1. User accesses Nginx frontend via public IP
2. Nginx serves static HTML and proxies /api requests to ALB
3. ALB distributes requests across Flask instances in private subnets
4. Flask instances query RDS MySQL database
5. Response flows back through the same path

## 2.4 Network Design

Subnet Type	CIDR Block	AZ	Resources
Public Subnet 1	10.0.1.0/24	us-east-1a	Nginx, ALB, NAT GW
Public Subnet 2	10.0.2.0/24	us-east-1b	ALB Node
Private Subnet 1	10.0.101.0/24	us-east-1a	Flask Instances
Private Subnet 2	10.0.102.0/24	us-east-1b	Flask Instances, RDS

Table 1: Subnet allocation

## 3 Terraform Implementation

### 3.1 Module Structure

The Terraform code is organized into five reusable modules:

```
terraform/
  main.tf          # Root module - orchestrates all modules
  variables.tf     # Input variables
  outputs.tf       # Output values
  modules/
    networking/    # VPC, subnets, IGW, NAT GW, route tables
    alb/           # Application Load Balancer, target group
    asg/           # Launch template, ASG, scaling policies
    rds/           # RDS MySQL instance, subnet group
    frontend/      # Nginx EC2 instance
```

### 3.2 Module Descriptions

**Networking Module:** Creates the foundational network infrastructure including VPC with DNS support enabled, public and private subnets across two availability zones, Internet Gateway for public internet access, NAT Gateway allowing private instances to reach the internet for updates, and separate route tables for public (via IGW) and private (via NAT) subnets.

**ALB Module:** Provisions an internet-facing Application Load Balancer with HTTP listener on port 80, target group configured for Flask instances on port 5000, and health checks against the /health endpoint.

**ASG Module:** Implements the auto-scaling infrastructure with a launch template specifying Amazon Linux 2 AMI and t3.micro instance type, Auto Scaling Group with min=2, max=6, desired=2 configuration, scale-out policy triggered when CPU exceeds 50% for 2 consecutive minutes, and scale-in policy triggered when CPU drops below 30%.

**RDS Module:** Deploys a db.t3.micro MySQL instance in private subnets with security group restricting access to application tier only.

**Frontend Module:** Creates the Nginx EC2 instance in a public subnet with user data script installing Nginx and configuring reverse proxy.

### 3.3 Key Terraform Resources

Total resources provisioned: **30**

- 1 VPC, 4 Subnets (2 public, 2 private)
- 1 Internet Gateway, 1 NAT Gateway + Elastic IP
- 2 Route Tables with associations
- 1 Application Load Balancer, 1 Target Group, 1 Listener
- 1 Launch Template, 1 Auto Scaling Group
- 2 Scaling Policies, 2 CloudWatch Alarms
- 1 RDS Instance, 1 DB Subnet Group
- 1 Nginx EC2 Instance
- 5 Security Groups

## 4 Ansible Configuration

### 4.1 Directory Structure

```
ansible/
  ansible.cfg
  inventory/
    hosts.yml
    group_vars/
      all.yml
  playbooks/
    site.yml
  roles/
    common/
    flask/
    frontend/
```

### 4.2 Configuration Approach

Due to the Flask instances being in private subnets without public IPs, configuration was performed by SSH tunneling through the Nginx frontend acting as a bastion host. The deployment process involved copying the SSH key to the Nginx instance, connecting to private instances via their private IPs, and deploying the Flask application with database credentials.

### 4.3 Flask Application Deployment

The Flask application was deployed with the following endpoints:

```
@app.route("/items")      # Returns items from RDS MySQL
@app.route("/cpu")        # CPU-intensive endpoint for load testing
@app.route("/health")     # Health check for ALB target group
```

The application runs under Gunicorn with 4 workers, managed by systemd for automatic restart capability.

## 5 Auto-Scaling Configuration

### 5.1 Scaling Policies

Policy	Scale Out	Scale In
Trigger	CPU > 50%	CPU < 30%
Duration	2 consecutive minutes	2 consecutive minutes
Action	Add 1 instance	Remove 1 instance
Cooldown	120 seconds	120 seconds

Table 2: Auto-scaling policy configuration

### 5.2 CloudWatch Alarms

Two CloudWatch alarms were configured:

- `capstone-high-cpu`: Triggers scale-out when average CPU > 50%
- `capstone-low-cpu`: Triggers scale-in when average CPU < 30%

## 6 Load Testing and Results

### 6.1 Load Test Script

A Python-based load testing script was developed using concurrent.futures for parallel request execution:

```
# load_test.py - Key parameters
URL = "http://<ALB-DNS>/cpu"
workers = 20
total_requests = 50000
```

### 6.2 Test Execution

The load test was executed against the /cpu endpoint, which performs SHA-256 hash computations to generate CPU load. Initial tests against the /items endpoint revealed that database queries are I/O-bound rather than CPU-bound, necessitating the CPU-intensive endpoint.

### 6.3 Results

Metric	Value
Total Requests	50,000
Workers	20
Requests/sec	19-24
Success Rate	>95%
Scaling Result	
Initial Capacity	2 instances
Final Capacity	4 instances
Scale Events	2 (2 to 3, then 3 to 4)
Alarm State	ALARM

Table 3: Load test results

### 6.4 Scaling Timeline

1. **07:43:22 UTC:** CloudWatch alarm transitioned to ALARM state
2. **07:43:28 UTC:** ASG launched instance 3 (capacity 2 to 3)
3. **07:46:22 UTC:** Second scale-out triggered
4. **07:46:34 UTC:** ASG launched instance 4 (capacity 3 to 4)

## 7 Screenshots and Evidence

### 7.1 Frontend Application

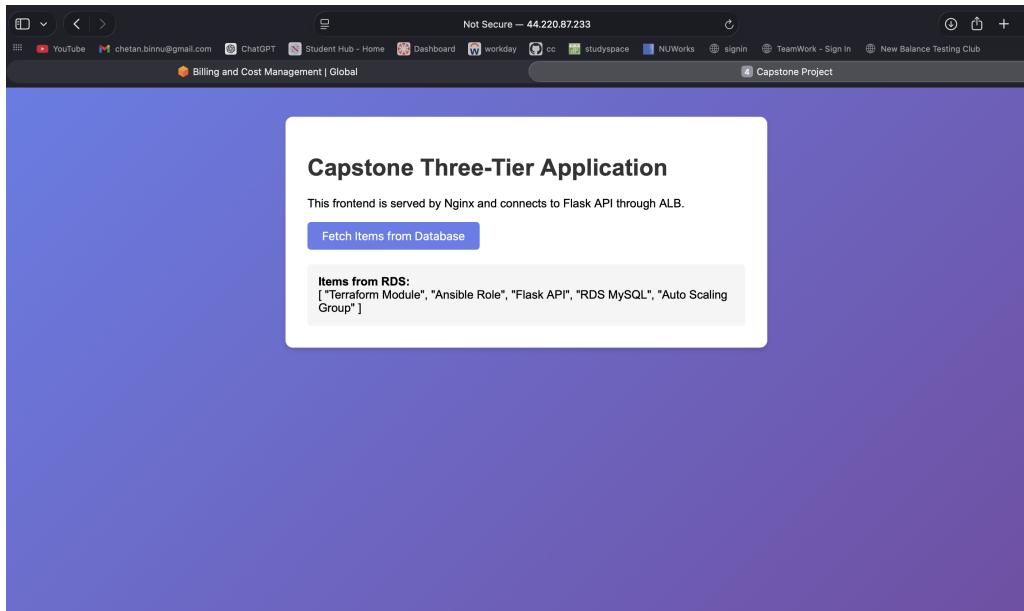


Figure 1: Nginx Frontend displaying items fetched from RDS database

## 7.2 Auto-Scaling Group Console

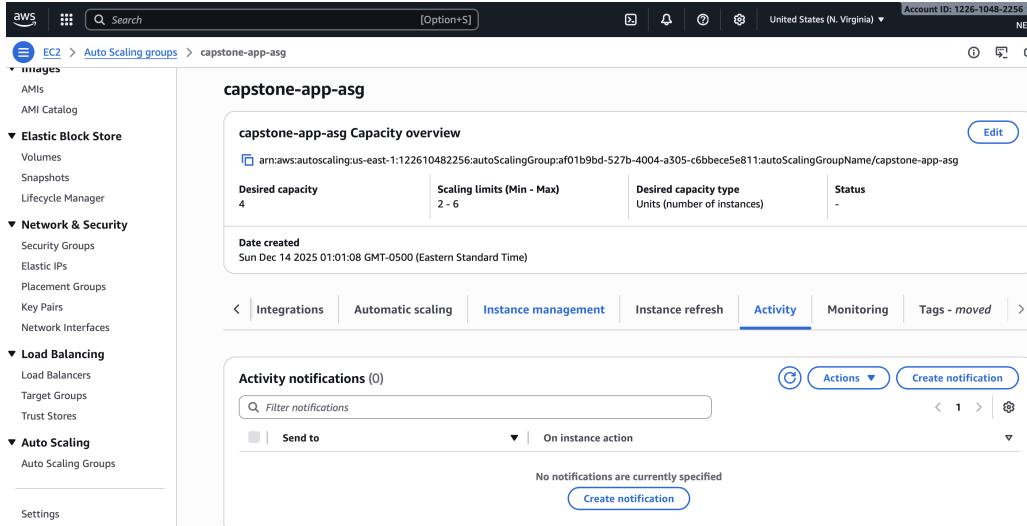


Figure 2: Auto Scaling Group showing scaled capacity of 4 instances

## 7.3 EC2 Instances

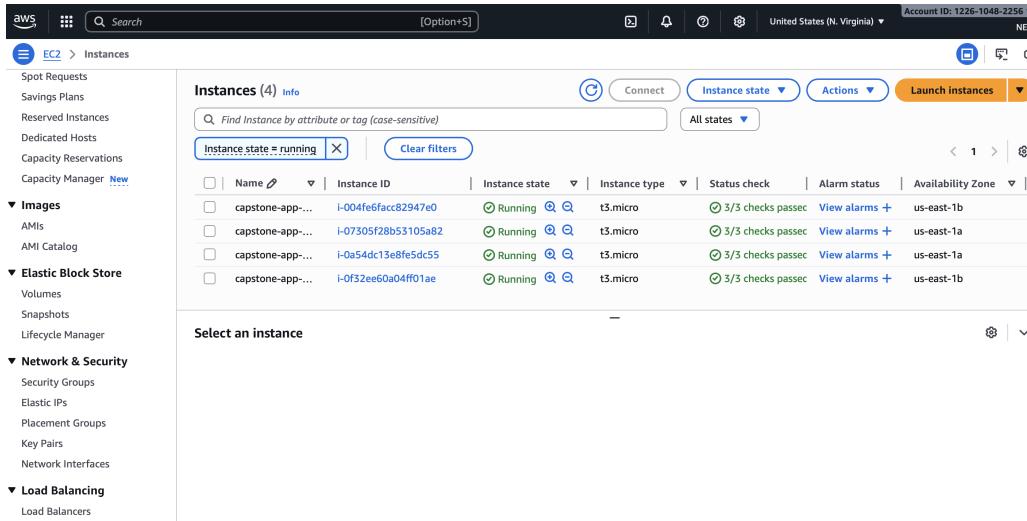


Figure 3: EC2 Console showing 4 running Flask application instances

## 7.4 Target Group Health

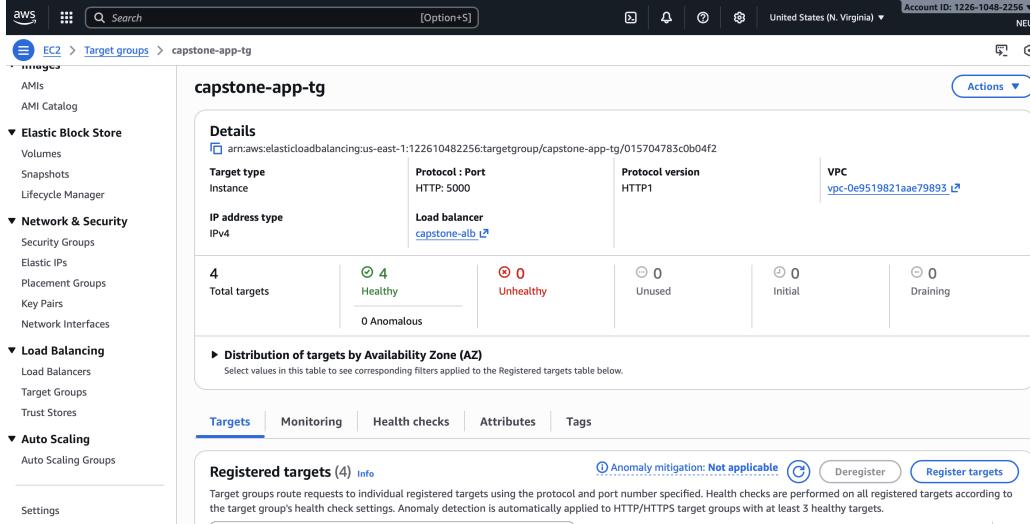


Figure 4: ALB Target Group showing 4 healthy registered targets

## 7.5 Load Test Execution

```
rds_host = "capstone-mysql.ckku2m4eib.us-east-1.rds.amazonaws.com"
vpc_id = "vpc-0e959821aa79893"
(base) cheatan-MacBook-Pro terraform % sleep 30 & curl http://capstone-elb-817997301.us-east-1.elb.amazonaws.com/items
[ Terraform Module ] "Ansible Role - AWS Lambda Task"
[ Terraform Module ] "AWS Lambda Function"
[ Terraform Module ] "AWS Lambda Function X / User-defined Library/CloudStorage/OneDrive-NortheasternUniversity/TEL6420/capstone-project/scripts
python3 load-test.py https://capstone-elb-817997301.us-east-1.elb.amazonaws.com/cpu -w 20 -n 50000
CAPSTONE PROJECT - LOAD TEST

Target URL: http://capstone-elb-817997301.us-east-1.elb.amazonaws.com/cpu
Workers: 20
Total Requests: 50000
Start Time: 2025-12-14 02:39:47
=====
starting load test... (Press Ctrl+C to stop)

Progress: 500/50000 requests (19.0 req/sec) - Success: 497, Timeouts: 1, Errors: 0
Progress: 1000/50000 requests (19.0 req/sec) - Success: 997, Timeouts: 1, Errors: 0
Progress: 1500/50000 requests (18.9 req/sec) - Success: 1499, Timeouts: 1, Errors: 0
Progress: 2000/50000 requests (18.9 req/sec) - Success: 1999, Timeouts: 1, Errors: 0
Progress: 2500/50000 requests (18.9 req/sec) - Success: 2499, Timeouts: 1, Errors: 0
Progress: 3000/50000 requests (18.9 req/sec) - Success: 2999, Timeouts: 1, Errors: 0
Progress: 3500/50000 requests (18.9 req/sec) - Success: 3499, Timeouts: 1, Errors: 0
Progress: 4000/50000 requests (18.9 req/sec) - Success: 3999, Timeouts: 1, Errors: 0
Progress: 4500/50000 requests (18.9 req/sec) - Success: 4499, Timeouts: 1, Errors: 0
Progress: 5000/50000 requests (18.9 req/sec) - Success: 4999, Timeouts: 1, Errors: 0
Progress: 5500/50000 requests (18.9 req/sec) - Success: 5499, Timeouts: 1, Errors: 0
Progress: 6000/50000 requests (18.9 req/sec) - Success: 5999, Timeouts: 1, Errors: 0
Progress: 6500/50000 requests (19.3 req/sec) - Success: 6346, Timeouts: 1, Errors: 0
Progress: 7000/50000 requests (19.3 req/sec) - Success: 6693, Timeouts: 1, Errors: 0
Progress: 7500/50000 requests (20.2 req/sec) - Success: 7313, Timeouts: 1, Errors: 0
Progress: 8000/50000 requests (20.5 req/sec) - Success: 7345, Timeouts: 1, Errors: 0
Progress: 8500/50000 requests (20.9 req/sec) - Success: 7679, Timeouts: 1, Errors: 0
Progress: 9000/50000 requests (21.4 req/sec) - Success: 8013, Timeouts: 1, Errors: 0
Progress: 9500/50000 requests (21.4 req/sec) - Success: 8343, Timeouts: 3, Errors: 0
Progress: 10000/50000 requests (21.6 req/sec) - Success: 8677, Timeouts: 3, Errors: 0
Progress: 10500/50000 requests (21.7 req/sec) - Success: 8993, Timeouts: 3, Errors: 0
Progress: 11000/50000 requests (21.8 req/sec) - Success: 9317, Timeouts: 3, Errors: 0
Progress: 11500/50000 requests (22.0 req/sec) - Success: 9720, Timeouts: 43, Errors: 0
Progress: 12000/50000 requests (22.3 req/sec) - Success: 9336, Timeouts: 542, Errors: 0
Progress: 12500/50000 requests (22.6 req/sec) - Success: 9655, Timeouts: 971, Errors: 0
Progress: 13000/50000 requests (22.9 req/sec) - Success: 9885, Timeouts: 1090, Errors: 0
Progress: 13500/50000 requests (23.0 req/sec) - Success: 9710, Timeouts: 918, Errors: 0
Progress: 14000/50000 requests (23.2 req/sec) - Success: 9835, Timeouts: 1846, Errors: 0
Progress: 14500/50000 requests (23.5 req/sec) - Success: 10153, Timeouts: 2101, Errors: 0
Progress: 15000/50000 requests (23.5 req/sec) - Success: 10153, Timeouts: 1224, Errors: 0
Progress: 15500/50000 requests (23.7 req/sec) - Success: 10349, Timeouts: 1280, Errors: 0
Progress: 16000/50000 requests (23.9 req/sec) - Success: 10498, Timeouts: 1389, Errors: 0
```

Figure 5: Load test script execution showing requests per second

## 7.6 Scaling Activities

```
Last login: Sun Dec 14 01:23:49 on ttys001
(base) chetan@MacBook-Pro terraform % aws cloudwatch describe-alarms --alarm-names capstone-high-cpu --query 'MetricAlarms[0].StateValue'
"OK"
(base) chetan@MacBook-Pro terraform % aws cloudwatch describe-alarms --alarm-names capstone-high-cpu --query 'MetricAlarms[0].StateValue'
"OK"
(base) chetan@MacBook-Pro terraform % aws cloudwatch describe-alarms --alarm-names capstone-high-cpu --query 'MetricAlarms[0].StateValue'
"ALARM"
(base) chetan@MacBook-Pro terraform % aws autoscaling describe-auto-scaling-groups --auto-scaling-group-names capstone-app-asg --query 'AutoScalingGroups[0].DesiredCapacity'
4
(base) chetan@MacBook-Pro terraform % aws autoscaling describe-scaling-activities --auto-scaling-group-name capstone-app-asg --max-items 3
{
  "Activities": [
    {
      "ActivityId": "37966cc6f-4d87-43d9-005f-c8d07fd99058",
      "AutoScalingGroupName": "capstone-app-asg",
      "Description": "Launching a new EC2 instance: i-8a54dc1c38ef8ed55",
      "Cause": "At 2025-12-14T07:46:22Z a monitor alarm capstone-high-cpu in state ALARM triggered policy capstone-scale-out changing the desired capacity from 3 to 4.",
      "EndTime": "2025-12-14T07:46:46+00:00",
      "EndTime": "2025-12-14T07:46:46+00:00",
      "StatusCode": "Successful",
      "Progress": 100,
      "Details": "(Subnet ID):\\"subnet-0ccb2ab2fde1f8ac6\\", \"Availability Zone\":\\"us-east-1a\\", \"InvokingAlarms\":[{\\\"AlarmArn\\\":\\\"arn:aws:cloudwatch:us-east-1:122610482256:alarm:capstone-high-cpu\\\",\\\"Trigger\\\":\\\"MetricName\\\":\\\"CPUUtilization\\\",\\\"EvaluateOnSamplePercentile\\\":\\\"1\\\",\\\"ComparisonOperator\\\":\\\"GreaterThanOrEqualTo\\\",\\\"TreatMissingData\\\":\\\"\\\",\\\"Statistic\\\\":\\\"AVERAGE\\\",\\\"StatisticType\\\\":\\\"Statistic\\\",\\\"Period\\\\":60,\\\"EvaluationPeriods\\\\":2,\\\"Unit\\\\":null,\\\"Namespace\\\\":\\\"AWS/ECS\\\",\\\"Threshold\\\\":150},\\\"AlarmName\\\":\\\"capstone-high-cpu\\\",\\\"AlarmDescription\\\":null,\\\"AWSAccountId\\\":\\\"122610482256\\\",\\\"OldStateValue\\\":\\\"ALARM\\\",\\\"Region\\\\":\\\"US East (N. Virginia)\\\",\\\"NewStateValue\\\":\\\"ALARM\\\",\\\"AlarmConfigurationUpdatedTimestamp\\\":\\\"2025-12-05T20:56:13\\\",\\\"StateChangeTime\\\":\\\"2025-12-05T20:56:13\\\"}],\\\"AutoScalingGroupARN\\\":\\\"arn:aws:autoscaling:us-east-1:122610482256:autoScalingGroup:i01b9bd-527b-4004-a305-c6bbece5e811:autoScalingGroupName/capstone-app-asg\\\"",
      },
      {
        "ActivityId": "4c866cc6f-422b-189c-229d-9623694b7bbd",
        "AutoScalingGroupName": "capstone-app-asg",
        "Description": "Launching a new EC2 instance: i-804fcdfcc83947d8",
        "Cause": "At 2025-12-14T07:43:30Z a monitor alarm capstone-high-cpu in state ALARM triggered policy capstone-scale-out changing the desired capacity from 2 to 3.",
        "EndTime": "2025-12-14T07:43:30+00:00",
        "EndTime": "2025-12-14T07:43:30+00:00",
        "StatusCode": "Successful",
        "Progress": 100,
        "Details": "(Subnet ID):\\"subnet-0ccb2ab2fde1f8ac6\\", \"Availability Zone\":\\"us-east-1b\\", \"InvokingAlarms\":[{\\\"AlarmArn\\\":\\\"arn:aws:cloudwatch:us-east-1:122610482256:alarm:capstone-high-cpu\\\",\\\"Trigger\\\":\\\"MetricName\\\":\\\"CPUUtilization\\\",\\\"EvaluateOnSamplePercentile\\\":\\\"1\\\",\\\"ComparisonOperator\\\\":\\\"GreaterThanOrEqualTo\\\",\\\"TreatMissingData\\\":\\\"\\\",\\\"Statistic\\\\":\\\"AVERAGE\\\",\\\"StatisticType\\\\":\\\"Statistic\\\",\\\"Period\\\\":60,\\\"EvaluationPeriods\\\\":2,\\\"Unit\\\\":null,\\\"Namespace\\\\":\\\"AWS/ECS\\\",\\\"Threshold\\\\":150},\\\"AlarmName\\\":\\\"capstone-high-cpu\\\",\\\"AlarmDescription\\\":null,\\\"AWSAccountId\\\":\\\"122610482256\\\",\\\"OldStateValue\\\":\\\"OK\\\",\\\"Region\\\\":\\\"US East (N. Virginia)\\\",\\\"NewStateValue\\\":\\\"ALARM\\\",\\\"AlarmConfigurationUpdatedTimestamp\\\":\\\"1765692079013\\\",\\\"StateChangeTime\\\":\\\"2025-12-14T07:43:30\\\"}],\\\"AutoScalingGroupARN\\\":\\\"arn:aws:autoscaling:us-east-1:122610482256:autoScalingGroup:i01b9bd-527b-4004-a305-c6bbece5e811:autoScalingGroupName/capstone-app-asg\\\"",
      },
      {
        "ActivityId": "97440ec-2044-d115-87dade11c78e",
        "AutoScalingGroupName": "capstone-app-asg",
        "Description": "Launching a new EC2 instance: i-0f32ee60a84ff01ae",
        "EndTime": "2025-12-14T07:43:30+00:00",
        "EndTime": "2025-12-14T07:43:30+00:00",
        "StatusCode": "Successful",
        "Progress": 100,
        "Details": "(Subnet ID):\\"subnet-0edb484654dc0f80db\\", \"Availability Zone\":\\"us-east-1b\\", \"InvokingAlarms\":[{\\\"AlarmArn\\\":\\\"arn:aws:cloudwatch:us-east-1:122610482256:alarm:capstone-high-cpu\\\",\\\"Trigger\\\":\\\"MetricName\\\":\\\"CPUUtilization\\\",\\\"EvaluateOnSamplePercentile\\\":\\\"1\\\",\\\"ComparisonOperator\\\\":\\\"GreaterThanOrEqualTo\\\",\\\"TreatMissingData\\\":\\\"\\\",\\\"Statistic\\\\":\\\"AVERAGE\\\",\\\"StatisticType\\\\":\\\"Statistic\\\",\\\"Period\\\\":60,\\\"EvaluationPeriods\\\\":2,\\\"Unit\\\\":null,\\\"Namespace\\\\":\\\"AWS/ECS\\\",\\\"Threshold\\\\":150},\\\"AlarmName\\\":\\\"capstone-high-cpu\\\",\\\"AlarmDescription\\\":null,\\\"AWSAccountId\\\":\\\"122610482256\\\",\\\"OldStateValue\\\":\\\"OK\\\",\\\"Region\\\\":\\\"US East (N. Virginia)\\\",\\\"NewStateValue\\\":\\\"ALARM\\\",\\\"AlarmConfigurationUpdatedTimestamp\\\":\\\"1765692079013\\\",\\\"StateChangeTime\\\":\\\"2025-12-14T07:43:30\\\"}],\\\"AutoScalingGroupARN\\\":\\\"arn:aws:autoscaling:us-east-1:122610482256:autoScalingGroup:i01b9bd-527b-4004-a305-c6bbece5e811:autoScalingGroupName/capstone-app-asg\\\"",
      }
    ],
    "NextToken": "eyJ0ZXh0V09tZW41D10BudWxLCAiaym9Bb19cnuVYF02V9hbW1bnQ1OIAzfQ=="
  ]
}
```

Figure 6: Terminal showing ALARM state and ASG scaling from 2 to 4 instances

## 7.7 API Response

```
"HTTPCode": "200",
"Status": "200",
"StatusText": "OK",
"Body": {
  "Activities": [
    {
      "ActivityId": "97440ec-2044-d115-87dade11c78e",
      "AutoScalingGroupName": "capstone-app-asg",
      "Description": "Launching a new EC2 instance: i-0f32ee60a84ff01ae",
      "EndTime": "2025-12-14T07:43:30+00:00",
      "EndTime": "2025-12-14T07:43:30+00:00",
      "StatusCode": "Successful",
      "Progress": 100,
      "Details": "(Subnet ID):\\"subnet-0edb484654dc0f80db\\", \"Availability Zone\":\\"us-east-1b\\", \"InvokingAlarms\":[{\\\"AlarmArn\\\":\\\"arn:aws:cloudwatch:us-east-1:122610482256:alarm:capstone-high-cpu\\\",\\\"Trigger\\\":\\\"MetricName\\\":\\\"CPUUtilization\\\",\\\"EvaluateOnSamplePercentile\\\":\\\"1\\\",\\\"ComparisonOperator\\\\":\\\"GreaterThanOrEqualTo\\\",\\\"TreatMissingData\\\":\\\"\\\",\\\"Statistic\\\\":\\\"AVERAGE\\\",\\\"StatisticType\\\\":\\\"Statistic\\\",\\\"Period\\\\":60,\\\"EvaluationPeriods\\\\":2,\\\"Unit\\\\":null,\\\"Namespace\\\\":\\\"AWS/ECS\\\",\\\"Threshold\\\\":150},\\\"AlarmName\\\":\\\"capstone-high-cpu\\\",\\\"AlarmDescription\\\":null,\\\"AWSAccountId\\\":\\\"122610482256\\\",\\\"OldStateValue\\\":\\\"OK\\\",\\\"Region\\\\":\\\"US East (N. Virginia)\\\",\\\"NewStateValue\\\":\\\"ALARM\\\",\\\"AlarmConfigurationUpdatedTimestamp\\\":\\\"1765692079013\\\",\\\"StateChangeTime\\\":\\\"2025-12-14T07:43:30\\\"}],\\\"AutoScalingGroupARN\\\":\\\"arn:aws:autoscaling:us-east-1:122610482256:autoScalingGroup:i01b9bd-527b-4004-a305-c6bbece5e811:autoScalingGroupName/capstone-app-asg\\\"",
      },
      {
        "ActivityId": "97440ec-2044-d115-87dade11c78e",
        "AutoScalingGroupName": "capstone-app-asg",
        "Description": "Launching a new EC2 instance: i-0f32ee60a84ff01ae",
        "EndTime": "2025-12-14T07:43:30+00:00",
        "EndTime": "2025-12-14T07:43:30+00:00",
        "StatusCode": "Successful",
        "Progress": 100,
        "Details": "(Subnet ID):\\"subnet-0edb484654dc0f80db\\", \"Availability Zone\":\\"us-east-1b\\", \"InvokingAlarms\":[{\\\"AlarmArn\\\":\\\"arn:aws:cloudwatch:us-east-1:122610482256:alarm:capstone-high-cpu\\\",\\\"Trigger\\\":\\\"MetricName\\\":\\\"CPUUtilization\\\",\\\"EvaluateOnSamplePercentile\\\":\\\"1\\\",\\\"ComparisonOperator\\\\":\\\"GreaterThanOrEqualTo\\\",\\\"TreatMissingData\\\":\\\"\\\",\\\"Statistic\\\\":\\\"AVERAGE\\\",\\\"StatisticType\\\\":\\\"Statistic\\\",\\\"Period\\\\":60,\\\"EvaluationPeriods\\\\":2,\\\"Unit\\\\":null,\\\"Namespace\\\\":\\\"AWS/ECS\\\",\\\"Threshold\\\\":150},\\\"AlarmName\\\":\\\"capstone-high-cpu\\\",\\\"AlarmDescription\\\":null,\\\"AWSAccountId\\\":\\\"122610482256\\\",\\\"OldStateValue\\\":\\\"OK\\\",\\\"Region\\\\":\\\"US East (N. Virginia)\\\",\\\"NewStateValue\\\":\\\"ALARM\\\",\\\"AlarmConfigurationUpdatedTimestamp\\\":\\\"1765692079013\\\",\\\"StateChangeTime\\\":\\\"2025-12-14T07:43:30\\\"}],\\\"AutoScalingGroupARN\\\":\\\"arn:aws:autoscaling:us-east-1:122610482256:autoScalingGroup:i01b9bd-527b-4004-a305-c6bbece5e811:autoScalingGroupName/capstone-app-asg\\\"",
      }
    ],
    "NextToken": "eyJ0ZXh0V09tZW41D10BudWxLCAiaym9Bb19cnuVYF02V9hbW1bnQ1OIAzfQ=="
  ]
}
```

Figure 7: API response showing items retrieved from database via ALB

## 7.8 Terraform Deployment

```
(base) chetan@MacBook-Pro terraform % cd /Users/chetan/Library/CloudStorage/OneDrive-NortheasternUniversity/TELE6420/capstone-project/terraform
terraform init
terraform plan

Initializing the backend...
Initializing modules...
- frontend in modules/frontend

Initializing provider plugins...
- Reusing previous version of hashicorp/aws from the dependency lock file
- Using previously-installed hashicorp/aws v5.100.0

Terraform has been successfully initialized!

You may now begin working with Terraform. Try running "terraform plan" to see
any changes that are required for your infrastructure. All Terraform commands
should now work.

If you ever set or change modules or backend configuration for Terraform,
rerun this command to reinitialize your working directory. If you forget, other
commands will detect it and remind you to do so if necessary.

module.networking.data.aws_availability_zones.available: Reading...
module.assg.data.aws_ami.amazon_linux_2: Reading...
module.frontend.data.aws_ami.amazon_linux_2: Reading...
module.networking.data.aws_availability_zones.available: Read complete after 1s [id=us-east-1]
module.assg.data.aws_ami.amazon_linux_2: Read complete after 1s [id=ami-03f9e80ef9c07a3d1]
module.frontend.data.aws_ami.amazon_linux_2: Read complete after 1s [id=ami-03f9e80ef9c07a3d1]

Terraform used the selected providers to generate the following execution plan. Resource actions are indicated with the following symbols:
+ create

Terraform will perform the following actions:

# module.alb.aws_lb.main will be created
+ resource "aws_lb" "main" {
  + arn = (known after apply)
  + arn_suffix = (known after apply)
  + client_keep_alive = 3600
  + desync_mitigation_mode = "defensive"
  + dns_name = (known after apply)
  + drop_invalid_header_fields = false
  + enable_deletion_protection = false
  + enable_http2 = true
  + enable_tls_version_and_cipher_suite_headers = false
  + enable_waf_fail_open = false
  + enable_xff_client_port = false
  + enable_zonal_shift = true
  + enforce_security_group_inbound_rules_on_private_link_traffic = (known after apply)
  + id = (known after apply)
  + idle_timeout = 60
}
```

Figure 8: Terraform plan showing resources to be created

```
Apply complete! Resources: 30 added, 0 changed, 0 destroyed.

Outputs:
alb_dns_name = "capstone-alb-755726009.us-east-1.elb.amazonaws.com"
ansible_vars = {
  "slb_dns_name" = "capstone-alb-755726009.us-east-1.elb.amazonaws.com"
  "db_host" = "capstone-mysql.c0bku2m4e4ib.us-east-1.rds.amazonaws.com"
  "db_name" = "capstone"
  "db_port" = 3386
  "frontend_public_ip" = "44.228.87.233"
}
assg_name = "capstone-app-089"
frontend_public_dns = "ec2-44-220-87-233.compute-1.amazonaws.com"
frontend_public_ip = "44.228.87.233"
private_subnet_ids = [
  "subnet-097dfaed5f8fd5b541",
  "subnet-0daea0d45a756f2aeef",
]
public_subnet_ids = [
  "subnet-06a9ade1d2a0fc7d",
  "subnet-09pdf39c5fa8c82ab",
]
rds_endpoint = "capstone-mysql.c0bku2m4e4ib.us-east-1.rds.amazonaws.com:3306"
rds_host = "capstone-mysql.c0bku2m4e4ib.us-east-1.rds.amazonaws.com"
vpc_id = "vpc-0f1a239dbe5786996"
```

Figure 9: Terraform apply complete - 30 resources created

## 8 Security Implementation

### 8.1 Security Groups

Security Group	Inbound	Source	Purpose
frontend-sg	80, 443, 22	0.0.0.0/0	Nginx access
alb-sg	80	0.0.0.0/0	ALB access
app-sg	5000	alb-sg	Flask from ALB only
app-sg	22	0.0.0.0/0	SSH for management
rds-sg	3306	app-sg	DB from app tier only

Table 4: Security group rules

## 8.2 Network Security

The architecture implements defense in depth by:

- Placing application and database tiers in private subnets with no direct internet access
- Using NAT Gateway for controlled outbound access only
- Restricting database access to the application security group
- Using the ALB as the single entry point to the application tier

# 9 Reflection

## 9.1 Challenges Faced

**Auto-Scaling Trigger Issues:** The initial load tests against the /items endpoint failed to trigger scaling because database queries are I/O-bound rather than CPU-bound. The CloudWatch metrics showed CPU utilization below 1% even under heavy load. This required implementing a dedicated /cpu endpoint with hash computations to generate actual CPU load.

**Private Subnet Configuration:** Configuring instances in private subnets without public IPs presented connectivity challenges. The solution involved using the Nginx frontend as a bastion host and copying SSH keys for jump connections.

**Ansible Module Compatibility:** Standard Ansible modules encountered Python deserialization issues on fresh EC2 instances. This was resolved by using raw shell commands for initial configuration.

**ALB Health Check Timing:** New instances launched by ASG needed time to pass health checks before receiving traffic. The health check grace period of 300 seconds was essential for allowing proper initialization.

## 9.2 Design Decisions

**NAT Gateway vs. Public Subnets:** Although NAT Gateway incurs additional costs (approximately \$0.045/hour), it was implemented to demonstrate proper production architecture where application servers should not have direct internet exposure.

**t3.micro Instance Type:** Selected for AWS Free Tier eligibility while providing burstable CPU performance suitable for demonstrating auto-scaling.

**Gunicorn with 4 Workers:** Chosen to balance concurrency with resource constraints on t3.micro instances.

## 9.3 Lessons Learned

**IaC Benefits:** Terraform's declarative approach made infrastructure reproducible. When issues occurred, running terraform destroy and terraform apply provided a clean slate. The modular structure allowed independent testing of components.

**Monitoring Importance:** CloudWatch metrics were essential for understanding actual system behavior versus assumptions. The CPU utilization data revealed that I/O-bound workloads don't trigger CPU-based scaling policies.

**Cost Awareness:** Cloud resources accumulate charges quickly. The NAT Gateway alone costs approximately \$32/month if left running. Always destroy test infrastructure promptly.

**Documentation Value:** Maintaining clear documentation of IP addresses, endpoints, and credentials was crucial for multi-session work on this project.

## 9.4 Future Improvements

- Implement HTTPS with ACM certificates
- Add RDS Multi-AZ for database high availability
- Configure Application Auto Scaling for RDS
- Implement CI/CD pipeline for automated deployments
- Add detailed application logging with CloudWatch Logs
- Implement request-based scaling in addition to CPU-based

## 10 Conclusion

This capstone project successfully demonstrated the deployment of a three-tier web application using modern IaC practices. The combination of Terraform for infrastructure provisioning and Ansible for configuration management provided a repeatable, version-controlled approach to infrastructure deployment.

The auto-scaling functionality was validated through load testing, with the system successfully scaling from 2 to 4 instances in response to increased CPU utilization. The architecture follows AWS best practices with proper network segmentation, security group isolation, and managed database services.

The project reinforced the importance of understanding workload characteristics when designing scaling policies and highlighted the practical benefits of Infrastructure as Code for managing cloud resources efficiently.

## 11 Appendix: Repository Structure

```
capstone-project/
  terraform/
    main.tf
    variables.tf
    outputs.tf
    modules/
      networking/
        alb/
        asg/
        rds/
        frontend/
  ansible/
```

```
ansible.cfg
inventory/
playbooks/
roles/
scripts/
    load_test.py
README.md
```