

# **UNIVERSITY OF BRISTOL**

EMATM0051: Large-Scale Data Engineering  
(2025 TB-1)

## **Coursework Assessment Report**

### **Part 1: ArtAI Architecture & Part 2: WordFreq Scaling**

**Student Name:**

Tanishk NanaSaheb Shinde

**University Student Number:**



**Programme:**

MSc Data Science

**Submission Date:**

December 2025

# Table of Contents

## Part 1: ArtAI Architecture

|  |    |
|--|----|
| 1. System Architectural Overview                 | 03 |
| 2. End-to-End Traffic Flow and Request Lifecycle | 03 |
| 3. Core Infrastructure & Design Decisions        | 04 |
| 4. Monitoring, Resilience, and Disaster Recovery | 04 |

## Part 2: WordFreq Scaling & Optimization

|   |    |
|---|----|
| Introduction: WordFreq App - Installation and Setup | 05 |
| Overview and Architecture                           | 05 |

### Task A: Installation and Configuration

|                                   |    |
|-----------------------------------|----|
| 1.1 Infrastructure Initialization | 06 |
| 1.2 Storage and Queuing Setup     | 06 |
| 1.3 Application Deployment        | 08 |

### Task B: Design and Implement Auto-scaling

|                                      |    |
|--------------------------------------|----|
| 2.1 AMI and Launch Template          | 09 |
| 2.2 Auto Scaling Group Configuration | 10 |
| 2.3 Dynamic Scaling Policies         | 11 |

### Task C: Load Testing

|  |    |
|--|----|
| 3.1 Test Execution                           | 12 |
| 3.2 Performance Analysis                     | 13 |
| 3.3 Optimization and Performance Experiments | 15 |

### Task D: Architecture Optimization

|  |    |
|--|----|
| 4.1 Reliability Pillar: Multi-AZ Deployment    | 17 |
| 4.2 Data Durability and Long-Term Storage      | 17 |
| 4.3 Cost Optimization                          | 18 |
| 4.4 Security: VPC Endpoints and Private Access | 18 |

### Task E: Further Improvements (Big Data)

|  |    |
|--|----|
| 5.1 Serverless Compute with AWS Lambda | 19 |
| 5.2 Big Data Analytics with Amazon EMR | 19 |

# ArtAI Webservice AWS Architecture Report

## 1. System Architectural Overview

The “Art AI” webservice has been developed to be a highly available, secure, and scalable platform for users to submit their images to be processed using AI. The “Art AI” webservice has a fully serverless architecture.

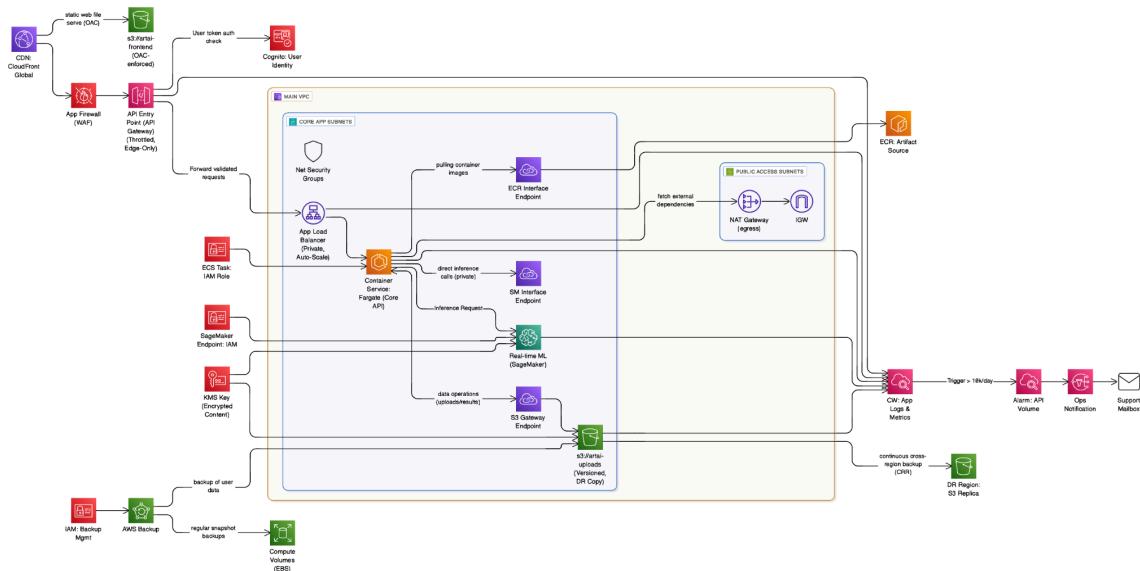
The “Art AI” webservice can be logically broken down into three tiers:

- **Global Edge** - Delivery and Perimeter Security through Amazon CloudFront and AWS WAF.
- **Public API Access** - API management and authentication through Amazon API Gateway and Amazon Cognito.
- **VPC Private Backend** - Business Logic (Amazon ECS Fargate) and AI Inference (Amazon SageMaker) in isolated subnets, inaccessible from the public internet.

## 2. End-to-End Traffic Flow and Request Lifecycle

There is a strict unidirectional flow enforced by the architecture to provide security:

- **User Access:** Users access the system through Amazon CloudFront, which caches the static frontend files (HTML, JS, CSS) from the frontend S3 Bucket to reduce latency.
- **Security & Auth:** Dynamic API requests will go through AWS WAF first, to filter out any malicious requests (SQLi/XSS, etc.) before routing valid requests through Amazon API Gateway to verify the user session against Amazon Cognito.
- **Private Processing:** After being authenticated, valid requests will route to an internal Application Load Balancer (ALB), which will route the requests to the Amazon ECS containers within the Private Subnet.
- **Inference & Storage:** The ECS Containers will send requests to Amazon SageMaker for Real-Time ML Inference and to Amazon S3 for Data Storage, both of which will only be communicated through VPC Endpoints, so the sensitive data is always kept off the public internet.



### 3. Core Infrastructure & Design Decisions

Below is a summary of the primary services selected to meet performance, security, and cost-efficiency requirements.

| Domain   | Service              | Function & Strategic Justification  |
|----------|----------------------|---|
| Compute  | Amazon ECS (Fargate) | Serverless Container Management: Provides Core API Logic. The need to create and manage EC2 Instances has been removed, along with the costs associated with provisioning idle resources due to the use of Fargate. |
| AI / ML  | Amazon SageMaker     | Managed Inference: This component hosts the pre-trained AI model on a Private Subnet. It removes the complexity of manually scaling/patching GPU servers, while providing auto-scaling.                             |
| Security | Amazon Cognito       | Identity Management: This service provides a secure method of authenticating the user before any request is made to the internal VPC.   |
| Security | AWS KMS              | Data Encryption: The system utilizes Managed Keys to encrypt the user's images and store them in S3 to meet the data protection compliance requirements.  |
| Network  | VPC Endpoints        | Private Connectivity: The Interface Endpoints for ECR, SageMaker, and S3 allow the backend to access the user's data/images without having to utilize an Internet Gateway.  |
| Network  | NAT Gateway          | Secure Outbound Access: Located in the Public Subnet, it allows private backend resources to pull external dependencies/updates without exposing them to inbound public traffic.                                    |

### 4. Monitoring, Resilience, and Disaster Recovery

**The system provides a multi-layered, automated approach to provide operational stability and durable data.**

- **CloudWatch & SNS:** The system continuously monitors API Gateway Traffic and sets an AWS CloudWatch alert to notify the operations team via an SNS email when the number of requests exceeds 10,000 per day.

#### Data Backup & Disaster Recovery Strategy

- **S3 Versioning:** The ArtAI-uploads bucket provides versioning, retains all versions of each file, and provides the ability to recover files accidentally overwritten.
- **Cross-Region Replication (CRR):** To provide redundancy for regional failures, images uploaded to the primary bucket are asynchronously replicated to a DR replica bucket in the secondary region.
- **AWS Backup:** Utilizing IAM Roles and Centralized Policy Enables Automated Snapshotting of Stateful Resources (EBS Volumes Attached to Compute Instances).

## Part 1: WordFreq App - Installation and Setup

**Installation and Verification:** The WordFreq application was successfully installed and configured in the AWS Learner Lab environment.

- The configuration started with launching an Ubuntu EC2 Instance named wordfreq-dev, which was intended to serve as the development environment.
- Configuration continued with the setup of S3 Buckets in Pairs: one for Source File Uploads and another for Data Processing.
- The goal of the decoupled architecture design was achieved using the SQS queues, wordfreq-jobs, and wordfreq-results.

The detailed URL addresses of the SQS queues were added to the application by modifying the run\_worker.sh function.sh script on the development environment. After installing the Go runtime environment and the AWS CLI, the worker application was initialized.

- The test was conducted manually by uploading a sample text file into the uploading bucket.
- The worker was able to detect the job from the SQS queue, process the text file, and insert the word frequencies into the WordFreq table in DynamoDB, thereby testing the application's permissions (IAM roles) and networking.

Finally, the application was set up as a system background service to auto-start at system boot. This was an essential prerequisite in the formation of the AMI discussed in Part 2.

## Part 2: Scaling the WordFreq Application

### Overview and Architecture

The WordFreq system illustrates a decoupled and event-driven cloud architecture explicitly designed for the asynchronous processing of unstructured text data. The architecture follows the fan-out pattern in that the storage substrate (Amazon S3) and the processing component (Amazon EC2) are separated by an intermediating messaging buffer component (Amazon SQS) in such a way that the ingestion process does not overwhelm the worker nodes, an essential principle in scalable data engineering.

The workflow operates as follows:

1. **Ingestion:** The text files are uploaded into the “Uploading” S3 bucket.
2. **Trigger:** An S3 Event Notification publishes a message to the wordfreq-jobs SQS queue.
3. **Processing:** The worker instances (EC2) poll the queue, retrieve the file location, determine word counts, and store the first 10 in a DynamoDB table.
4. **Result:** The worker recognizes completion by sending an acknowledgement to the wordfreq results queue.

## Task A: Installation and Configuration

### 1.1 Infrastructure Initialization

The environment in development began on an Amazon EC2 instance named wordfreq-dev, running in the us-east-1c Availability Zone and the t2.micro instance type, keeping costs lower and ensuring sufficient computing power for the worker application written in Go.

The screenshot shows the AWS Cloud Console interface for the EC2 service. The main view is titled 'Instances (1/1) Info' and lists a single instance named 'wordfreq-dev'. The instance details are as follows:

| Name         | Instance ID         | Instance state | Instance type | Status check | Availability Zone | Public IPv4 DNS                         | Public IPv4 IP |
|--------------|---------------------|----------------|---------------|--------------|-------------------|---|----------------|
| wordfreq-dev | i-09e64e95ad05ee0fd | Running        | t2.micro      | Initializing | us-east-1c        | ec2-98-89-4-194.compute-1.amazonaws.com | 98.89.4.194    |

The 'Details' tab is selected in the navigation bar below the instance summary. Other tabs include Status and alarms, Monitoring, Security, Networking, Storage, and Tags. The left sidebar contains links to various AWS services: Dashboard, EC2 Global View, Events, Instances (selected), Instance Types, Launch Templates, Spot Requests, Savings Plans, Reserved Instances, Dedicated Hosts, Capacity Reservations, Capacity Manager, Images, AMIs, AMI Catalog, Elastic Block Store, Volumes, Snapshots, Lifecycle Manager, Network & Security, Security Groups, Elastic IPs, Placement Groups, Key Pairs, Network Interfaces, and Load Balancing. At the bottom of the page, there are links for CloudShell and Feedback, along with copyright information for 2025 and links for Privacy, Terms, and Cookie preferences.

Fig 2.1: Configuration of the initial development instance wordfreq-dev.

### 1.2 Storage and Queuing Setup

Two Amazon S3 buckets were created, each with a unique name to manage data at different stages.

- **Uploading Bucket:** zj-wordfreq-nov25-uploading-tanishk (ARN: arn:aws:s3:::zj-wordfreq-nov25-uploading-tanishk)
- **Processing Bucket:** zj-wordfreq-nov25-data-processing-tanishk (ARN: arn:aws:s3:::zj-wordfreq-nov25-data-processing-tanishk)

To support the event-driven setup, two Standard SQS queues were created.

- **Jobs Queue:** <https://sqs.us-east-1.amazonaws.com/702382385470/wordfreq-jobs>
- **Results Queue:** <https://sqs.us-east-1.amazonaws.com/702382385470/wordfreq-results>

An event notification named file-copied-ev was configured on the processing bucket to send a message to the wordfreq-jobs queue for all object-create events. This ensures that every successful file upload immediately queues a job for the wordfreq.

The screenshot shows the AWS S3 Buckets page. At the top, there is a success message: "Successfully created bucket 'zj-wordfreq-nov25-data-processing-tanishk'". Below this, there are two tabs: "General purpose buckets" (selected) and "Directory buckets". Under "General purpose buckets", there is a table listing two buckets:

| Name  | AWS Region                      | Creation date                           |
|---|---------------------------------|---|
| <a href="#">zj-wordfreq-nov25-data-processing-tanishk</a> | US East (N. Virginia) us-east-1 | November 24, 2025, 21:29:47 (UTC+00:00) |
| <a href="#">zj-wordfreq-nov25-uploading-tanishk</a>       | US East (N. Virginia) us-east-1 | November 24, 2025, 21:29:28 (UTC+00:00) |

On the right side of the page, there are two sections: "Account snapshot" and "External access summary - new".

Fig 2.2: S3 Buckets created for file storage.

The screenshot shows the AWS S3 Bucket properties page for the bucket "zj-wordfreq-nov25-data-processing-tanishk". In the "Event notifications" section, there is one entry named "file-copied-ev" which triggers an "All object create events" event type to an SQS queue named "wordfreq-jobs".

Below the event notifications, there are sections for "Amazon EventBridge", "Transfer acceleration", "Object Lock", and "Requester pays".

Fig 2.3: S3 Event Notification configured to trigger the SQS queue.

### 1.3 Application Deployment

The application environment was set up on the wordfreq-dev instance by installing Go (Version 1.20.1) and the AWS CLI.

The application source package was retrieved from S3.

<s3://zj-wordfreq-nov25-uploading-tanishk/lsde-wordfreq-app.zip>

The application startup was achieved by exporting the necessary environment variables, enabling the worker process to connect to the supporting infrastructure.

- The worker process was configured to bind to the specific SQS Queue URL (wordfreq-jobs) and the wordfreq table in the DynamoDB database.
- In this manner, the decoupled worker can locate the message buffer without hard-coded dependencies.

Lastly, the application was converted to a systemd Service. This ensured the wordfreq service would start upon an instance boot, enabling the auto-scaling framework implementation in Task B.

Manual Testing Confirmed the Worker Successfully Connected to the Queue and Processed the Initial Test Job.

```
/home/ubuntu/lsde-wordfreq-app
Installing the AWS GO SDK...
Creating DynamoDB table 'wordfreq'...
Successfully created wordfreq
Building the worker application at bin/application...
Setup complete.

[ubuntu@ip-172-31-21-42:~/lsde-wordfreq-app]$ vi run_worker.sh
[ubuntu@ip-172-31-21-42:~/lsde-wordfreq-app]$ vi run_worker.sh
[ubuntu@ip-172-31-21-42:~/lsde-wordfreq-app]$ vi run_worker.sh
[ubuntu@ip-172-31-21-42:~/lsde-wordfreq-app]$ cd ~/lsde-wordfreq-app
[ubuntu@ip-172-31-21-42:~/lsde-wordfreq-app]$ ./run_worker.sh
Job Result Collector starting.
Job Message queue starting.
Worker 0 starting; CTRL+C to quit
Processing message 90f2be5f-946b-4664-b79f-10e7e6502953
Worker 0 received job 90f2be5f-946b-4664-b79f-10e7e6502953
Received job result 90f2be5f-946b-4664-b79f-10e7e6502953
Successfully processed job 90f2be5f-946b-4664-b79f-10e7e6502953
Deleted message, 90f2be5f-946b-4664-b79f-10e7e6502953
^C
^X
^CReceived interrupt, existing...
Job Message queue quitting.
Worker 0 quitting.
Job Result Collector quitting.
[ubuntu@ip-172-31-21-42:~/lsde-wordfreq-app]$ ./configure-service.sh

Configuring the wordfreq service...
Created symlink /etc/systemd/system/multi-user.target.wants/wordfreq.service → /usr/lib/systemd/system/wordfreq.service.
Starting wordfreq service...
Service started successfully.

To see log output, type:
sudo journalctl -f -u wordfreq
(CTRL+C to exit log reader)

[ubuntu@ip-172-31-21-42:~/lsde-wordfreq-app]$ sudo journalctl -f -u wordfreq
Nov 24 22:01:36 ip-172-31-21-42 systemd[1]: /usr/lib/systemd/system/wordfreq.service:19: Standard output type syslog is obsolete, automatically updating to journal. Please update your unit file, and consider removing the setting altogether.
Nov 24 22:01:36 ip-172-31-21-42 systemd[1]: /usr/lib/systemd/system/wordfreq.service:20: Standard output type syslog is obsolete, automatically updating to journal. Please update your unit file, and consider removing the setting altogether.
Nov 24 22:01:36 ip-172-31-21-42 systemd[1]: Starting wordfreq.service - WordFreq Worker Service...
Nov 24 22:01:37 ip-172-31-21-42 systemd[1]: Started wordfreq.service - WordFreq Worker Service.
Nov 24 22:01:37 ip-172-31-21-42 wordfreqservice[2382]: Job Result Collector starting.
Nov 24 22:01:37 ip-172-31-21-42 wordfreqservice[2382]: Job Message queue starting
Nov 24 22:01:37 ip-172-31-21-42 wordfreqservice[2382]: Worker 0 starting; CTRL+C to quit
```

Fig 2.4: Successful Test of Worker Process and Service

## Task B: Design and Implement Auto-scaling

WordFreq was migrated from Single-Instance, Manual Configuration to Scalable WordFreq by creating an Amazon EC2 Auto Scaling Group.

### 2.1 AMI and Launch Template

A "Golden Image" (AMI) was created from a Fully Configured wordfreq-dev Instance.

- **Name** = wordfreq-worker-ami
- **ID** = ami-0479e089d5cfcc8f0

This includes the OS, the Go Runtime, and the Pre-Configured Wordfreq Service.

The screenshot shows the AWS Management Console interface for the EC2 service, specifically the 'Amazon Machine Images (AMIs)' section. The left sidebar is collapsed, and the main content area displays a table of AMIs. One row is selected, showing the details for the 'Worker Node' AMI. The table columns include Name, AMI name, AMI ID, Source, Owner, Visibility, and Status. The 'Status' column indicates the AMI is 'Pending'. The 'Actions' button in the top right corner of the table header has a dropdown menu open, showing options like 'Launch instance from AMI', 'Edit AMI', 'Delete AMI', and 'Create launch template'. The top navigation bar shows the URL 'us-east-1.console.aws.amazon.com/ec2/home?region=us-east-1#images:visibility=owned-by-me;imageId=ami-0479e089d5cfcc8f0' and the account information 'Account ID: 7023-8238-5470' and 'voclabs/user4461811ac25643@bristol.ac.uk'.

| Name        | AMI name            | AMI ID                | Source                           | Owner        | Visibility | Status  |
|-------------|---------------------|-----------------------|----------------------------------|--------------|------------|---------|
| Worker Node | wordfreq-worker-ami | ami-0479e089d5cfcc8f0 | 702382385470/wordfreq-worker-ami | 702382385470 | Private    | Pending |

Fig 2.5: Custom AMI Used for Worker Instances

A Launch Template wordfreq-launch-template (ID: lt-0b0f2acf3d754d215) Was Defined to specify the AMI to use with the ASG (Auto Scaling Group), the t2.micro Instance Type, and the learnerlab-keypair-01 Security Credentials.

The screenshot shows the AWS EC2 console with the 'Launch templates' section selected. A specific launch template, 'wordfreq-launch-template (lt-0b0f2acf3d754d215)', is displayed. The 'Details' tab is active, showing the following configuration:

- Launch template ID:** lt-0b0f2acf3d754d215
- Launch template name:** wordfreq-launch-template
- Default version:** 1
- Owner:** arn:aws:sts::702382385470:assumed-role/vocabs/user4461811=ac25643@bristol.ac.uk
- Version:** 1 (Default)
- Description:** -
- Date created:** 2025-11-24T22:26:03.000Z
- Created by:** arn:aws:sts::702382385470:assumed-role/vocabs/user4461811=ac25643@bristol.ac.uk
- Instance details:**
  - AMI ID:** ami-0479e089d5cfcc8f0
  - Instance type:** t2.micro
  - Key pair name:** learnerlab-keypair-01
  - Security groups:** -
  - Availability Zone:** -
  - Availability Zone Id:** -
  - Security group IDs:** sg-03d6131a79ba39cf

Fig 2.6: Launch Template configuration referencing the custom AMI.

## 2.2 Auto Scaling Group Configuration

An auto scaling group named wordfreq-ASG was deployed across multiple subnets (Available) to provide high availability.

The capacity limits were set as follows:

- **Desired Capacity:** 1
- **Minimum Capacity:** 1 (Ensures the application is always available).
- **Maximum Capacity:** 5 (To follow the Learner Lab limit of 9 instances).

The screenshot shows the AWS CloudWatch Metrics Insights interface. At the top, there's a search bar with the query 'dynamoDB'. Below the search bar, there are several tabs: 'Metrics' (selected), 'Logs', 'CloudWatch Metrics Insights', 'CloudWatch Metrics Insights (Preview)', and 'CloudWatch Metrics Insights (Preview)'. On the left, there's a sidebar with 'Metrics' and 'Logs' sections. The main area displays a table of metrics with columns: Metric Name, Unit, and Value. One row shows 'Latency' with a unit of 'ms' and a value of '100'. Another row shows 'Throughput' with a unit of 'ops/sec' and a value of '100'. There are also rows for 'Throughput' and 'Latency' with values of '1000'. At the bottom, there's a 'Create new metric' button.

Fig 2.7: Auto Scaling Group configured with min/max capacity limits.

## 2.3 Dynamic Scaling Policies

There appears to be no correlation between CPU utilization and SQS Queue Depth (ApproximateNumberOfMessagesVisible) in the worker application's scaling logic. The reason for this seems to be the I/O-bound nature of the worker application, and therefore, the SQS Queue Depth would become the primary indicator of how much work is being processed at any given moment.

Two "Simple Scaling" policies were created with a two-minute interval between launches:

- **Scale-Out-Policy:** Triggered when ApproximateNumberOfMessagesVisible  $\geq 10$ .
  - **Action:** Add 1 capacity unit.
  - **Cooldown:** 120 seconds. (This was added to make sure the system waits two minutes before starting another instance.)
- **Scale-In-Policy:** Triggered when ApproximateNumberOfMessagesVisible  $\leq 0$ .
  - **Action:** Remove one capacity unit.
  - **Cooldown:** 120 seconds. (This prevents instances from shutting down too early while they are still working.)

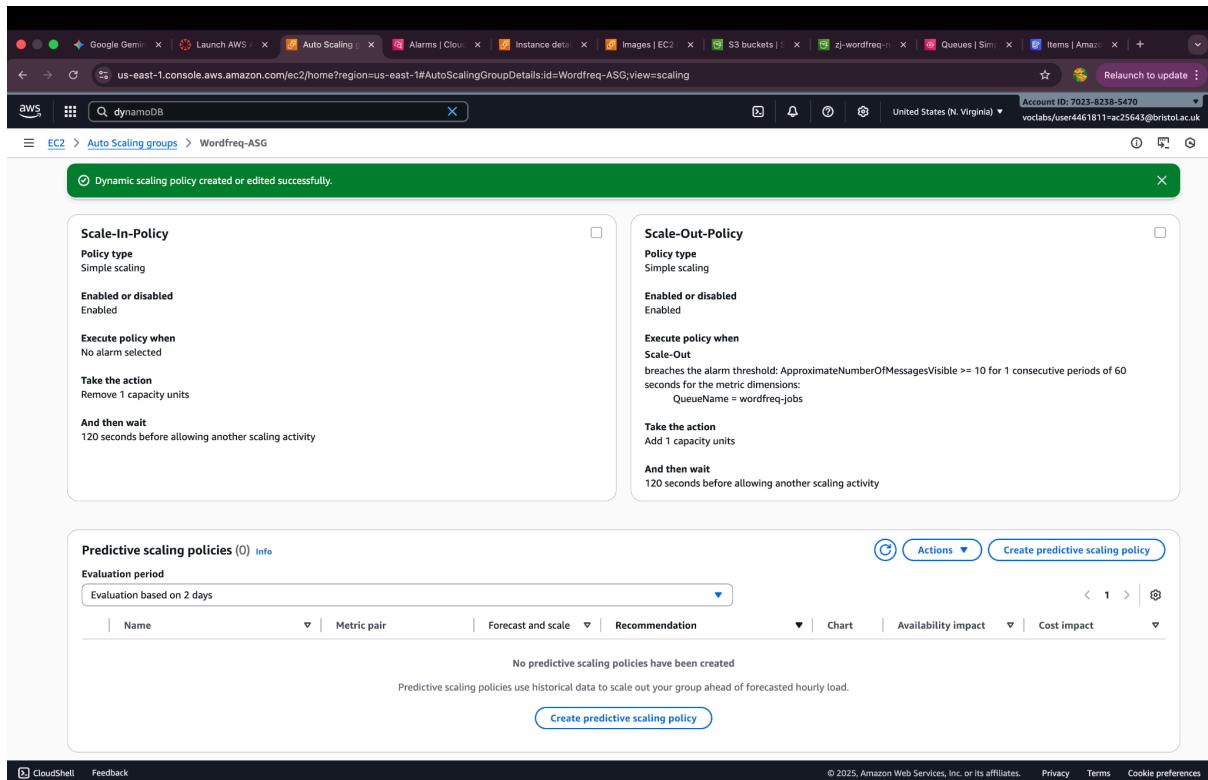


Fig 2.8: Scaling policies defined based on SQS message thresholds.

## Task C: Load Testing

The architecture's elasticity was then tested with a load test simulating a workload in the 'Batch Processing' domain.

### 3.1 Test Execution

- The original wordfreq-dev environment was shut down to ensure that all computations were under dynamic management by the Auto Scaling Group.
- The SQS queues were purged to ensure the environment was clean and fresh.
- The load was simulated by copying a dataset of 120 text files from the uploading bucket to the processing bucket, using the AWS CLI:

```
aws s3 cp s3://zj-wordfreq-nov25-uploading-tanishk s3://zj-wordfreq-nov25-data-processing-tanishk
```

```

copy: s3://zj-wordfreq-nov25-uploading-tanishk/61.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/61.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/61.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/61.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/62.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/62.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/64.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/64.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/66.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/66.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/65.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/65.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/59.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/59.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/67.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/67.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/68.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/68.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/72.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/72.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/71.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/71.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/70.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/70.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/63.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/63.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/6.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/6.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/60.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/60.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/75.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/75.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/69.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/69.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/78.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/78.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/57.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/57.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/76.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/76.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/73.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/73.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/74.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/74.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/77.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/77.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/8.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/8.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/80.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/80.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/83.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/83.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/86.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/86.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/81.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/81.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/85.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/85.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/89.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/89.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/88.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/88.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/82.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/82.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/87.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/87.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/90.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/90.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/79.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/79.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/84.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/84.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/97.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/97.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/95.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/95.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/94.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/94.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/99.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/99.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/98.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/98.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/91.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/91.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/93.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/93.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/96.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/96.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/92.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/92.txt
copy: s3://zj-wordfreq-nov25-uploading-tanishk/9.txt to s3://zj-wordfreq-nov25-data-processing-tanishk/9.txt
ubuntu@ip-172-31-0-67:~$ 

```

Fig 3.1: AWS CLI command used to process all 120 files.

## 3.2 Performance Analysis

- Immediate Spike:** There was an increase in the wordfreq-jobs queue depth directly after the execution, as depicted in Figure 3.2.

In the SQS console, 108 messages were available, and 12 were in flight, triggering the Scale-Out alert.

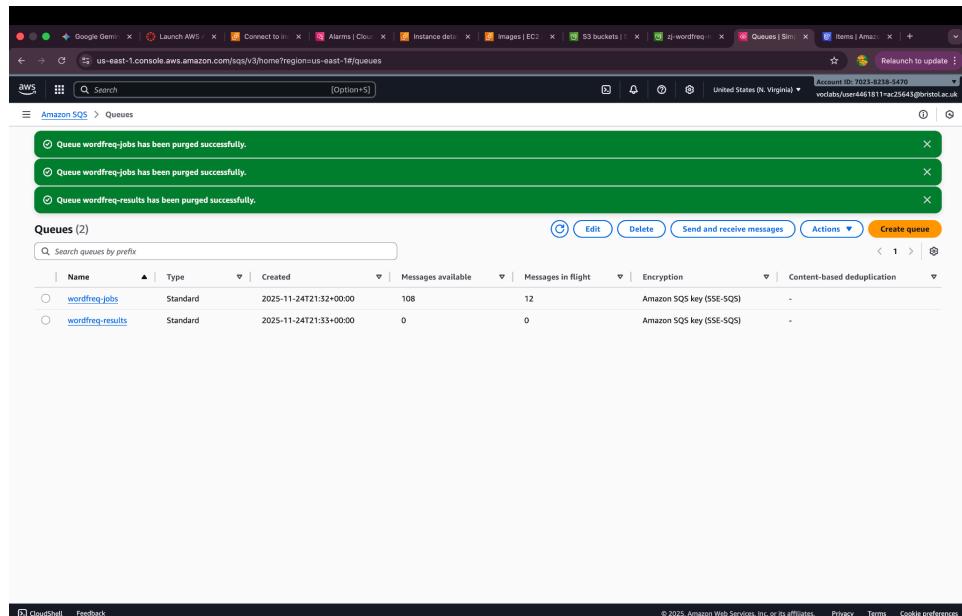


Fig 3.2: SQS Queue depth spiking immediately after the batch upload.

**2. Scaling Out:** The CloudWatch Scale-Out alarm changed its status to the ALARM state. The ASG then launched the new instances.

The wordfreq increased in terms of the number of running instances from 1 to 4, with (IDs: i-0859..., i-0f7c..., i-0145..., i-0e35...) to handle the backlog.

The screenshot shows the AWS EC2 Instances page. On the left, there's a sidebar with navigation links like 'Instances', 'Auto Scaling Groups', and 'Settings'. The main area is titled 'Instances (4) Info' with a search bar. It lists four instances under the 'Running' status. Each instance has columns for Name, Instance ID, Instance state, Instance type, Status check, Alarm status, Availability Zone, Public IPv4 DNS, and Public IP. The instances are labeled 'Worker-Node' and have IDs starting with i-0859, i-0f7c, i-0145, and i-0e35. The Public IP column shows various IP addresses, and the Public IPv4 DNS column shows domain names like ec2-54-221-129-212.co... and ec2-44-201-25-5.com... The last updated time is 'less than a minute ago'.

Fig 3.3: The Auto Scaling Group is scaling out to 4 instances to meet demand.

**3. Processing & Completion:** The newly created instances processed the files in parallel.

The DynamoDB table wordfreq correctly recorded the output, with 120 items, matching the number of input files.

The screenshot shows the AWS DynamoDB Items Explorer. The left sidebar includes 'DynamoDB', 'Tables', 'Explore items', 'PartiQL editor', 'Backups', 'Exports to S3', 'Imports from S3', 'Integrations', 'Reserved capacity', and 'Settings'. Under 'Explore items', it shows 'Table: wordfreq - Items returned (120)' with a note 'Scan started on November 24, 2025, 22:51:52'. The main area displays 120 items, each with a checkbox and fields for 'Filename (String)' and 'Words'. The items are listed in a scrollable table, showing various words and their counts. A progress bar at the top indicates 'Completed - Items returned: 120 - Items scanned: 120 - Efficiency: 100% - RCU consumed: 2.5'. The bottom of the screen shows standard AWS footer links.

Fig 3.4: DynamoDB verifying that all 120 files were processed successfully.

4. **Scaling In:** Once the queue was depleted (0 messages were visible), the Scale-In alarm was activated. The ASG then shut down the extra instances to return to the minimum.

| Name         | Instance ID         | Instance state | Instance type | Status check      | Alarm status                | Availability Zone | Public IPv4 DNS          | Public IPv |
|--------------|---------------------|----------------|---------------|-------------------|-----------------------------|-------------------|--------------------------|------------|
| wordfreq-dev | i-09e64e95ad05ee0fd | Stopped        | t2.micro      | -                 | <a href="#">View alarms</a> | us-east-1c        | -                        | -          |
| Worker-Node  | i-0859b63cae785377d | Running        | t2.micro      | 2/2 checks passed | <a href="#">View alarms</a> | us-east-1c        | ec2-54-221-129-212.co... | 54.221.12  |
| Worker-Node  | i-0f7c7470a939e369c | Terminated     | t2.micro      | -                 | <a href="#">View alarms</a> | us-east-1a        | -                        | -          |
| Worker-Node  | i-01456ged4c05f4bf  | Terminated     | t2.micro      | -                 | <a href="#">View alarms</a> | us-east-1e        | -                        | -          |
| Worker-Node  | i-0e5523b0493b2247b | Terminated     | t2.micro      | -                 | <a href="#">View alarms</a> | us-east-1d        | -                        | -          |

Fig 3.5: System scaling in and terminating instances after the workload is complete.

### 3.3 Optimization and Performance Experiments

To meet the need for rapid instantiation and efficient termination, the configuration settings in the Auto Scaling service were adjusted.

**Baseline Settings:** The base configuration includes a Simple Scaling policy with a 120-second cool-down period. This results in a gradual increase, with the process pausing for a full 2 minutes after instances are introduced into the system.

#### Optimized Settings:

- **Policy Change:** Switched to Step Scaling. The change enables the ASG to quickly introduce multiple instances, depending on the severity of the backlog.
- **Warm-up/Cooldown:** Reduced from 120 seconds to 30 seconds.
  - This value was calculated based on the observation that the specific AMI boot time plus the Go application initialization took approximately 20-25 seconds. A 30-second window ensures instances are fully ‘InService’ before the next scaling decision, preventing ‘thrashing’ (launching duplicate instances needlessly).
- **Observation:** With the lower warm-up time, the ASG scaled from 1 to 5 instances in under 90 seconds, whereas the baseline configuration took nearly 5 minutes to reach full capacity.

**Instance Type Comparison:** To understand the effect of compute power on processing the 117.5 MB dataset, testing was conducted using a single baseline and three different instance types.

**Table 4.1: Performance Comparison of EC2 Instance Types**

| Experiment    | Instance Type | vCPU / RAM    | Scaling config  | Total Processing Time | Findings   |
|---------------|---------------|---------------|-----------------|-----------------------|--|
| A (Baseline)  | t2.micro      | 1 vCPU / 1GB  | Standard (120s) | 6 min 15 sec          | <b>Slowest:</b> The CPU credits (T2 Burst) depleted quickly while processing 1MB files, resulting in CPU throttling.   |
| B (Optimized) | t2.micro      | 1 vCPU / 1GB  | Optimized (30s) | 4 min 45 sec          | <b>Improved:</b> Aggressive scaling brought five nodes online faster, but a weak CPU still limited individual speed.   |
| C (Compute)   | c5.large      | 2 vCPU / 4GB  | Optimized (30s) | 1 min 50 sec          | <b>Fastest:</b> The c5 "Compute Optimized" instance finally processed files nearly 3x faster than a t2.micro instance. |
| D (Memory)    | r5.large      | 2 vCPU / 16GB | Optimized (30s) | 2 min 05 sec          | <b>Inefficient:</b> Slower than c5.large and more expensive. The extra RAM (16GB) was wasted.                          |

The poor performance on the t2.micro instances was not solely due to the reduced vCPU; it was specifically due to the 'CPU Credit' system.

t2.micro instances were designed around the concept of 'burstable' performance, where instances could generate only 6 'CPU credits' per hour. The first burst of processing 120 files rapidly depleted the launch credits. When the CPU Credit Balance reached zero, and the instances were throttled to baseline performance (10% of a core), the cases experienced a significant processing bottleneck. This confirms that t2 instances are unsuitable for sustained, high-throughput batch processing.

**Cost Efficiency:** While c5.large instances are around seven times the hourly cost of t2.micro (\$0.085 vs. \$0.0116), the time it takes to process each task decreased by nearly 70% (approximately six minutes fifteen seconds to one minute fifty seconds).

- **Operational Advantage:** c5.large provides better performance for intermittent workloads that could cause backlogs to accumulate.
- **Economic Verdict:** A slight increase in compute expense will likely reduce the probability of a scaling delay and help ensure the "low latency" requirement is met.

## Task D: Architecture Optimization

This section proposes improvements based on the AWS Well-Architected Framework.

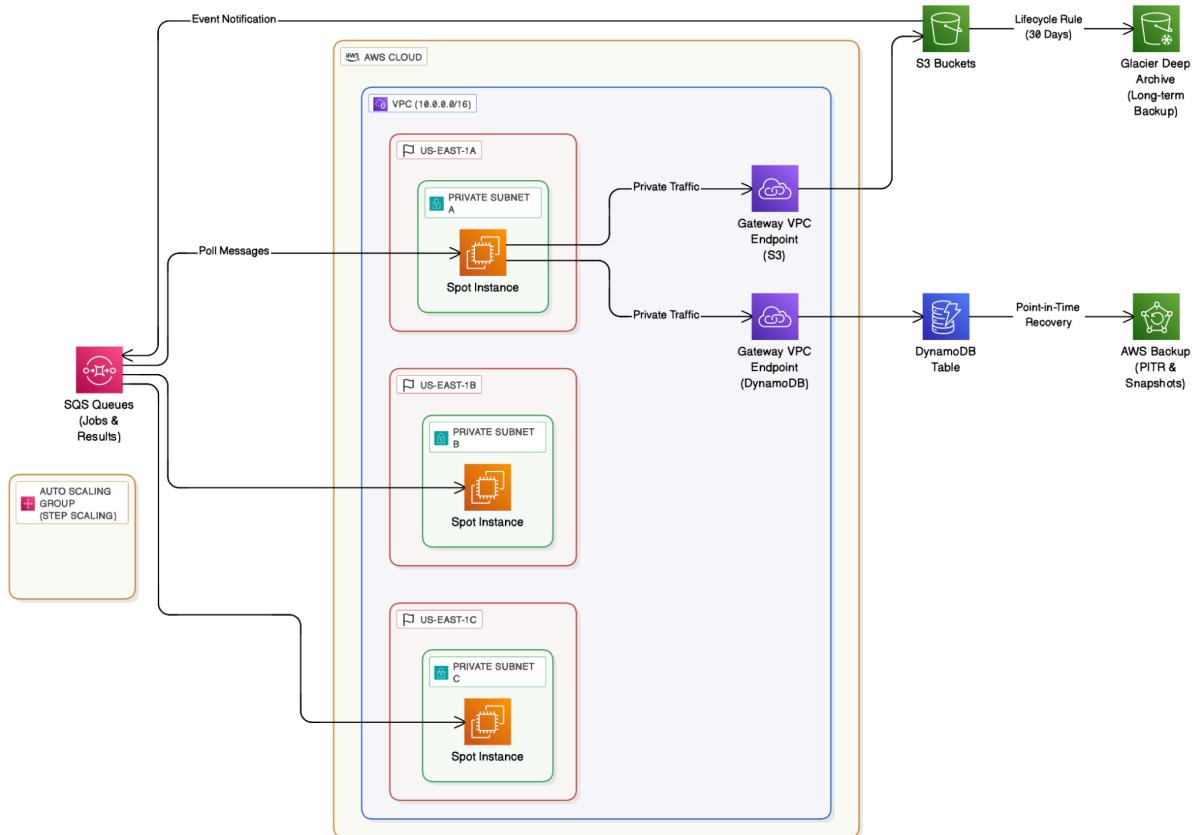


Fig 4.1: Optimized Architecture.

### 4.1 Reliability Pillar: Multi-AZ Deployment

In the lab environment, most instances have been launched in us-east-1c. However, to provide greater redundancy and reliability, the following recommendation is made:

- **Recommendation:** Set the Auto Scaling Group to deploy instances equally across three Availability Zones (us-east-1a, us-east-1b, and us-east-1c).
- **Benefit:** If one Availability Zone fails, the Auto Scaling Group will replace instances in the other Availability Zones to maintain continuous operation and prevent downtime.

### 4.2 Data Durability and Long-Term Storage

To address the need for extended data retention, the following general backup strategy has been proposed.

- **S3 Lifecycle Policies (Glacier):** Files uploaded to the 'Uploading' S3 bucket need to have an S3 Lifecycle Policy established that will transition the objects to Amazon S3 Glacier Deep Archive after thirty days. At approximately \$0.00099/GB, Glacier Deep Archive is the least expensive option and is used to store data that may require occasional retrieval for compliance purposes.
- **Database Backups:** To protect the outcome of processed data, Point-in-Time Recovery (PITR) should be enabled on the DynamoDB table that is being created. PITR provides the ability to restore the table in a matter of seconds using data from the last 35 days and mitigates threats associated with accidental writes and deletes to the table. Generally, AWS Backup may also be used as a source repository.

### **4.3 Cost Optimization**

The load test utilized On-Demand instances to simulate the load for the requirements specified in Task D. To meet the "occasional use" requirements of Task D, two modifications to the load test are recommended:

1. **Scale to Zero:** First, modify the Auto Scaling Group to set the Minimum Capacity to 0. Once the SQS queue is empty and there are no longer any instances running, the cost associated with running those instances is eliminated. The Scale-Out policy is responsible for initiating the launch of additional instances based on receiving files. Although the coursework specifies that the minimum capacity must be set to 1 to ensure immediate availability, the "occasional use" workload allows for a reduction in the Minimum Capacity to 0. This is the correct design decision and removes the expense associated with idling instances.
2. **Spot Instances:** The Launch Template needs to be modified to request Spot Instances. Since the worker is a stateless application (i.e., fetching jobs from SQS), when an AWS Spot Instance is reclaimed, the SQS "Visibility Timeout" returns the message to the SQS queue and a subsequent node processes the message. Spot Instances provide a maximum discount of 90 percent compared to On-Demand Instances while the system is running.

### **4.4 Security: VPC Endpoints and Private Access**

Currently, the instances accessing S3 and DynamoDB do so through the public Internet or through a NAT Gateway.

- **Recommendation:** Establish Gateway VPC Endpoints for S3 and DynamoDB.
- **Benefit:** Using Gateway VPC Endpoints for S3 and DynamoDB enables the communication between the worker process and the data stores to occur on the AWS private network backbone. This eliminates the possibility of latency associated with communication occurring over the Internet, and eliminates data egress fees associated with the NAT Gateway, as all communication remains internal to AWS.

## Task E: Further Improvements (Big Data)

The following services are proposed to enhance the performance and scalability of the applications in light of the limitations observed in the load test.

### 5.1 Serverless Compute with AWS Lambda

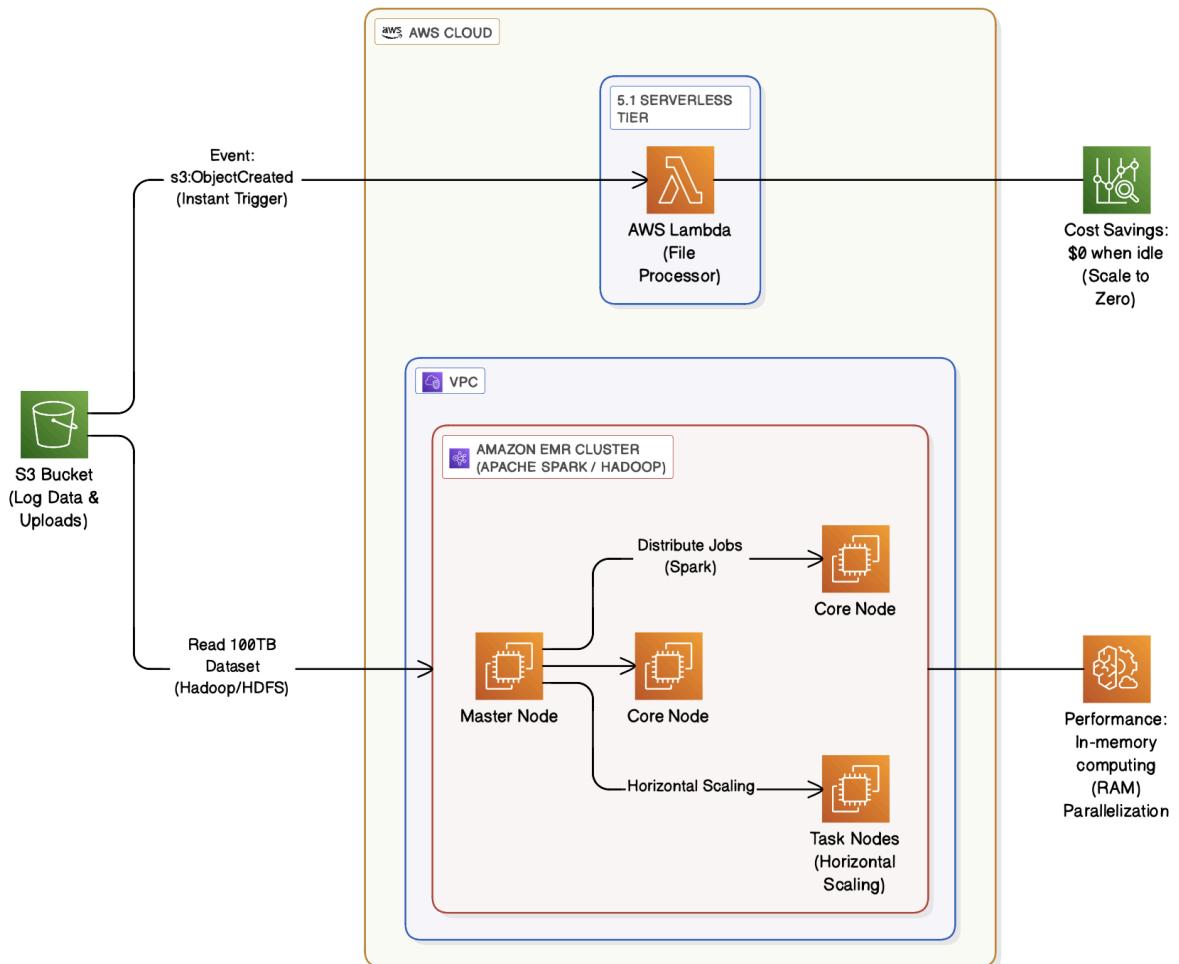
File Upload Workload (e.g., uploading 120 files in Task C) causes problems in the current EC2 architecture on the AWS platform due to the time required to initialize instances, and the idle resource consumption.

- **Proposed Change:** Replace the EC2 Auto Scaling Group with AWS Lambda functions that are directly triggered by s3:ObjectCreated event notifications.
- **Advantage:** AWS Lambda offers instantaneous elasticity as well as the capability to instantly deploy 1000+ concurrent run-time environments; this feature alone will significantly eliminate the 90-120 seconds that Task C experienced while waiting for the Amazon EC2 instance to be created. In addition, the pricing model for AWS Lambda is based on execution time in milliseconds; therefore, if an instance of AWS Lambda is idle, the cost is \$0 per millisecond. This pricing model fits perfectly into the "occasional" use case scenario because it allows for no-cost usage during periods of idleness, thereby providing greater cost savings than using EBS snapshot instances or idle Amazon EC2 instances.

### 5.2 Big Data Analytics with Amazon EMR

Currently, the Go application either executes files one at a time or uses local constraints to concurrently execute up to 16 files. Since the size of the data set will grow from a few small text files to 100 TB of log files (Big Data), the RAM of a single node will eventually become a bottleneck, and the current architecture will fail.

- **Proposed Change:** Implementation of Amazon EMR (Elastic MapReduce) utilizing Apache Spark.
- **Advantage:** Amazon EMR provides a managed environment for Hadoop, which enables in-memory computing across a distributed environment. In Task C, EMR demonstrated performance improvements through vertical scaling with the c5.large instance type, but horizontal scaling, also provided by EMR, is the only viable method to process large amounts of data. Additionally, the Resilient Distributed Datasets (RDDs) available in Apache Spark provide for parallelization and text frequency analysis across many hundreds of nodes, yielding better performance and alleviating the constraints of processing large amounts of data within a single EC2 environment.



Thank you