# Assignment 2: AWS Project

## COMP5349: Cloud Computing

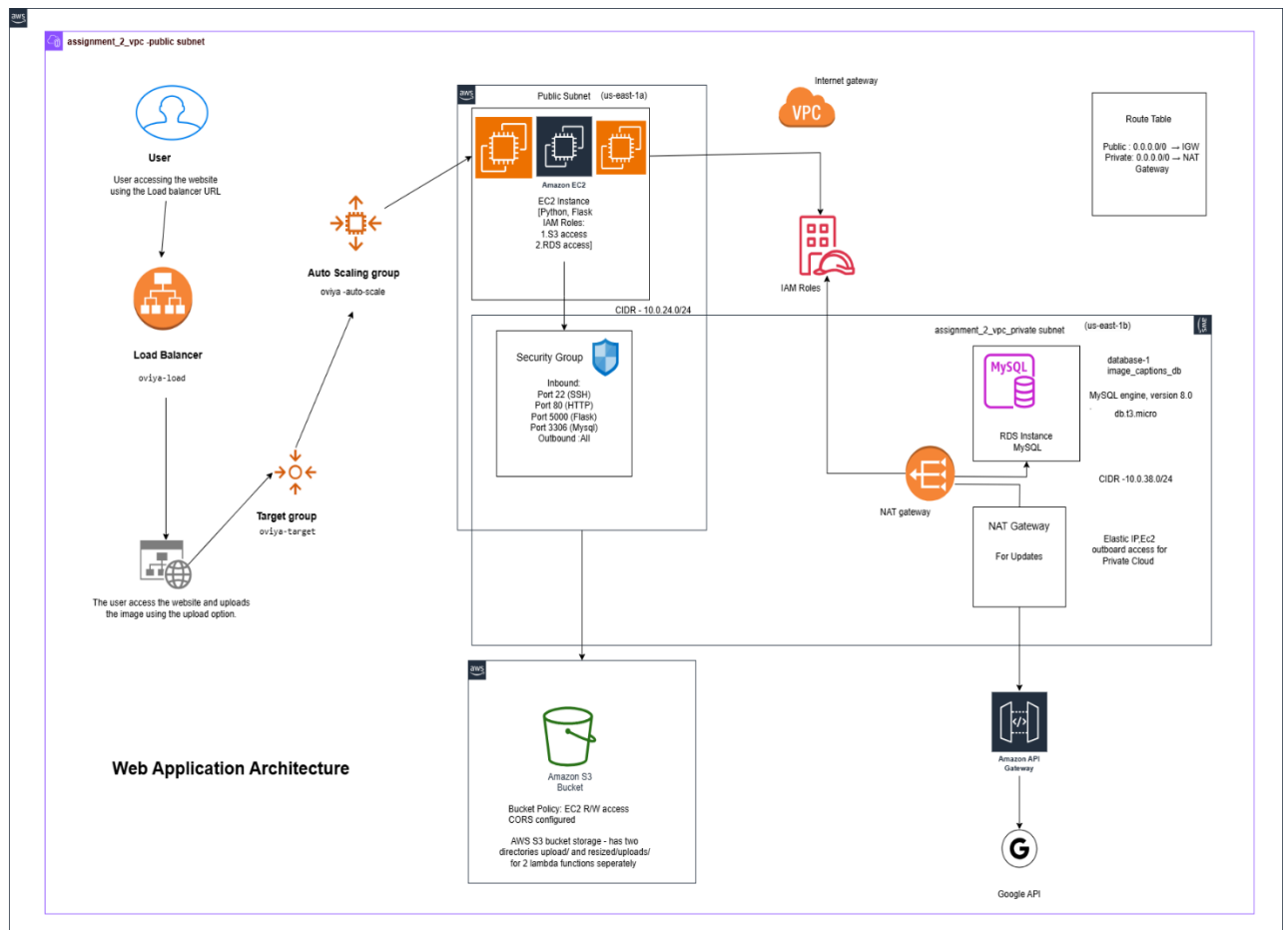SID: 540967223

Name: Oviya Balamurugan

**Introduction**

In this assignment I have demonstrated an end-to-end deployment which is featured with main functionality such as scalability, uploading the image to the AWS platform, integration of the web-based system and the modern serverless architecture which is achieved by EC2 web-based hosting in turn delivering a robust system and having a cloud native solution. One of the important goals was to upload the images from the user end through the web page and have a real time backend processing with AI generated captions and thumbnail creation. This task is accomplished by 2 lambda functions which are triggered automatically. As I have integrated Auto scaling groups, load balancer, S3, RDS and lambda I was able to attain the availability ensuring the automation and having the best security practices as well.

The website is deployment on the Ec2 instance by using the Auto Scaling groups and which is distributed by load balancer which ensures seamless scaling and traffic balance. The image storage is handled by S3 bucket along with the metadata and captions which is stored in private RDS (MySQL) instance. The deployment process comes under the secure VPC environment using both the public and private subnets so there is network ispolationa nd controlled access of the lambda functions (obal0980-lambda-captions and ovi-lambda-thumbnail) which is triggered by the S3 bucket using the Event Bridge. The layers are added in the particular lambda functions. Access control is maintained by the IAM roles, security groups, subnetting and connecting Lambda and RDS in a secure way.
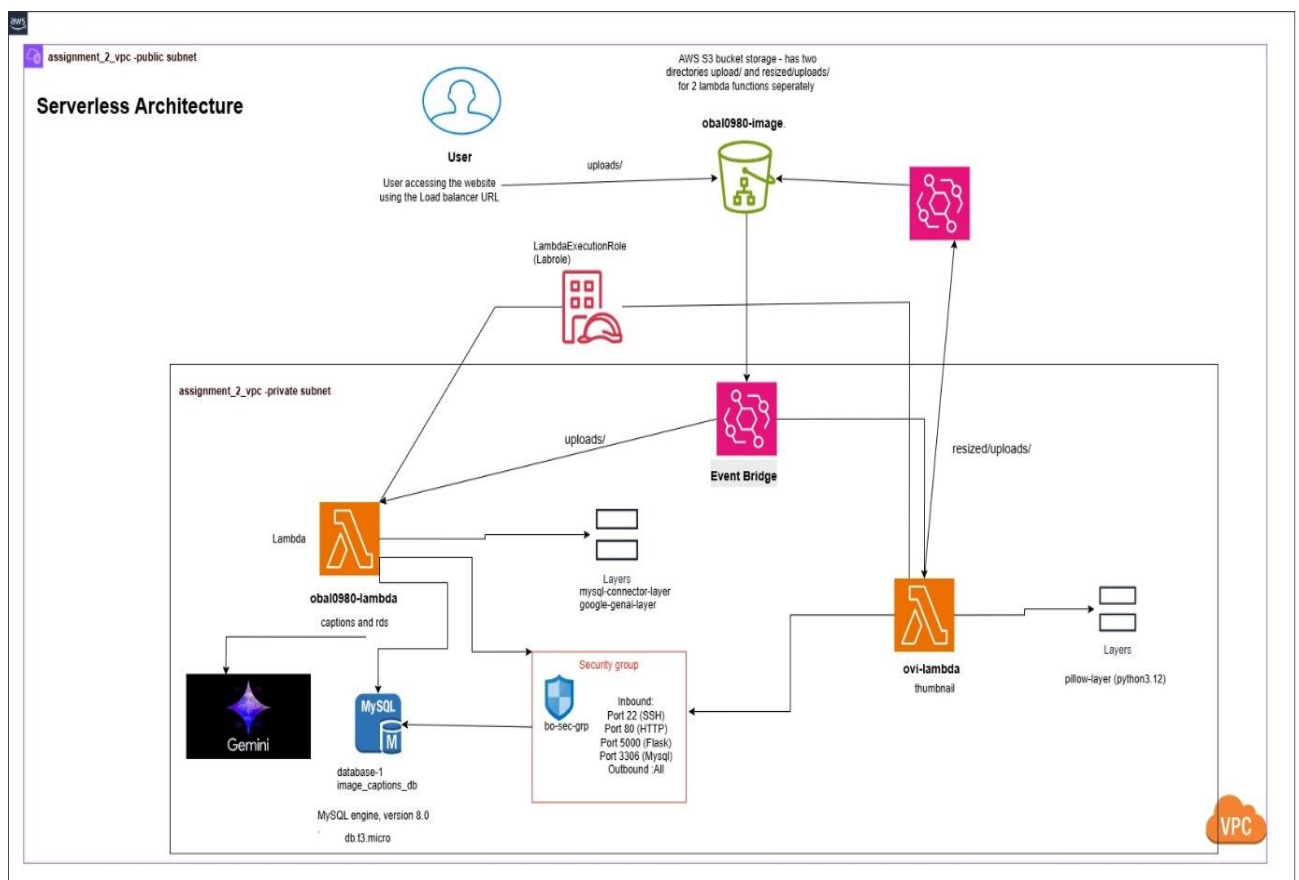
# Web Application Architecture:

The diagram effortlessly distinguishes the VPC with public and private subnets, clearly showing secure network partitioning between the web app (ALB) and the RDS database. It highlights the S3 bucket as an integration point of integration, used for uploads by users as well as an initiating point for backend Lambda processing.

## Serverless Architecture:

The diagram itself is a wonderful representation of the serverless event-driven process, showcasing how the S3 ObjectCreated events invoke the Lambda functions so gracefully through EventBridge orchestration within the secure private subnetted environment. It significantly illustrates separation of the processing tasks, with the caption-generating functionality leveraging external Gemini API integration and RDS connectivity and the thumbnail functionality being standalone with S3-to-S3 processing, both leveraging Lambda layers for optimal dependency management as well as cross-security group configuration.



## Web Application Deployment:

### Compute Environment:

### EC2 with Auto Scaling Group (ASG):

The web application is deployed on the EC2 instance and is managed by an Auto Scaling group (ASG) in my case I have named it as oviya-auto-scale. The ASG adjusts the number of instances based on the real time demand which ensures the high availability and cost effectiveness.

**Web Application Deployment:**

**Compute Environment:**

**Configuration of Auto Scaling group:**

| Parameter | Value / Description |
|---|---|
| Auto Scaling Group Name | oviya-auto-scale |
| Min / Max Capacity | 2 / 5 |
| Scaling Policy | Scale out >60% CPU, scale in <30% CPU |
| Subnets | Public subnets in us-east-1a, us-east-1b |

Launch Template: I have used this feature to setup my key features of EC2 which is detailed in the table below.

| Configuration | values |
|---|---|
| AMI ID | ami-0953476d60561c955 |
| Instance type | t2.micro |
| Key pair name | vockey |
| Security groups | bo-sec-grp |
| Security group IDs | sg-0841f40050f6f867e |

**VPC and Subnets:**

- The services are present in the assignment_2_vpc which is in the subnets public and private.
- **Public Subnets**: EC2, ALB.
- **Private Subnets**: RDS and Lambda. (Since the there is a connection between the rds and lambda)

Security groups:

| Security Group | Purpose | Key Rules |
|---|---|---|
| bo-sec-grp | Web server access | Allow HTTP (80), HTTPS (443) |
| RDS SG | DB protection | Allow MySQL (3306) from EC2/Lambda only |
| Bastion SG | Admin access | Allow SSH (22) from my personal IP |

**IAM Roles:**

- **EC2 Role:** Grants access to S3 (obal0980-image) and RDS (database-1) via Secrets Manager.

- **Lambda Role:** Allows S3 read/write, RDS access, and VPC attachment for private subnet connectivity.

- IAM roles provide temporary, fine-grained permissions that enable services to interact securely without using long-term credentials.
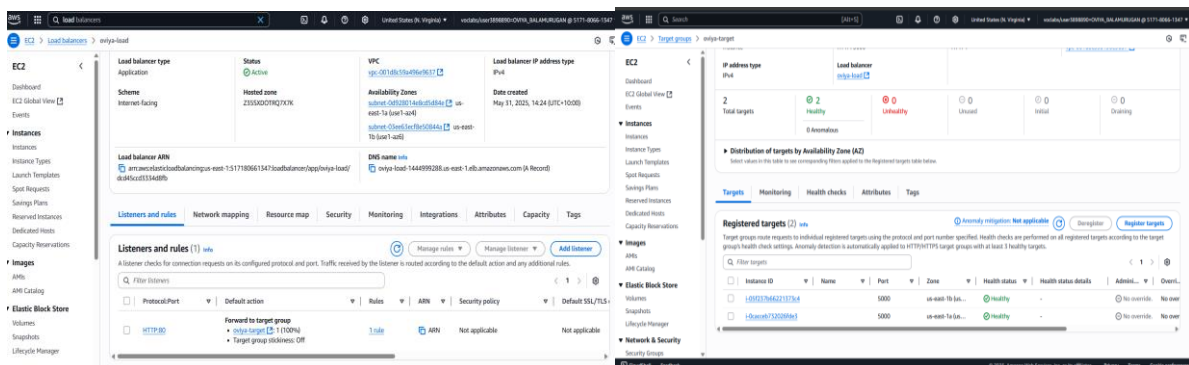
Admin Access:

The public subnet provides secure SSH access to EC2 instances and the private subnets for administration and troubleshooting.



**i-0cacceb732026fde3**
PublicIPs: 52.91.4.157   PrivateIPs: 10.0.9.127

| Component | Details |
|-----------|---------|
| Bastion Host | EC2 in public subnet |
| Key Pair | Same as web app EC2 |
| IP Restriction | 59.102.83.144/32 |

## Load Balancer Configuration

- The Application Load Balancer (ALB) listens on port 80 for HTTP traffic and directs it to a target group of EC2 instances which are spread across multiple Availability Zones, so there is no single server that gets bogged down.
- Health checks are configured so that only healthy EC2 instances get traffic, providing high availability and seamless failover for your app

| ALB Element | Configuration Details |
|---|---|
| Listener | Port 80 (HTTP) → Target Group (`oviya-target`) |
| Target Group | Registered with ASG, health checks on the endpoint |
| Subnets | Public subnets. |

## Database:

- MySQL database is provisioned through Amazon RDS, running in my own VPC's private subnets for security.
- I chose instance type (db.t3.micro), specified storage, and securely handled credentials with AWS Secrets Manager.
- Security groups permit EC2 instances and Lambda functions in the VPC to connect (3306), and public access is disabled.
- Application and Lambdas exchange information through the RDS endpoint and environment variable or Secrets Manager credentials, so that all information is contained within AWS's trusted network.

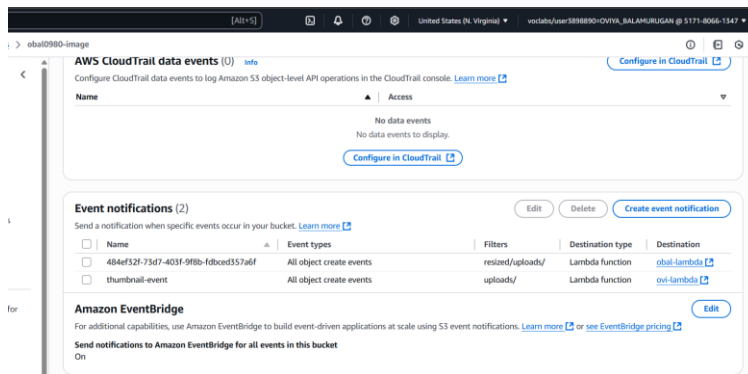| Database elements | Values |
|---|---|
| DB Engine | MySQL |
| Instance Class | db.t3.micro |
| Subnet Group | Private subnets |
| Security Group | Allows MySQL from Lambda |
| Authentication | IAM + Secrets Manager |

**Storage:**

The S3 bucket (obal0980-image) is used for storing user-uploaded images and Lambda-generated thumbnails.

| S3 Folder | Purpose |
|---|---|
| uploads/ | Raw user image uploads |
| resized/uploads/ | Thumbnails from Lambda |

- Event Triggers: S3 ObjectCreated events on uploads/ trigger two Lambda functions:

  - obal0980-lambda-captions (caption generation via Gemini API)

  - ovi-lambda-thumbnail (thumbnail creation)

Event Bridge and Event notification



**Serverless Component Deployment:**

**Event-Driven Architecture:**

The serverless components of my application are designed using an **event-driven architecture**. Specifically, the backend automation is triggered by **S3 ObjectCreated events**, which are routed via **Amazon EventBridge** to initiate two Lambda functions in response to image uploads.

Working of the Trigger**:**

- Whenever a user uploads an image to the S3 bucket obal0980-image, specifically under the uploads/ prefix, an **ObjectCreated event** is generated by S3.

- This event is picked up by **Amazon EventBridge**, which has two rules configured:

- **Rule 1** targets the Lambda function obal0980-lambda-captions, which handles image annotation.

- **Rule 2** targets the Lambda function ovi-lambda-thumbnail, which generates image thumbnails.

Resources Interacted With:

| Lambda Function | Triggered By | Resources Accessed |
|---|---|---|
| obal0980-lambda-captions | S3 ObjectCreated → EventBridge | - S3 bucket (download image) - Gemini API (generate captions), RDS MySQL (database-1), (store metadata) |
| ovi-lambda-thumbnail | S3 ObjectCreated → EventBridge | - S3 bucket (read original, write resized to resized/uploads/) |

- Both Lambda functions are configured with execution roles granting them access to the necessary AWS services securely via IAM, and obal0980-lambda-captions is additionally placed inside a private subnet within the VPC to connect to the private RDS database.

## Annotation Function:

(obal0980-lambda-captions)

- The function code is simply a .zip file with only the captions.py script, while all of the dependencies (like Pillow) are included in a separate Lambda custom layer attached to the function.
- Deployment was through the AWS Management Console (or CloudShell), and the function has proper IAM permissions to access S3 (to read images), RDS (to store captions), and the Gemini API (to make to AI-generated captions) with sensitive attributes and secret managed as environment variables and in Secrets Manager.
- The Lambda is configured to run inside the VPC for secure access to the RDS, and it is then automatically triggered on any S3 ObjectCreated event sent through EventBridge, so it creates a back-end processing system and works smoothly in an event-driven fashion.

| Attribute | Value |
|---|---|
| Runtime | Python 3.12 |
| Packaging Format | .zip archive with captions.py script |
| Deployment Method | AWS Console via CloudShell |
| Layers Used | pillow-layer for image handling |
| Event Trigger | S3 ObjectCreated to EventBridge Rule |
| IAM Role | Allows access to S3, RDS, and HTTP calls to Gemini |
| VPC Connection | Yes - private RDS access via subnet |
| Environment Variables | DB credentials (via Secrets Manager), Gemini API Key |

## Thumbnail Generator:

- The code for this Lambda is also zipped with only the thumbnail.py, while all dependencies are in an attached Lambda layer to keep the deployment package lean and maintainable.
- The function is deployed via the AWS Console with an IAM role that grants permissions to read & write to the S3 bucket so it can process uploaded images and generate thumbnails.
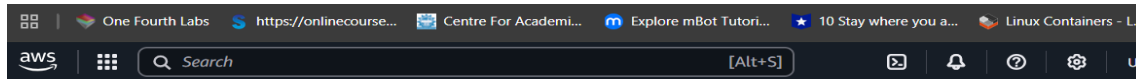
- The code is triggered by S3 ObjectCreated events (by using EventBridge), it processes the image and writes the thumbnail to a specified S3 folder, and there is no need for VPC access since it only communicates with S3.

| Attribute | Value |
|---|---|
| Runtime | Python 3.12 |
| Packaging Format | .zip archive containing thumbnail.py script |
| Deployment Method | AWS Console (uploaded ZIP from CloudShell terminal) |
| Layers Used | pillow-layer for image processing |
| Event Trigger | S3 ObjectCreated event via Amazon EventBridge |
| IAM Role | Lambda execution role with permission to read from and write to S3 |
| VPC Connection | Not required (no RDS access needed) |
| Environment Variables | S3 bucket name, thumbnail output path (e.g., resized/) |

**Auto Scaling:**

| Configuration | Value |
|---|---|
| Auto Scaling Group Name | oviya-auto-scale |
| Launch Template | AMI with Flask app, port 5000 exposed |
| Load Balancer | oviya-load (Application Load Balancer) |
| Target Group | oviya-target (Port 5000 health check) |
| Scaling Policy | Target tracking (CPU > 60%) |
| Initial Capacity | 2 instance |
| Max Capacity | 4 instances |

| Test Tool | ApacheBench (ab) |
|---|---|
| Test Command | ab -n 1000 -c 50 http://oviya-load-1444999288.us-east-1.elb.amazonaws.com/ |





Spikes were detected around 21:45–22:15, likely due to simulated load or application interactions, followed by stabilization.

The auto scaling feature worked as expected. During high traffic situations, when our application demand was particularly high, a new EC2 instance was launched automatically, and the ALB was load balancing traffic to both servers. Once demand calmed down again, the extra instance was terminated, and we confirmed that we are deploying and consuming resources efficiently and dynamically.

Test Steps

- Accessed web application via ALB endpoint.
- Ran ApacheBench: 1000 requests with concurrency of 50.
- Monitored EC2 scaling, ALB target health, and CPU usage in the AWS Console.

Outcomes & Evidence

- ASG scaled EC2 instances from 1 to 2 due to CPU threshold being exceeded.
- ALB distributed traffic equally between both healthy instances.
- ApacheBench: 1000 requests completed with very high throughput.

Target group health check:



Ec2 instance creation:

The scale in and scale out is shown as launch and terminated in this screenshot.



**Conclusion, challenges and lessons learned:**

- During this assignment, I was able to develop and deploy a scalable, event-driven image annotation platform on AWS. I did load testing using ApacheBench (ab) to see if the platform would scale and stayed observant of the system with CloudWatch. It was really cool to see the platform scale bi-directionally into/out of running services, which was almost synonymous with high availability, as well as have services utilized effectively.

- While doing this project, I did run into some hiccups along the way, like the package upload limits of Lambda, and I felt somewhat lost with EventBridge rules at the beginning. However, learning how to use Lambda layers eliminated my package size issue, and therefore deployment went a lot more smoothly. Once I learned how to use EventBridge, linking S3 events to my Lambda functions was much easier and intuitive.

- All in all, this project helped increase my confidence with AWS. I learned a lot about networking, permissions, and automating cloud workflows. More importantly, I now feel a lot more confident designing secure, cloud-native solutions that take performance, cost, and practicality into account.