# BACTERIA SEARCH ENGINE

*Cloud computing final project*
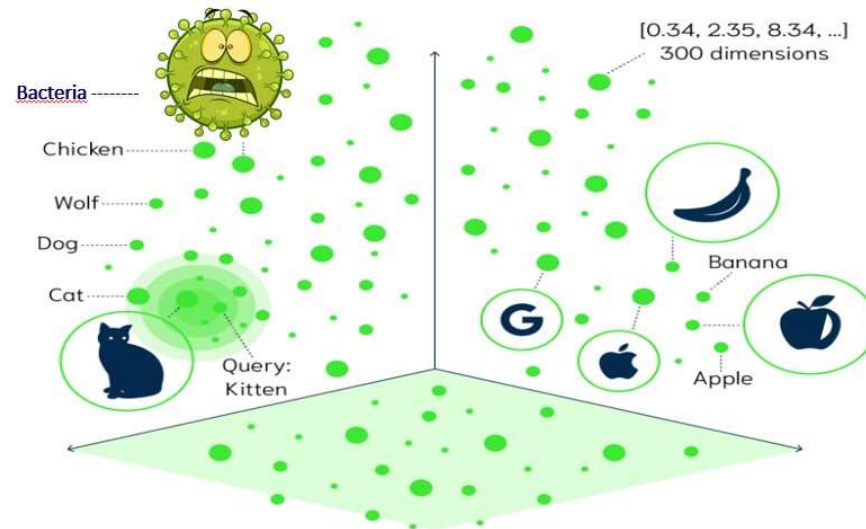
SAPIENZA
UNIVERSITÀ DI ROMA

*Matteo Migliarini  1886186*
*Giuseppe Di Poce  2072371*
*Enrico Grimaldi   1884443*

# Problem Description

## Weaviate vectorial database to implement an image search engine;

- The goal of the project is to develop an application that uses Weaviate to create a search engine for images of bacteria. *Weaviate* is known for its speed and flexibility, allowing for efficient search operations on large datasets.

- The process involves using a pre-trained neural network to extract features from the images and generate vector *embeddings.* These embeddings are then stored in the Weaviate database, enabling quick searches for similar images based on a given query image.

- When a user *uploads an image*, Weaviate computes the vector embedding for that image and compares it with the embeddings of reference images in the index to find the most similar matches.
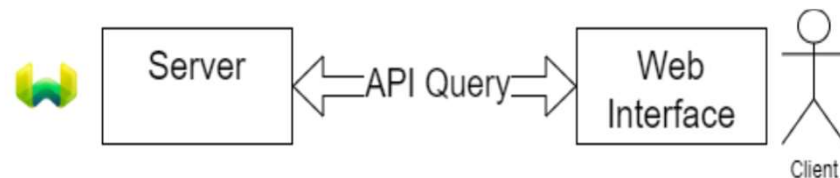
# Solution implementation

we want to implement a scalable web application, capable of rapid elasticity.

The application is divided into *microservices*:

- **Web Interface**: The client-side is a fat client, where a static HTML/CSS/JS webpage is served from an edge node. The webpage performs API calls to the server to query and retrieve results.The client performs computations to update the webpage, but it shouldn't be computationally intensive.



- **API Call**: The server-side implements a REST API, particularly using POST requests to handle the query image in the request payload and return a list of images in the response.

- **Vectorial Database**: Weaviate is chosen as the vector database. It uses vector representations of items and cosine similarity to find similar items to a query. The vector database consists of a pretrained neural network that transforms items into vectors and the database software itself, which handles queries, schema, and vector-item mappings.

# Docker Containers

- To Develop our web app we need to ensure consistent and reliable deployment across different environments.

- The application utilizes **Docker containers** for Weaviate and the pre-trained neural network **images**.

- Docker provides a lightweight and portable solution, encapsulating the entire application stack along with dependencies and configurations.

- In our implementation we will run on top of an AWS service two different Docker containers that will communicate each other (Weaviate and NN).

- Each container could be allocated on specific computing resources, <u>for example the NN needs GPU resources while the Weaviate container needs a fast memory.</u>
In this way we ensure optimal usage.

- The system can scale by deploying multiple containers in parallel, achieving high availability and load balancing.

```yaml
 3  services:
 4    weaviate:
 5      command:
 6      - --host
 7      - 0.0.0.0
 8      - --port
 9      - '8080'
10      - --scheme
11      - http
12      image: semitechnologies/weaviate:1.19.8
13      ports:
14      - 8080:8080
15      restart: on-failure:0
16      environment:
17        IMAGE_INFERENCE_API: 'http://i2v-neural:8080'
18        QUERY_DEFAULTS_LIMIT: 25
19        AUTHENTICATION_ANONYMOUS_ACCESS_ENABLED: 'true'
20        PERSISTENCE_DATA_PATH: '/var/lib/weaviate'
21        DEFAULT_VECTORIZER_MODULE: 'img2vec-neural'
22        ENABLE_MODULES: 'backup-filesystem,img2vec-neural,
23        CLUSTER_HOSTNAME: 'node1'
24        BACKUP_FILESYSTEM_PATH: '/tmp/backups'
25      volumes:
26      - /var/weaviate:/var/lib/weaviat
27      - backups-weaviate:/tmp/backups
28    i2v-neural:
29      image: semitechnologies/img2vec-pytorch:resnet50
30      environment:
31        ENABLE_CUDA: '0'
32      deploy:
33        resources:
34          reservations:
35            devices:
36              - driver: nvidia
37                count: 1
38                capabilities: [gpu]
39
40  volumes:
41    backups-weaviate:
```

# Lambda Function

- The lambda function get triggered only when an image is uploaded to the web app by the user, *through a REST request (POST).*

- The lambda will communicate by an API call with the web app and connect with Weaviate by the IPv4 address of the Application Load Balancer service.

- This decoupling allows for greater scalability and responsiveness. The Lambda function can be invoked as needed, providing real-time responses without continuous resource allocation, since it's managed by AWS.
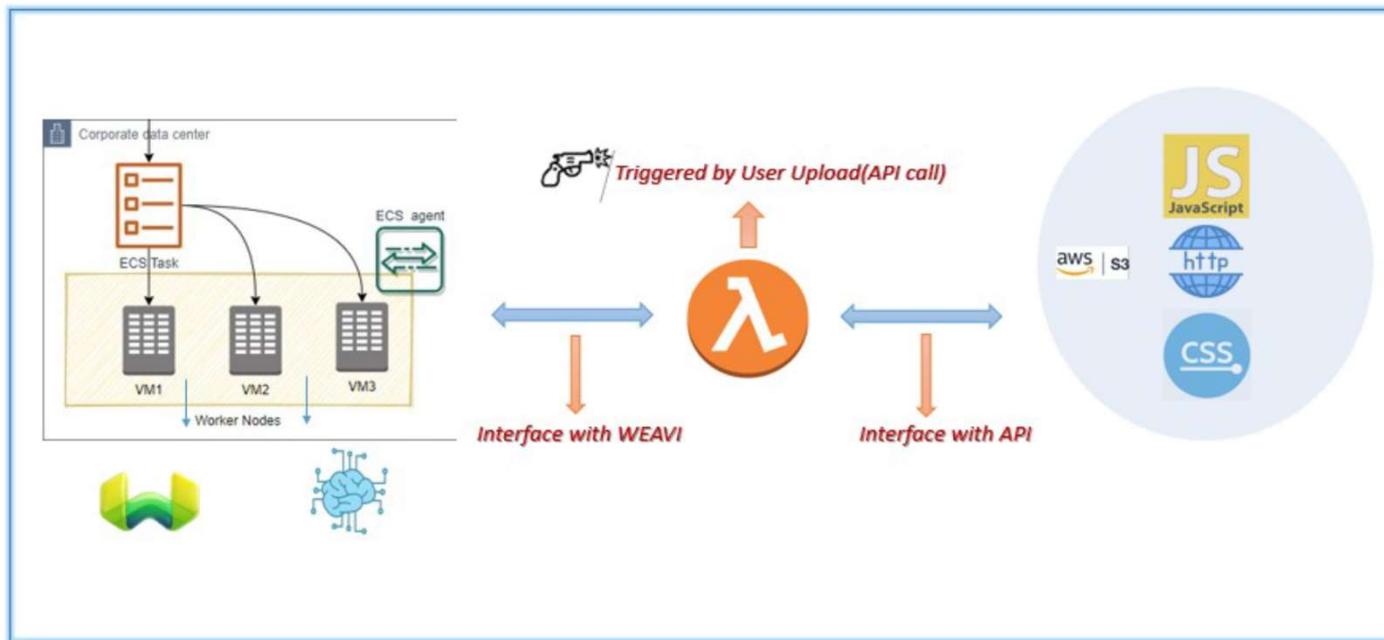


```python
import json
import weaviate
import numpy as np


def lambda_handler(event, context):
    print(event)
    try:
        image_b64 = event['body']
    except:
        return {
            'statusCode': 200,
            'body': json.dumps('ABORT

        }


    client = weaviate.Client(
        url="http://44.208.20.199:808
    )

    res = client.query.get(
        "Microrganism", ['species', 'i
        ).with_near_image(
            { "image": image_b64}, enc
        ).with_limit(15).with_addition


    return {
        'statusCode': 200,
        'body': json.dumps(res)
    }
```

# Why to scale?

## ECS *modus operandi*

This solution leverages Amazon ECS. Docker images are pushed to Amazon ECR, and ECS task definitions and services are created. Both service runs on the same EC2 instance, and load balancing is achieved through an Application Load Balancer (ALB).
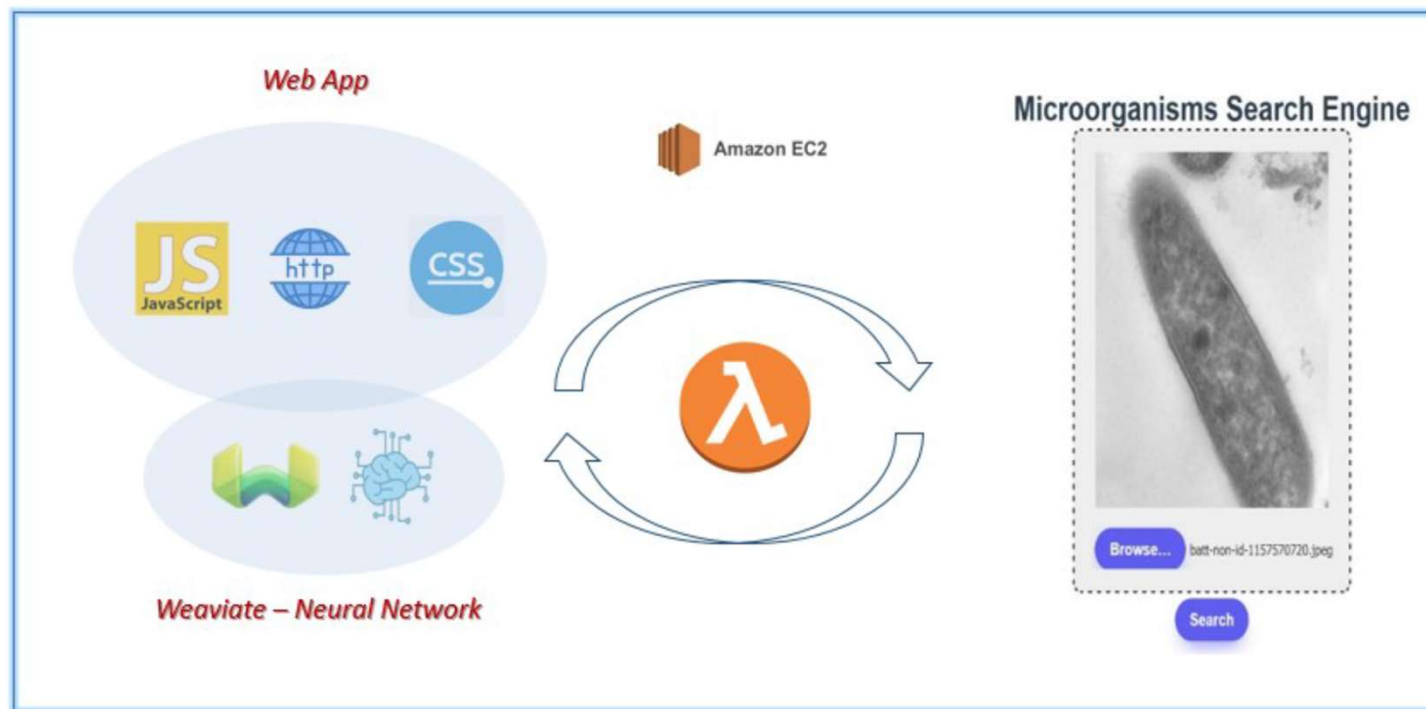


Scaling, monitoring, and adjustment of resources are performed for each ECS service based on specific requirements. This approach provides better performance optimization by allowing separate scaling policies for different components.

# Why to scale?

## EC2 *modus operandi*

Scaling the application enhances performance by distributing the workload across multiple EC2 instances, reducing response times, and improving query processing. It also enhances availability by using multiple instances and an elastic load balancer to redirect traffic in case of instance failure;



Due to _restrictions issues_ on the Learner Lab account we couldn't leverage ECS properly, so we resorted to **EC2 service**, running on multiple instances our web app tools.

# Scalable Architecture

To deploy the application and create a scalable architecture, the following steps are involved:

- Create Virtual Private Cloud (VPC) (create security group and configure inbound rules). Notice that we allow only the lambda function to access the ELB (and instances of ec2);

- Create Web Group & Define EC2 Rules;

- Define Launch template of EC2 Instance, also creating a custom AMI;

- Scale and Load Balance the Architecture:

- Configure the Auto Scaling Group with a minimum and maximum number of instances, such as a _minimum of 2 and a maximum of 6 EC2 instances_ for our purposes.

- The ALB makes intelligent routing decisions based on the content of the HTTP/HTTPS client requests, while also ensuring high availability and scalability.

# Scalable Architecture

# Web App results



**Microorganisms Search Engine**

- The client will upload here the query image of bacteria this event will act as a trigger for our lambda function.

- Remember that the lambda interface with the API with the Web app and with the EC2 instance to communicate With Weaviate.

Drop files here
or
Browse...  No file selected.



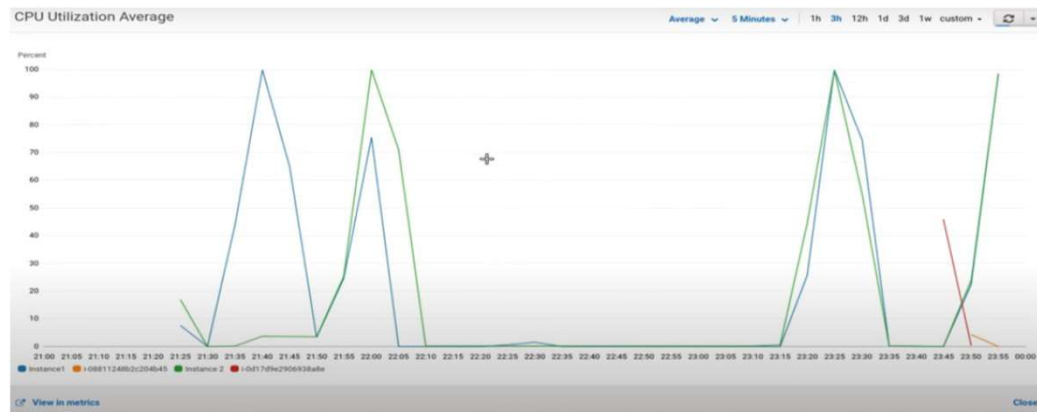**Microorganisms Search Engine**

New Search

**Paramecium**

Paramecium is a unicellular organism belonging to the group of protozoa. It has an elongated shape with movable cilia covering its entire surface for movement.

- **Size**: Typically ranging from 50 to 350 micrometers in length.
- **Diet**: Primarily feeds on bacteria, algae, and organic debris.
- **Temperature**: Lives in freshwater, and its temperature tolerance varies depending on the species.
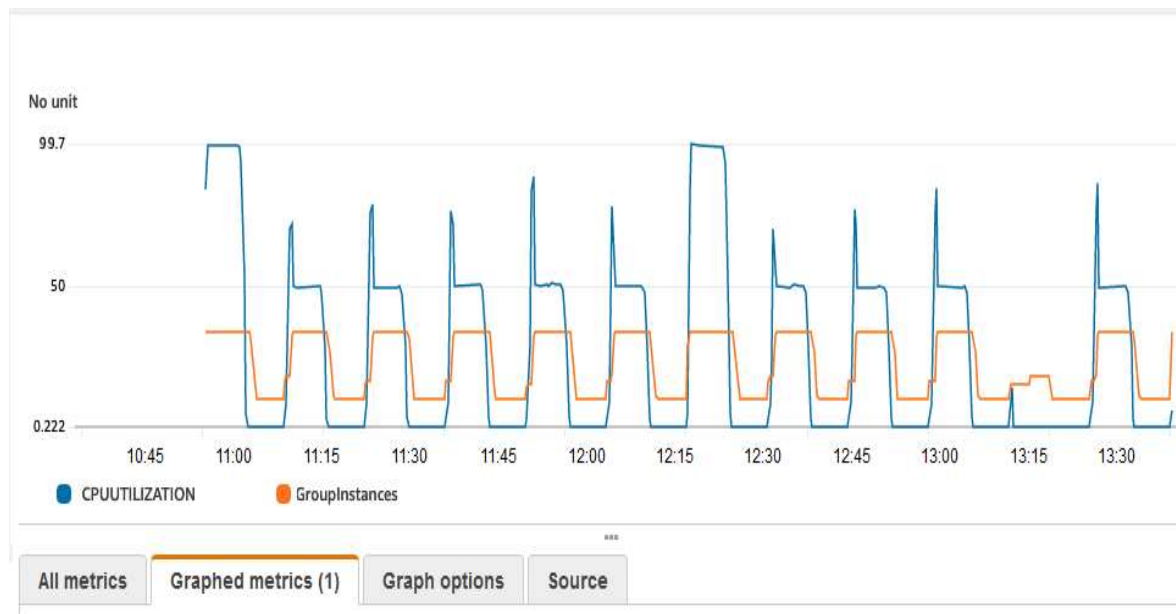- **Infectivity**: 1
- **Poisonus**: No

82%

- Each image will have a description of each type of bacteria.

- Notice that the web site is developed in a static way and it's just a bunch of files onto a S3 bucket.

- This ensures that the website is highly available since the distribution of the web site is menaged by AWS.

# Scalability results



- Peaks of the CPU utilization of the two baseline instance when stressed.



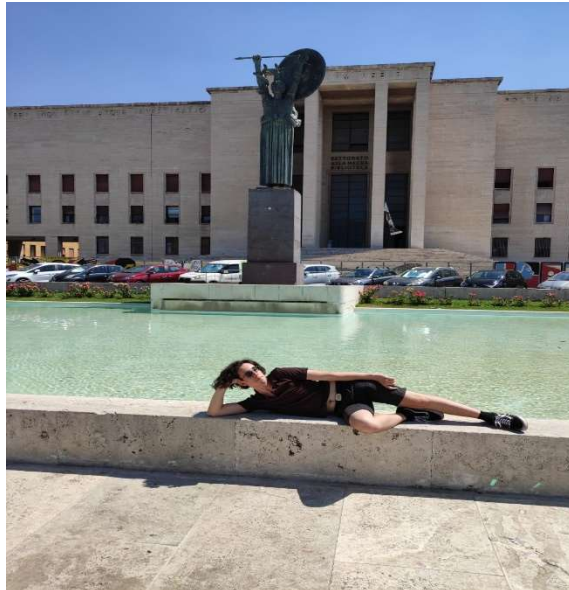- CPU usage versus number of instances running in Autoscaling Group.

# Thanks for your attention



ENRICO GRIMALDI



MATTEO MIGLIARINI



GIUSEPPE DI POCE