

# AWS Three-Tier Image Recognition Application

*Nishanth Murali*

## 1. Problem statement

In this project, we aim to develop a scalable image recognition service on the cloud which has all the software components loosely coupled. Every component of AWS used is elastic and hence, the project can be scaled up and down as per the number of requests coming in.

- Our cloud application performs image recognition on user-provided images using the provided deep learning model. It then returns the recognition result as output to the users.
- Each image input through the website is a request to our application and our application can handle multiple requests concurrently. It automatically scales out when the request demand increases, and automatically scales in when the demand drops.
- All the inputs (images) and outputs (recognition results) are stored in S3 for persistency.
- The application should handle all the requests as fast as possible, and it should not miss any requests. The recognition requests should all be correct.

## 2. Design and implementation

### 2.1 Architecture

- The user accesses the web tier to upload the images. The web tier always runs on an AWS instance. This web-tier instance then uploads the image to an S3 bucket ('input' bucket) and gets the 'presigned URL' from it. The URLs are sent to the app instance through the Request queue (FIFO SQS), where each app instance polls for a message from it and processes the same to get the classification results. Each of these results is then sent back to the web tier using a Response queue (FIFO SQS) and these results are then merged in order to write to the output file that stores all the outputs in the S3 bucket, in order to achieve persistence.
- The SQS is responsible for maintaining two communication channels between the web tier instance and application tier instances. The FIFO queues will result in each image URL being processed exactly once by the app instances so that there is no duplication of the classification results.
- The bash script configured in each of the app instances will be used to run the image classification code automatically when an instance starts. This will avoid human intervention while processing the images.

The high-level architecture of the cloud-application is shown in Fig. 1.

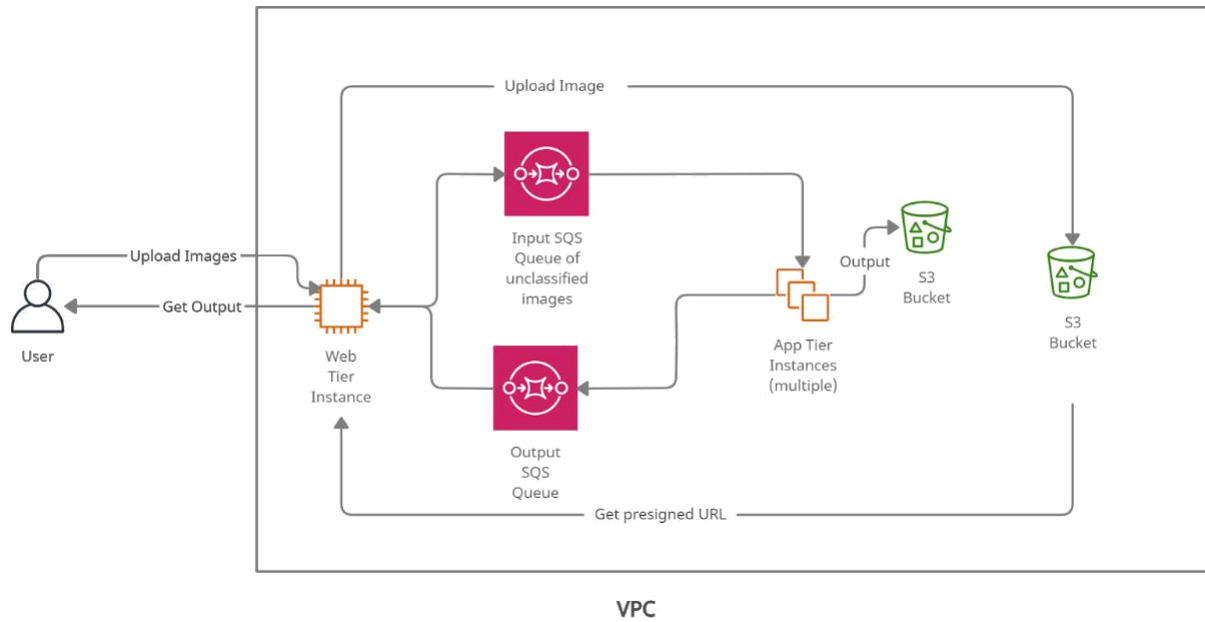


Fig. 1

## 2.2 Autoscaling

Autoscaling is one of the main reasons for hosting applications on the cloud. This project also utilizes the Autoscaling property of the cloud, so that the application resources are used efficiently. The web tier instance takes care of scaling up of the instances, based on the number of image-classification requests made. We will be using at most 19 app instances in this project and each instance processes the images to get the classification result. If the number of requests is more than 19 at a particular time, the request-messages are held in the Request queue until there is an app instance available to pick the request. The code embedded in each of the app instances will long poll the SQS request queue, in order to check if there are messages (image URLs) in the Request queue. If there are no messages left in the Request queue, the app instance will stop itself automatically, so that the resources are not wasted and thus the application is scaled down by the app instance themselves.

## 3. Testing and evaluation

Testing is an important part of building a cloud application, where you check if the application handles the requests concurrently and also if it auto-scales, so that minimum resources are used to maximize the efficiency of the final product.

In this project, we tested our application for stress and also efficient auto-scaling of the EC2 instances. The images below show that the application works as expected and the deployed product is efficient and fast.

Firstly, we uploaded 100 images through our web-tier instance (fig. 2), We then obtained the images' object URLs from the 'input' S3 bucket, where the images were stored. These URLs were then sent to the SQS 'Request' queue ('send\_queue.fifo' in fig. 3). Initially there was just one EC2 instance running,

which was the web-tier instance (fig. 4) and as per the length of the Request queue, the app instances got started (fig. 5). The output messages (image classification results) were transferred to the SQS Response queue ('receive\_queue.fifo' in fig. 6) and the messages in the Request queue were deleted. Once all the messages were received by the Response queue, the messages (results) were fed to the 'outputfile.txt' file, which was then sent to the 'output' S3 bucket, and the messages in the Response queue were also deleted (fig. 7). Meanwhile, when the EC2 app-instance had no messages to receive, they automatically shut themselves down (fig. 8). The classification results were then shown to the user on the UI (fig. 9). The 'input' bucket (fig. 10) and the 'output' bucket (fig. 11) had all the input image files and the output file respectively. The output file is shown in fig. 12.

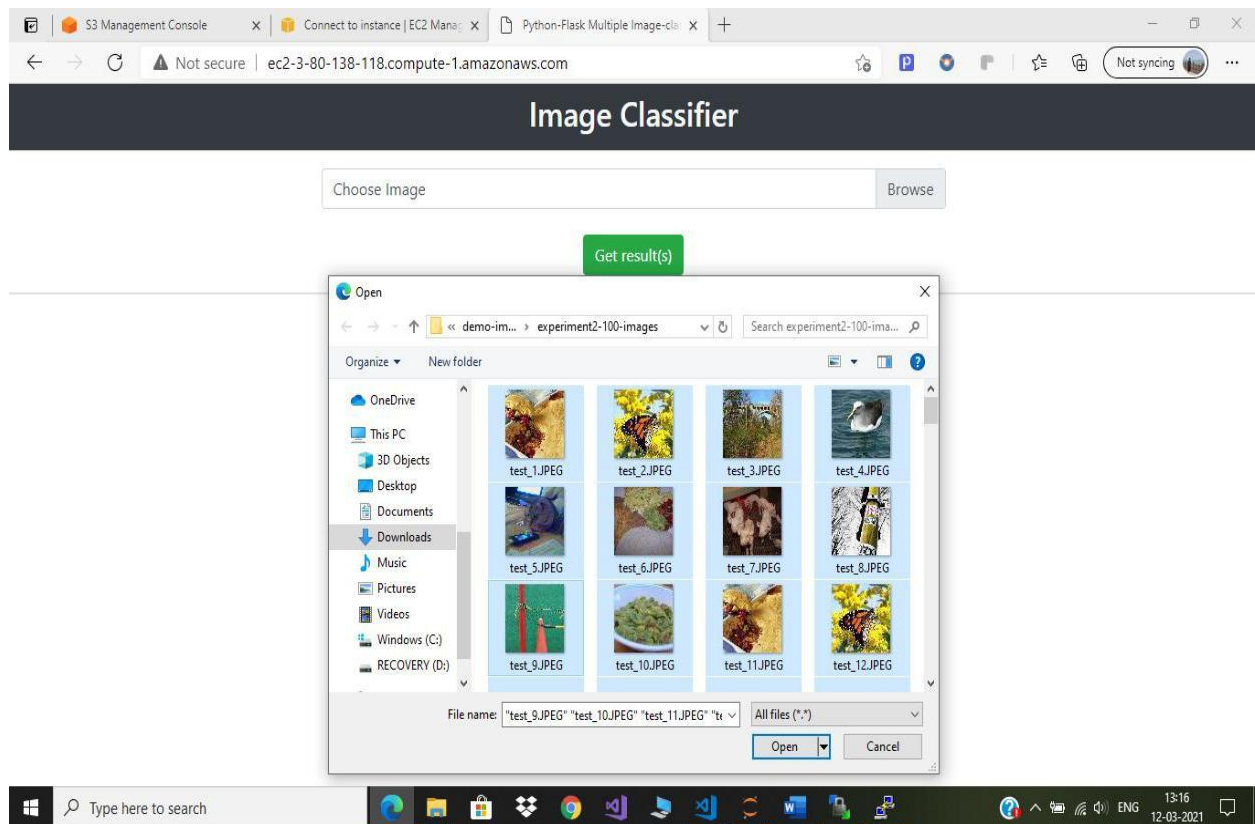


Fig. 2

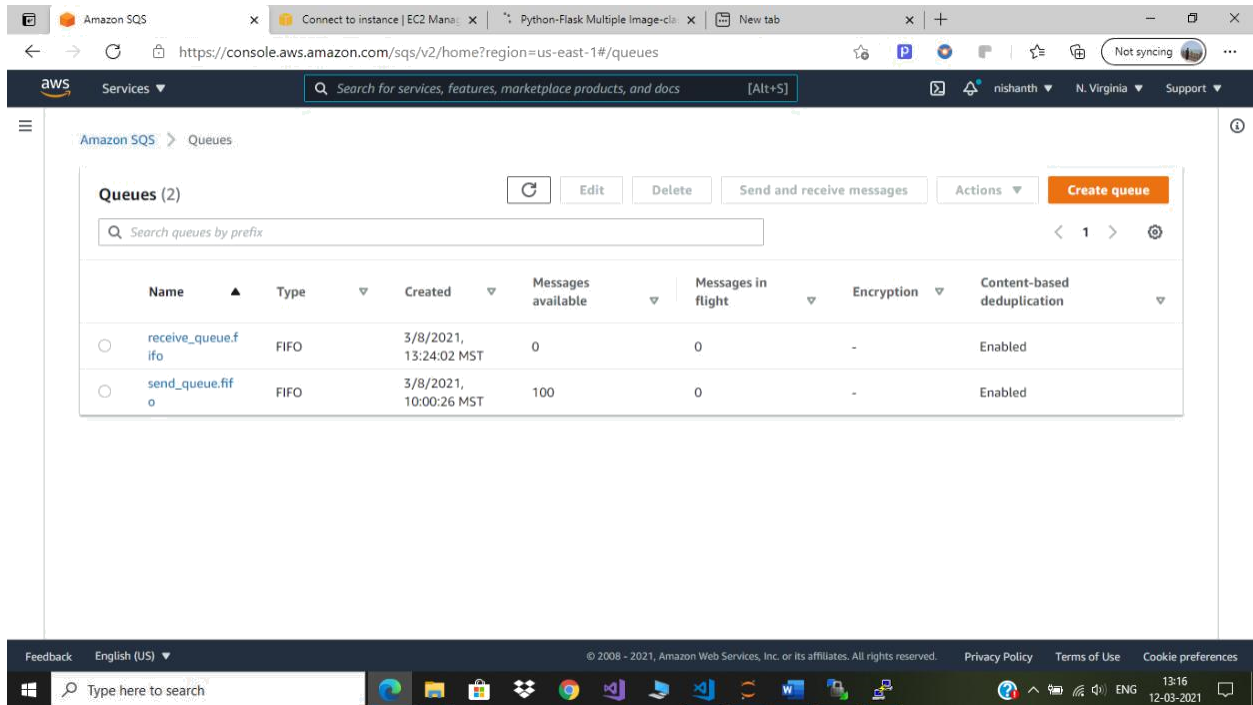


Fig. 3

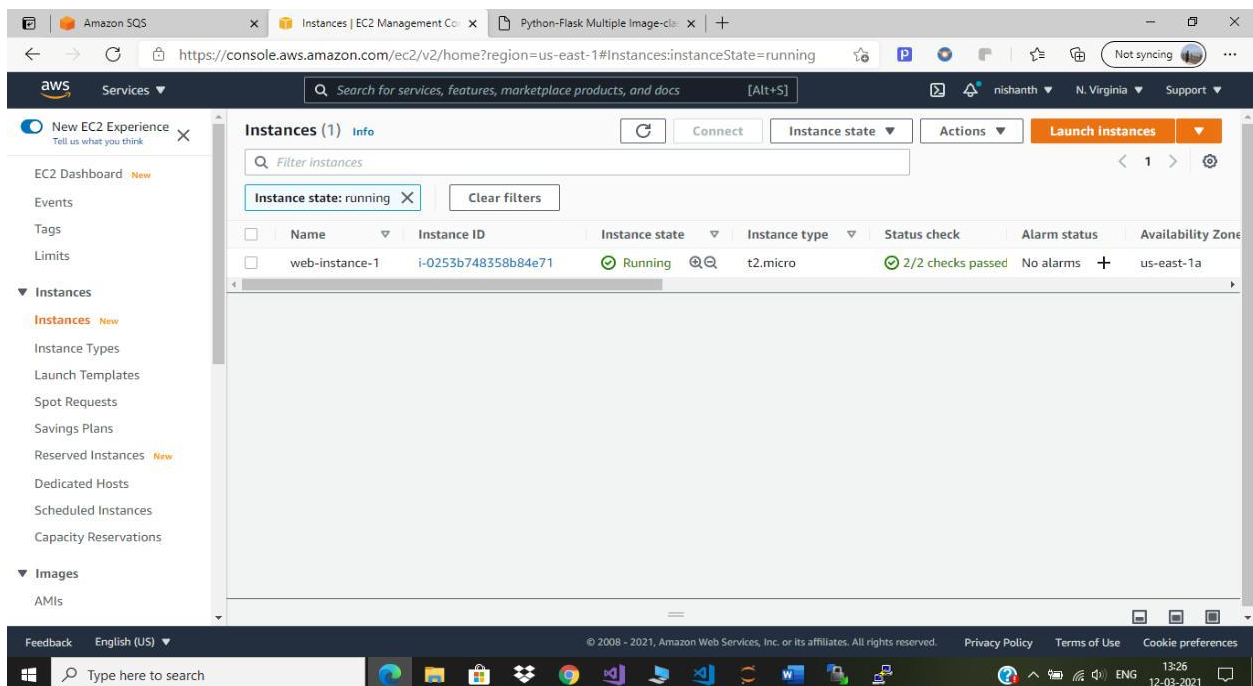


Fig. 4

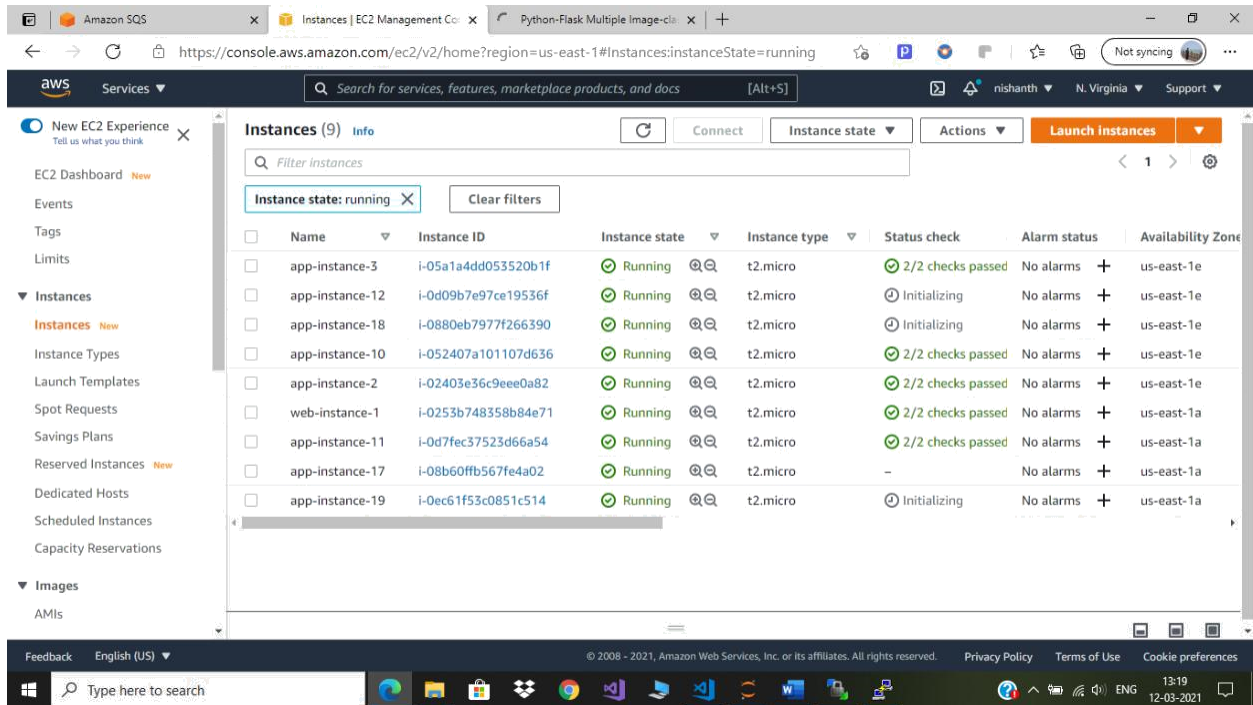


Fig. 5

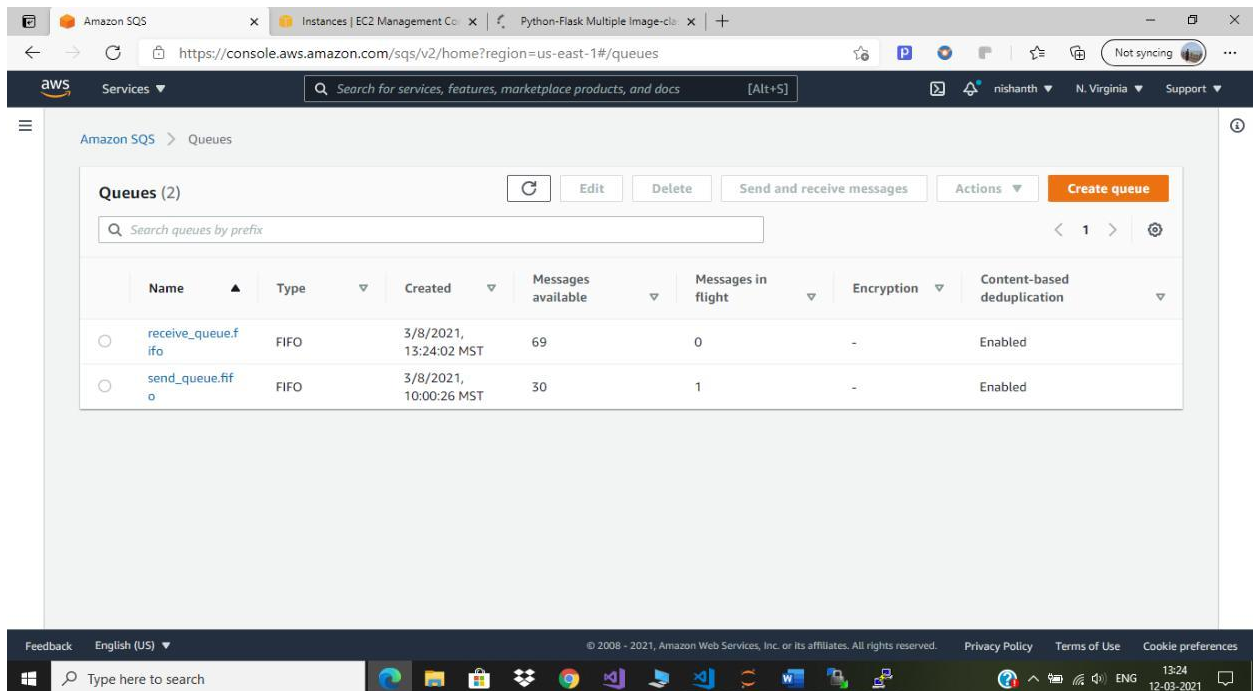


Fig. 6

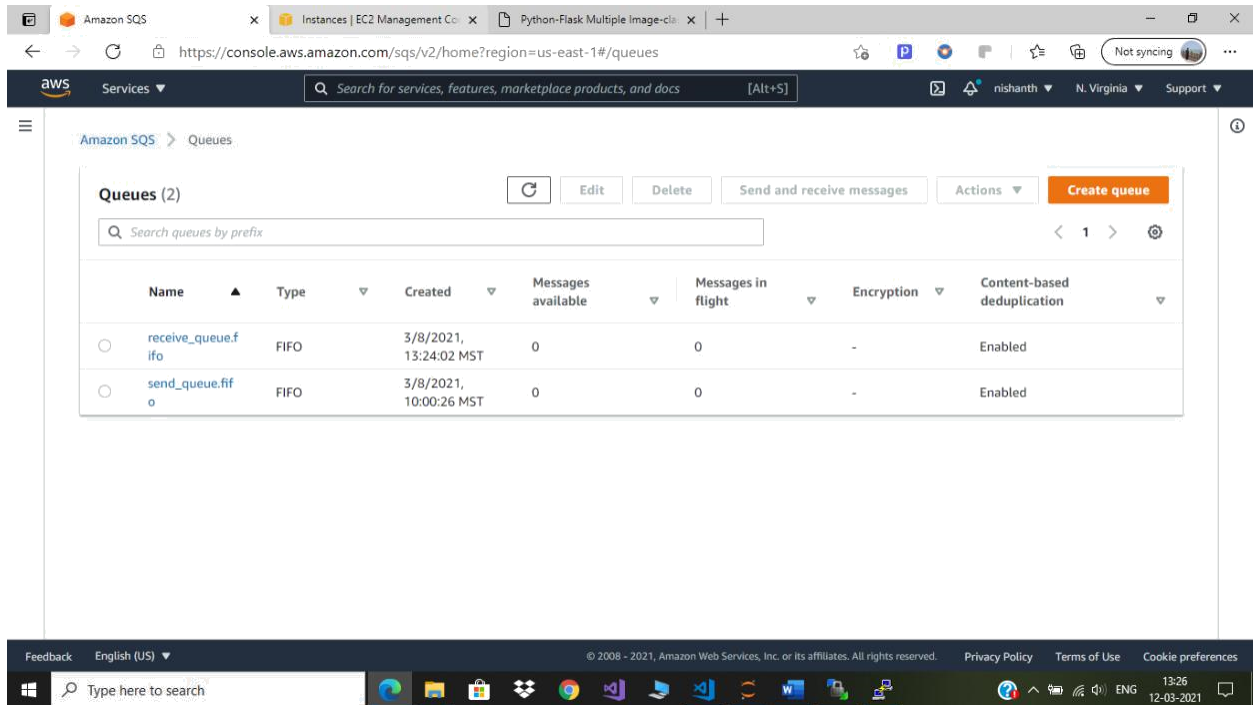


Fig. 7

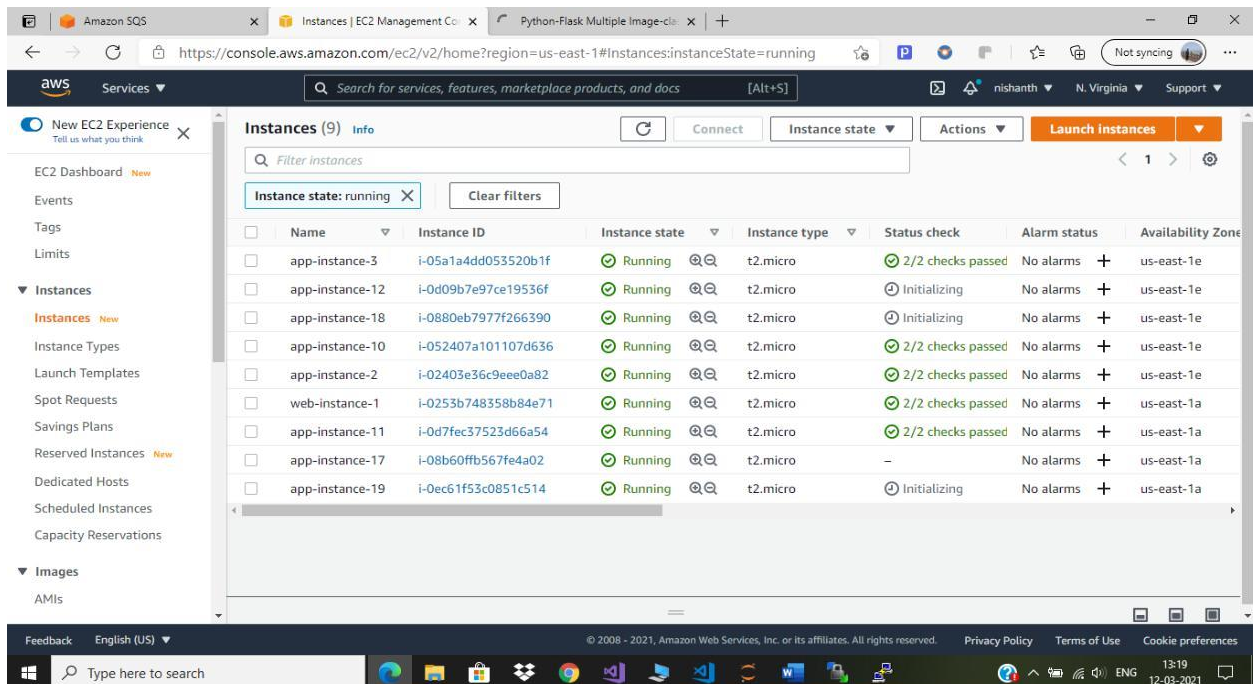


Fig. 8



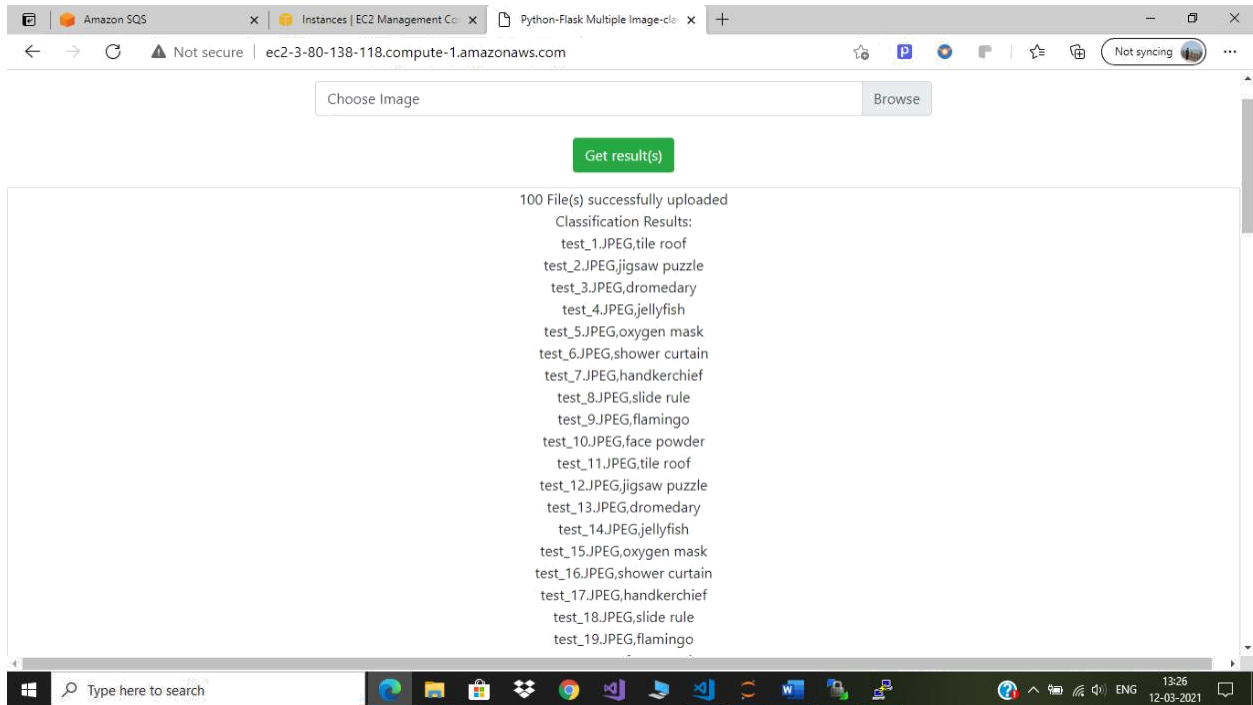


Fig. 9

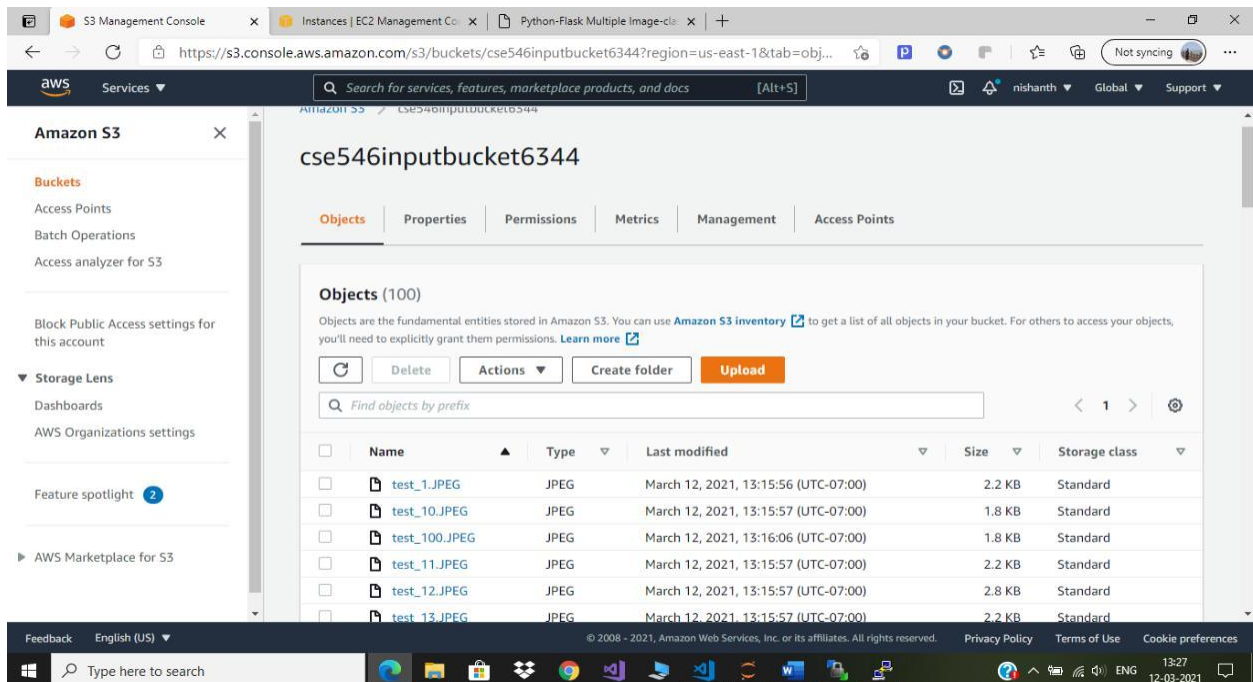


Fig. 10

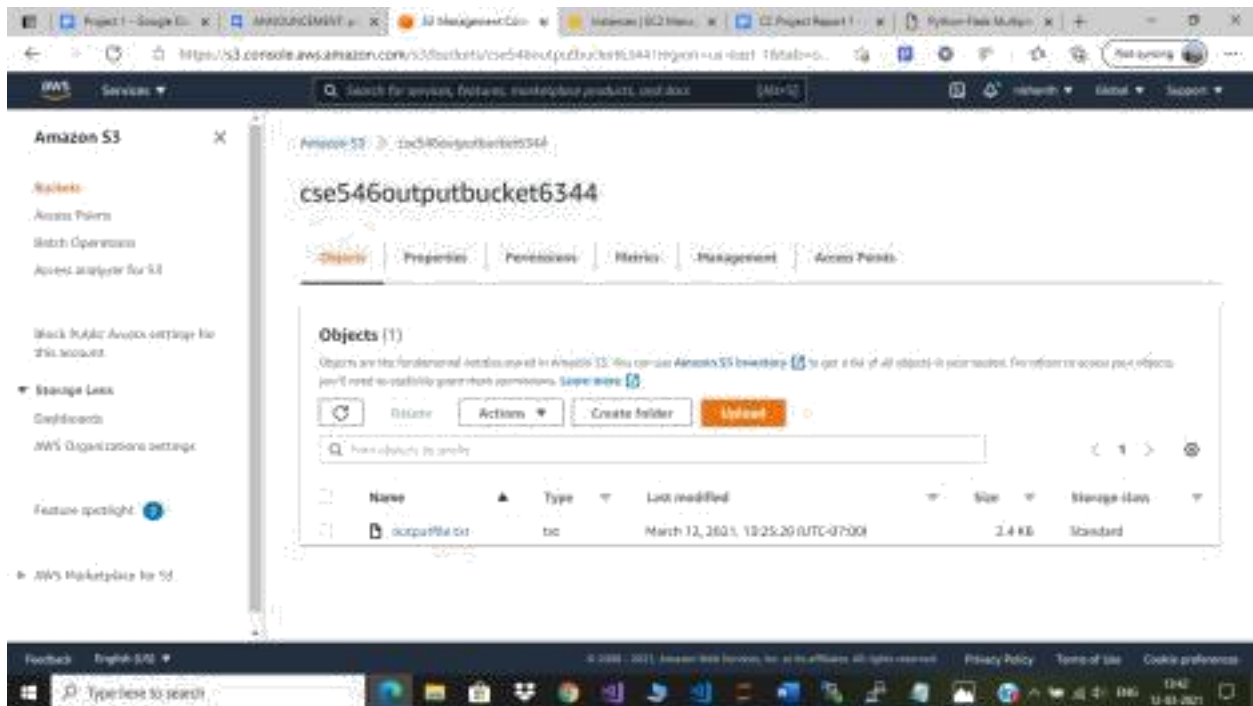


Fig. 11

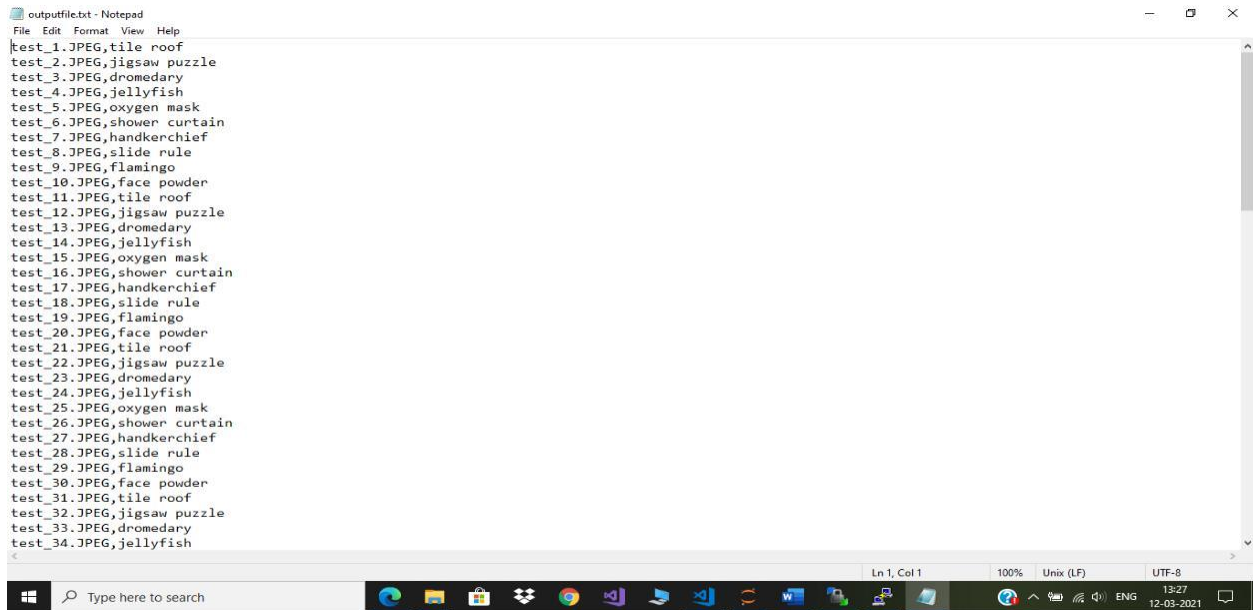


Fig. 12



The application's output, i.e., the image classification results, was evaluated with respect to the provided actual image tags and the classification results were totally accurate.

## 4. Code

The project uses three python programs that implement the different abstractions of the boto3 services, that have been used in the Flask app, which would run in the 'app.py' file. All the Three services are independent of each other and from the main module. They are loosely coupled such that they could be reused in other projects too with minimal or no changes. The three python modules using AWS SDKs in this project are;

1. SQSService: 'SQSService' contains methods that are required to perform basic operation on SQS Service like creating, deleting, sending, and receiving messages using 'sqs' resource of boto3 SDKs in Python.
2. S3Service: The boto3 SDKs also provide services to handle S3 buckets. 'S3Service' contains the methods that represent the Simple Storage Service (S3) at higher level client visibility. The methods help in uploading and downloading files to and from the S3 buckets.
3. EC2Service: 'EC2Service' contains methods that are required to perform basic operations on EC2, such as starting or terminating (stopping) an instance.

### 4.1 Web-Tier Module

Web-tier Module consists of a Flask code 'app.py' which acts as driving functionality of the web-tier. In this module, we would be mapping the URI to the backend API through Flask. The Frontend HTML and Javascript are also mapped using the rendering template of the flask. The main function takes input from the user through the Html form which is built using the bootstrap library. The files are then sent to the S3 and its Object URL is generated. The generated File Urls are pushed into the SQS queue (Request queue) to be later read by app-tier instances. After pushing the input image url to the SQS, the Web-Tier polls for messages from the output queue (Response queue), which are then stored in the S3 bucket and displayed to the user using Jinja template.

### 4.2 App-Tier Module

For each input image, an app-tier Instance processes it to get the classification result using the 'image\_classification.py' python code. It is the heart of the application as it takes care of delivering the output from the deep learning model to the user. The App-Tier Module First polls the message from the SQS, classifies the image using a deep learning model, and sends the output back to the web-tier through SQS. It also scales down the application, by shutting down the instance it runs on, after it sees no message in the Request queue.

### 4.3 Installation of the application

1. Web-tier

We need to type in the following commands in the Ubuntu shell, before we can run the Flask app that would run on an Apache server. Before that, we should create a folder called 'flaskproject' in /home/ubuntu directory and add all the application files in this folder. The commands to be executed later are given below;

1. `sudo apt-get update`
2. `sudo apt-get python3-pip`
3. `pip 3 install virtualenv`
4. `sudo apt-get install apache2 libapache2-mod-wsgi-py3`
5. `sudo ln -sT ~/flaskproject /var/www/html/flaskproject`
6. `sudo a2enmod wsgi`

Then, we need to edit a configuration file whose path is /etc/apache2/sites-enabled/000-default.conf, wherein we will add the following lines after the line containing "var/www/html";

```
WSGIDaemonProcess flaskproject threads=5
WSGIScriptAlias / /var/www/html/flaskproject/app.wsgi
```

```
<Directory flaskproject>
    WSGIProcessGroup flaskproject
    WSGIApplicationGroup %{GLOBAL}
    Order deny,allow
    Allow from all
</Directory>
```

Now, the environment for the Flask application is setup and we need to execute the below given three commands to run the application on EC2 instance's localhost, so that we can type in the public DNS of the EC2 instance in our local browser to access the web-tier, i.e the frontend of our application.

## 2. App-tier

We need to have a folder named 'classifier' in /home/ubuntu of our app-instance EC2 and add the files 'image\_classification.py' and 'imagenet-labels.json' in that folder. After this, we add the 'pystart.sh' shell script in /home/ubuntu and edit the crontab, in order to run the classification code automatically when the instance starts. The crontab can be edited by typing the commands;

1. `crontab -e`
2. `@reboot /home/ubuntu/pystart.sh`

With these configurations installed in the EC2 instance, we can build an AMI from this instance, so that we can create instances from this AMI, for autoscaling. All the app-instances created while auto-scaling the application will be using this AMI ID, so that all our configurations are up and running.

