# Operationalizing an AWS ML Project
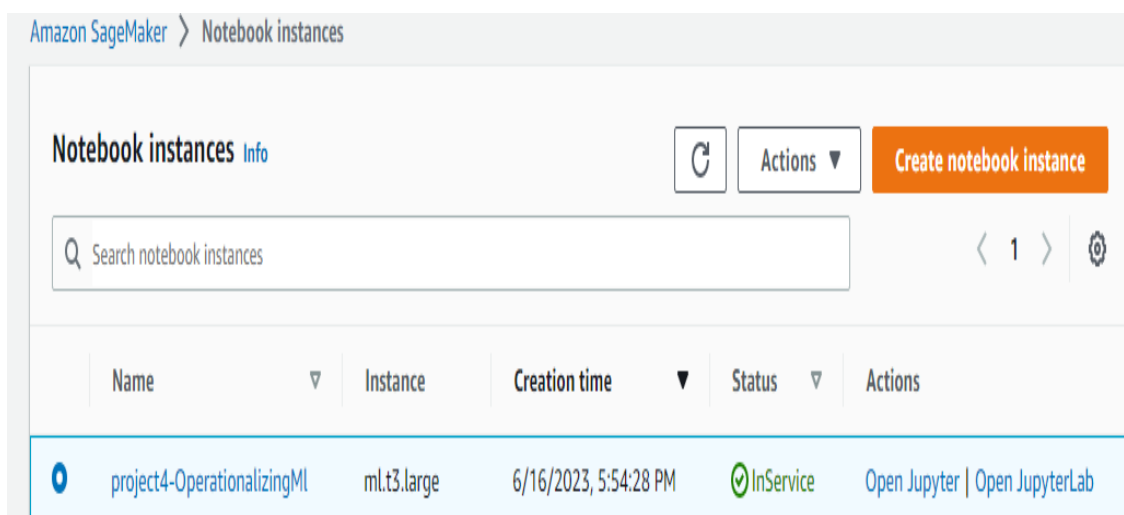
## Dog Image Classification

**The objective of the project is to finish the following steps:**

**1.** Train and deploy a model on Sagemaker, using the most appropriate instances. Set up multi-instance training in your Sagemaker notebook.

**2.** Adjust your Sagemaker notebooks to perform training and deployment on EC2.

**3.** Set up a Lambda function for your deployed model. Set up auto-scaling .ν for your deployed endpoint as well as concurrency for your Lambda function.

**4.** Ensure that the security on your ML pipeline is set up properl.

## Step 1: Training and deployment on Sagemaker:

**Created Sagemaker notebook instance** I have used ml.t3.large as this is sufficient and has a low cost to run my notebook.

## S3 bucket for the job (dog–images–mlop–project)



## Single instance training

## Multi-Instance training

| | Name | Creation time | Duration | Job status | Warm pool status | Time left |
|---|---|---|---|---|---|---|
| ○ | dog-pytorch-multi-instance-2023-06-17-15-11-24-196 | 6/17/2023, 5:11:24 PM | 21 minutes | ⊘ Completed | - | - |

**Training jobs** Info

🔍 Search training jobs

Actions ▼    **Create training job**

## Deployment

**Endpoints**

🔍 Search endpoints

Update endpoint    Actions ▼    **Create endpoint**

| | Name | ARN | Creation time | Status | Last updated |
|---|---|---|---|---|---|
| ○ | pytorch-inference-2023-06-17-15-55-03-786 | arn:aws:sagemaker:us-east-1:372206764755:endpoint/pytorch-inference-2023-06-17-15-55-03-786 | 6/17/2023, 5:55:04 PM | ⊘ InService | 6/17/2023, 5:57:24 PM |

## Step 2: EC2 Training

I have used t2.micro as it is a free tier





The adjusted code in ec2train1.py is very similar to the code in train_and_deploy–solution.ipynb. But there are few differences between the modules used – some modules can only be used in SageMaker. Much of the EC2 training code has also been adapted from the functions defined in the hpo.py starter script. Ec2train.py trains model with specific arguments while hpo.py takes argument for model by parsing through command line. The later code can train multiple model with different hyper parameters.

**Advantages and Disadvantages of EC2**

The advantages of EC2 Instances are less expensive than SageMaker instances, but the disadvantage of them is that they offer fewer managed.

**Step 3: Step 3: Lambda function setup**

After training and deploying your model, setting up a Lambda function is an important next step. Lambda functions enable your model and its inferences to be accessed by API's and other programs, so it's a crucial part of production deployment.
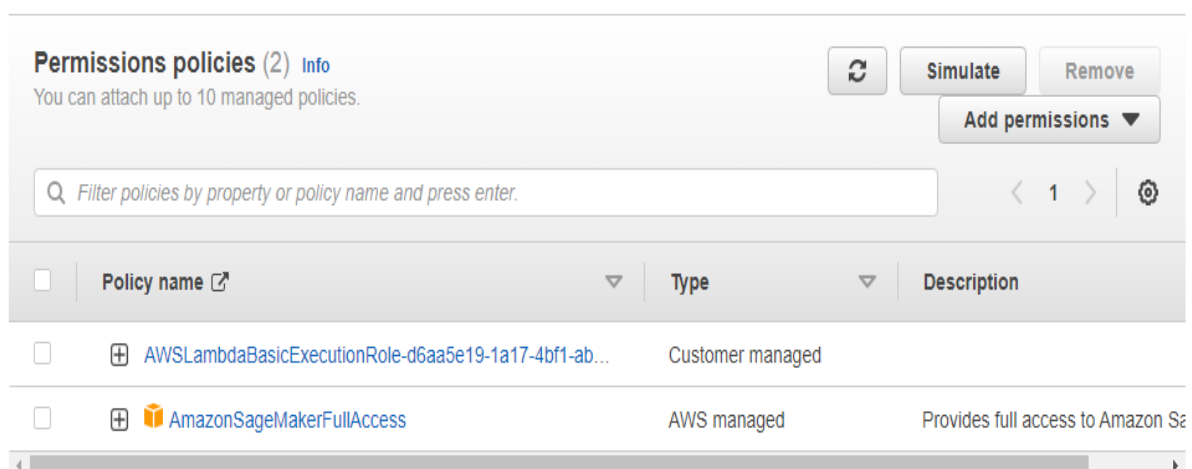
Thus, I deployed the lambda function with my endpoint name.

**Step 4: Lambda security setup and testing**

**Adding endpoints permission to lambda functions**

The lambda function will only be able to invoke endpoint if it has the proper security policies attached to it. Therefore I attached **AmazonSageMakerFullAccess** to the role of the lambda function



**Testing lambda function**

Now with the right permission we can create a new test to test our Lambda Function.



## Step 5: Lambda concurrency setup and endpoint auto-scaling

**Adding concurrency to Lambda Function**

By default, Lambda functions process one request at a time. If they receive multiple requests at the same time, they process one while the other ones wait. After the first request is processed, then the lambda function can move on to process the second request, then the third, and so on.

If you want to decrease the latency of your project in these high-traffic situations, you can implement concurrency, because concurrency so that the Lambda Function can respond to multiple requests at once.

**Reserved instances: 5/1000 Provisioned instances: 3/5**

**Auto-scaling**

Sagemaker endpoints require automatic scaling to respond to high traffic.

I enabled auto-scaling:

**Minimum instances: 1**

**Maximum instances: 3**

**Target value: 20 number of simultaneous requests which will trigger scaling**

**scale-in time: 30 s**

**scale-out time: 30 s**

**scale-in cool down:**

The scale-in period is the amount of time AWS will wait before deploying more instances for your endpoint.

- If you choose a high number, then AWS will wait a long time before deploying more instances. This helps you avoid incurring costs for momentary spikes in traffic.
- If you choose a low number, then AWS will deploy instances more quickly, but this responsiveness will be more costly.

**scale-out cool down**

- The scale-out cool down period is the amount of time AWS will wait before deleting extra deployed instances. If you choose a low number, then AWS will wait only a short time before deleting extra deployed instances. This helps you avoid incurring costs for momentary spikes in traffic. If you choose a high number, then AWS will keep extra instances deployed longer, but this extra capacity will be more costly.

**Therefore I choose an optimal value for both of them as well the low number of instances because of the low budget.**