# Operationalizing an AWS ML Project
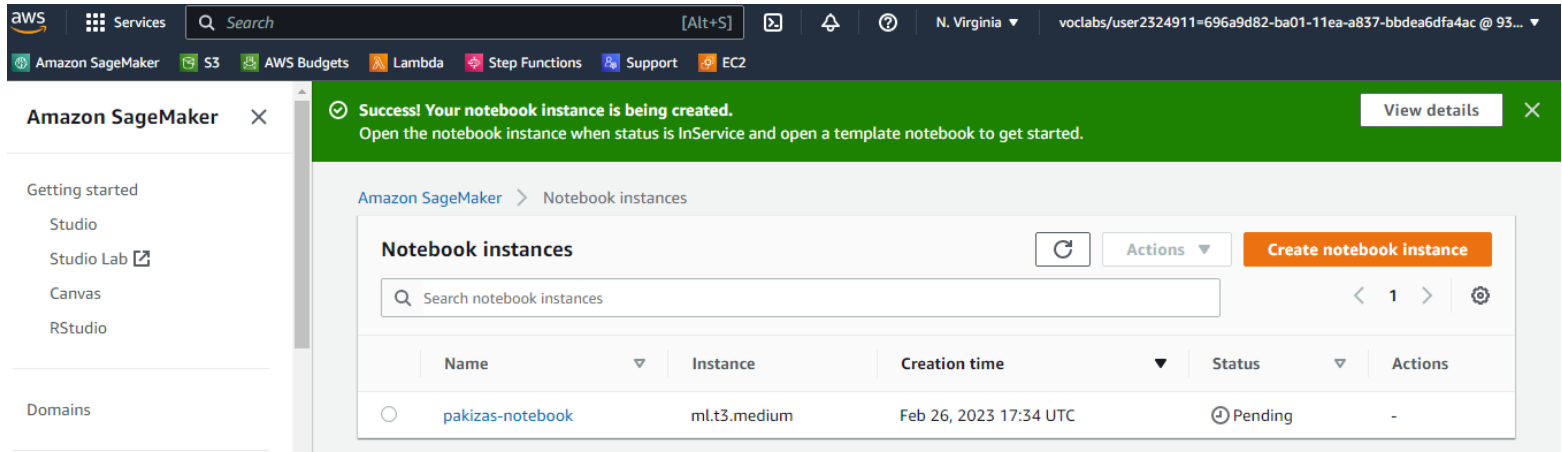
## Step 1: Training and Deployment on Sagemaker
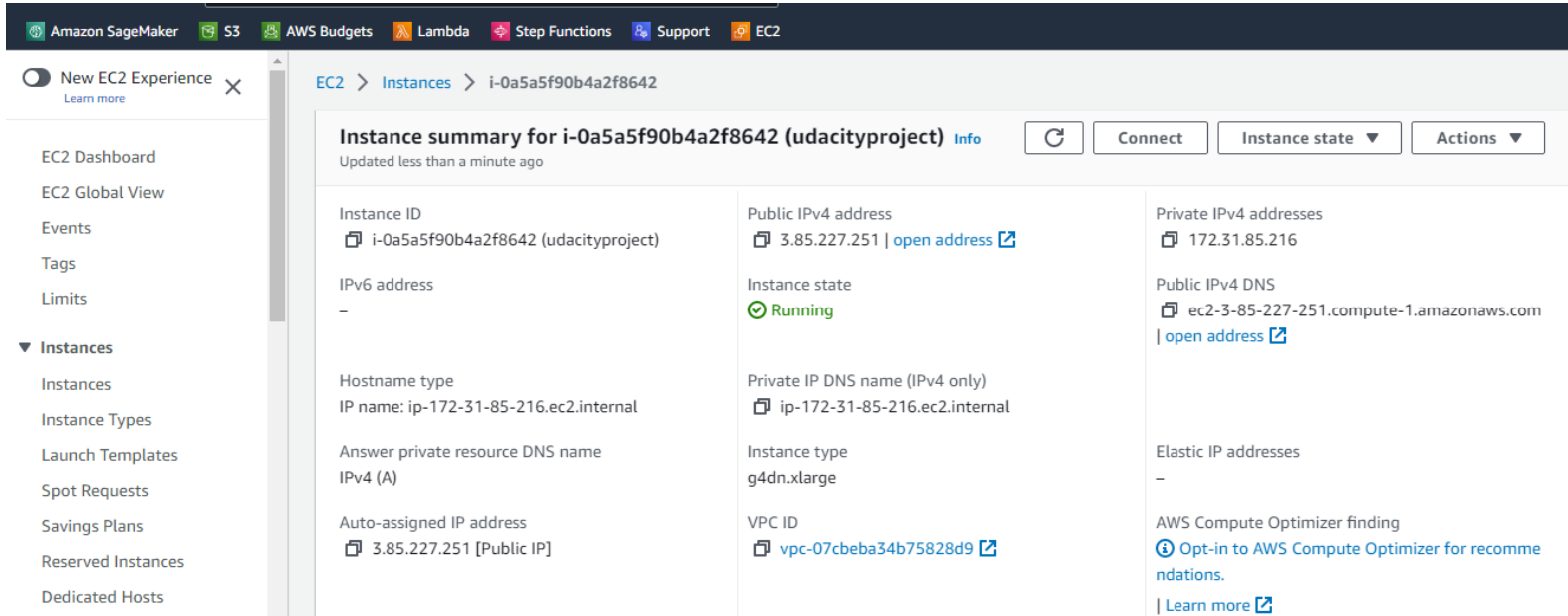
1) Created a new notebook instance



2) Used 'ml.t3.medium' instance type for notebook, but later changed the instance type to 'ml.m5.2xlarge' in order to train the models without getting resource limit error.
3) Fixed some bugs in the code
4) Trained using multiple instances and prepared results
5) Multi instance training ended after 15 minutes.

## Step 2: EC2 Setup

1) Launched an EC2 instance. Used 'ml.g4dn.xlarge' instance type in order to train the model without getting memory error. I chose "Deep Learning AMI GPU PyTorch 1.13.1 (Amazon Linux 2) 20230221" as a system type.



2) Connected to the instance and trained the model successfully.

Before uploading the data, I used "source activate pytorch" in order not to get any error related to python libraries. After that, I run the following commands in my EC2 terminal:

```
wget https://s3-us-west-1.amazonaws.com/udacity-aind/dog-project/dogImages.zip

unzip dogImages.zip

mkdir TrainedModels

cat > solution.py
```

After that, I copied and pasted the file content in "ec2train1.py". After pasting the code, I used Ctrl+D command to exit.

There are some similarities between the Amazon Sagemaker Studio and EC2 instances. However, the code for Amazon EC2 instance should be a self-contained Python script. By using the EC2 instance, we can train a model by using a local file system to access the data, as opposed to Amazon Sagemaker script that need methods for accessing Amazon S3 datasets.

## Step 3,4:  Lambda function setup and Security and testing

1) Created a new lambda function by using the "lambdafunction.py".
2) After creating a new lambda function, additional polices were added in order to run the lambda function successfully.

"Invoke Endpoint" works in the same way as a "Predict" in Amazon Sagemaker Studio.

**Output of the lambda function:**

"body": "[[-0.21229296922683716, -0.15309667587280273, -0.3061373829841614,
0.11011774092912674, -0.16774681210517883, -0.008066248148679733, -
0.23197175562381744, 0.15635429322719574, -0.12622188031673431, -
0.12885969877243042, 0.3695796728134155, -0.002062767744064331, -
0.055869728326797485, 0.14742261171340942, -0.11725631356239319, -
0.3161870837211609, -0.03053201735019684, -0.3684312105178833, -
0.32753172516822815, 0.15400850772857666, 0.07220902293920517,
0.09954917430877686, 0.013533521443605423, -0.06362387537956238, -
0.2291933298110962, -0.07303585857152939, -0.06227679178118706, -
0.39104604721069336, 0.029530785977840424, -0.18164679408073425, -
0.11869032680988312, -0.34357741475105286, 0.02373727224767208, -
0.042474523186683655, 0.2240193486213684, -0.005224950611591339, -
0.09473496675491333, -0.036600545048713684, 0.08635075390338898, -
0.303500235080719, 0.08778636157512665, -0.025020167231559753,
0.017018944025039673, 0.10680541396141052, -0.10075695067644119, -
0.15924464166164398, 0.23671609163284302, 0.0304688960313797, -
0.04243922978639603, -0.046819064766168594, 0.05256631597876549, -
0.03278496116399765, 0.006142046302556992, -0.01153213158249855,
0.027360720559954643, 0.08312027156352997, -0.09118440002202988, -
0.08264446258544922, 0.11408921331167221, -0.09587959945201874, -
0.048805855214595795, 0.058021046221256256, 0.13345670700073242, -
0.13694187998771667, -0.19396939873695374, -0.3879741430282593, -
0.32485464215278625, -0.10602971166372299, -0.2708476781845093, -
0.0530611053109169, 0.23138290643692017, -0.016494084149599075, -
0.04421060532331467, -0.2205985188484192, -0.16145184636116028,
0.0401601567864418, -0.167718306183815, -0.0998743399977684,
0.028368674218654633, 0.09563405811786652, -0.3112568259239197, -
0.11798204481601715, -0.14478018879890442, -0.12296410650014877, -
0.3944595754146576, 0.16801868379116058, -0.008254840970039368, -

# Step 5: Concurrency and auto-scaling



I set a minimum acceptable value of 2 concurrent lambda and endpoint executions for this project, the expected endpoint traffic is zero for this project as it will be removed.