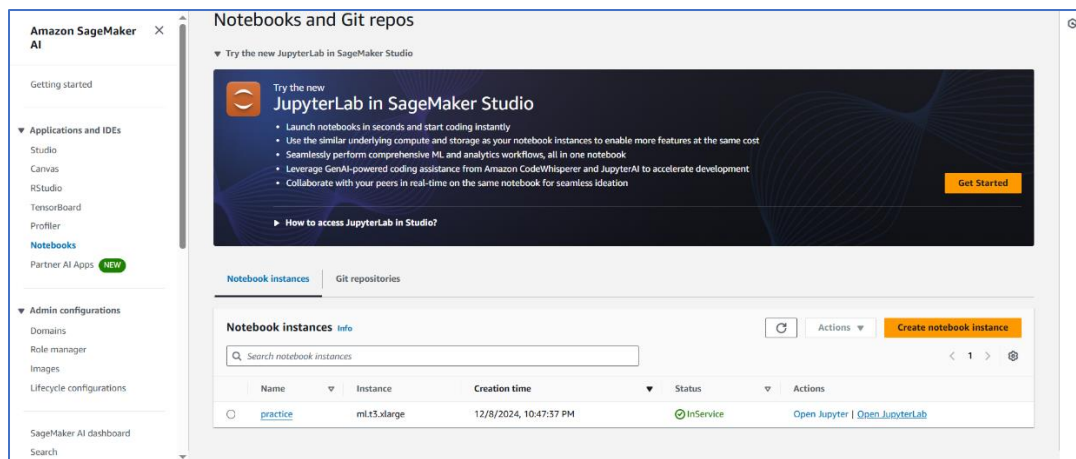# Operationalizing an AWS ML Project

# Dog Image Classification
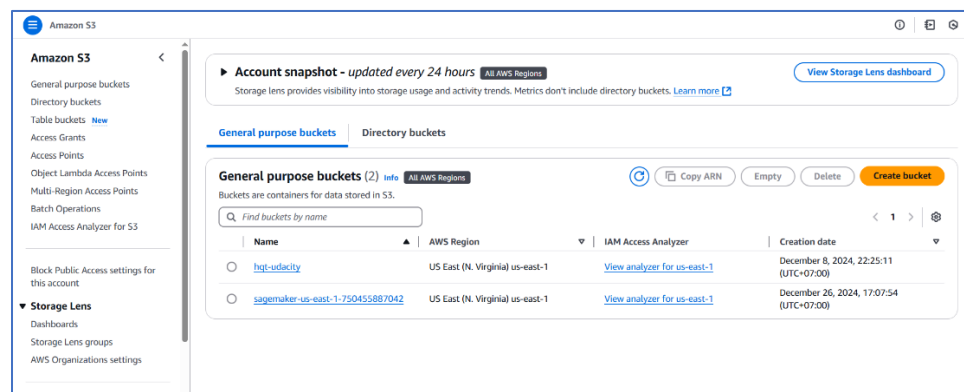
## 1. Training and deployment

- **Setup a Sagemaker instance**

    I chose the **ml.t3.xlarge** instance for running the notebook because it offers a good balance of cost and performance for development tasks. This instance provides sufficient CPU and memory resources to handle data preprocessing, model training, and basic hyperparameter tuning in SageMaker without incurring high costs, making it suitable for this stage of the project.
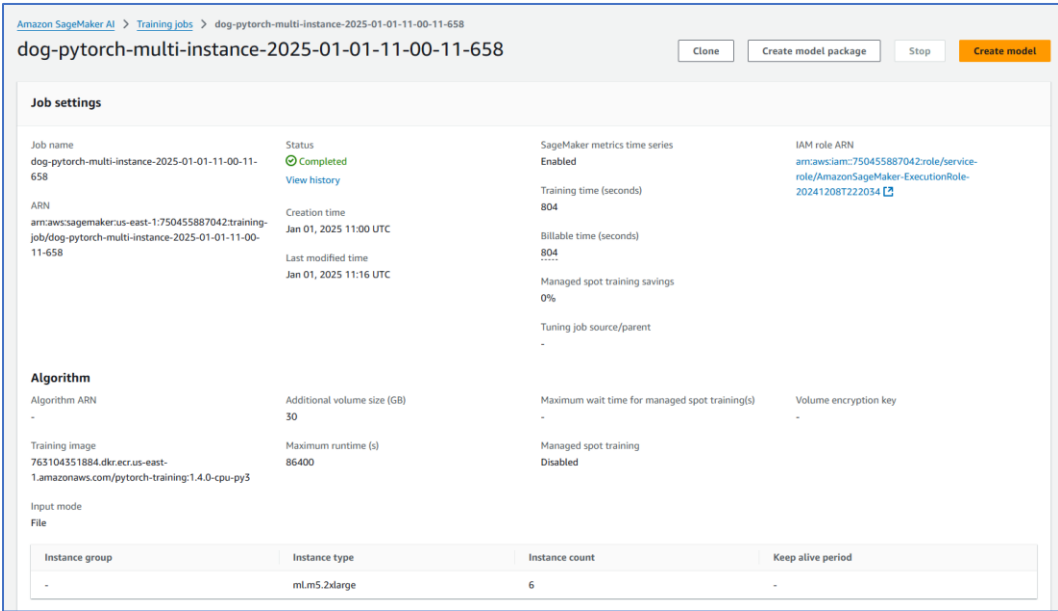


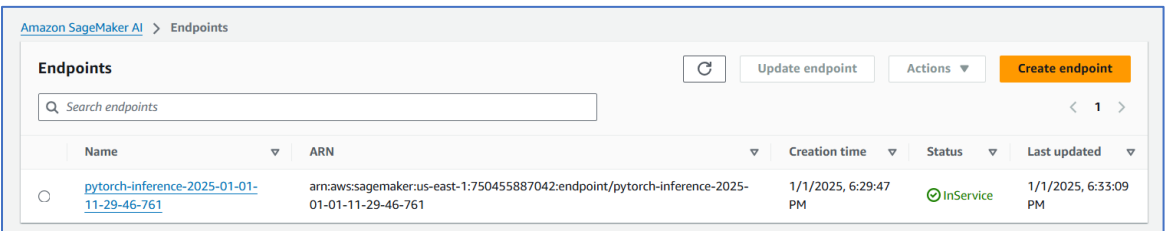- **Setup data storage on S3 (hqt-udacity bucket)**

- **Multi-instance training**

  The training job uses 6 **ml.m5.2xlarge** instances with SageMaker's PyTorch estimator for distributed training. This instance type was chosen for its cost-effectiveness and balanced CPU and memory resources, suitable for workloads not requiring GPUs. SageMaker Debugger and Profiler were enabled to monitor and optimize the training process. This setup ensures scalability and reduces training time efficiently.
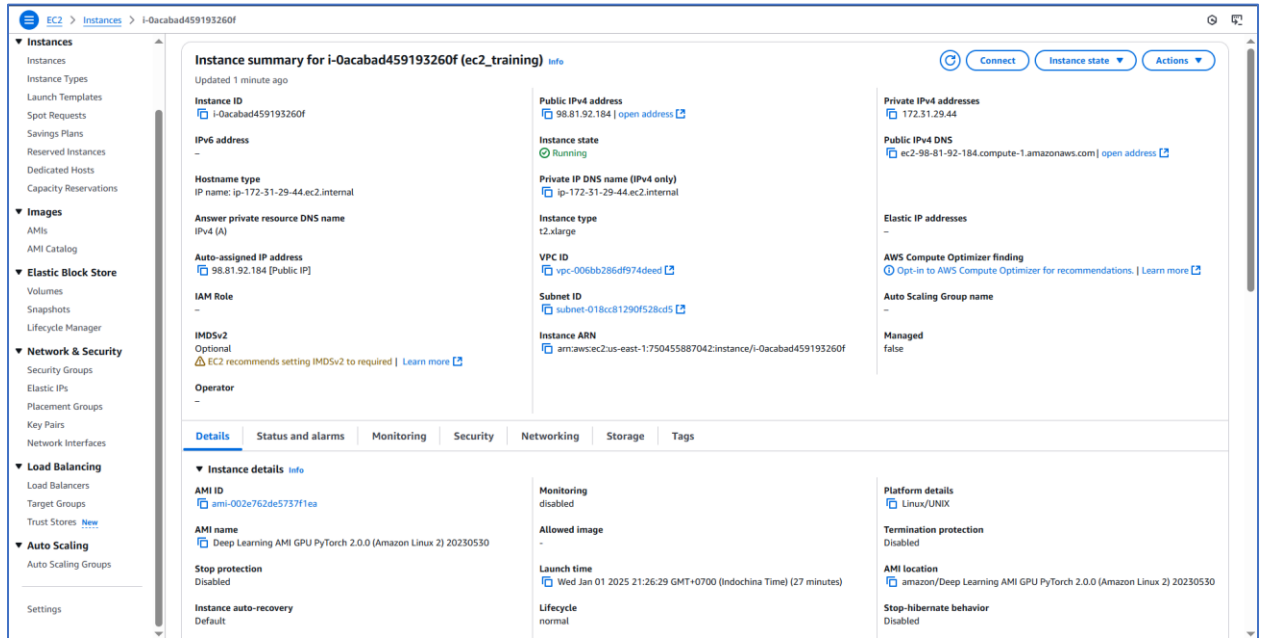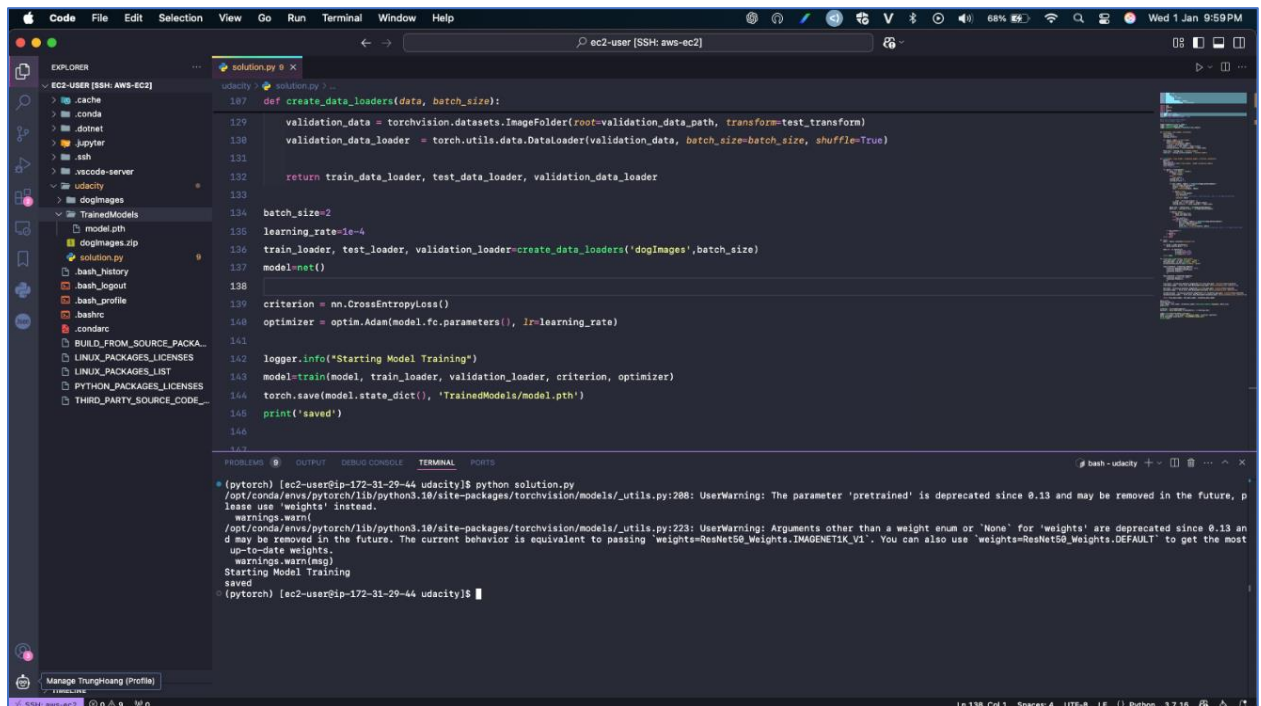


- **Deploy endpoint**



## 2. EC2 Training

- **Create and open an EC2 instance**

  I used **Deep Learning AMI GPU PyTorch 2.0.0 (Amazon Linux 2)** which has built-in Pytorch, Numpy libraries. I setup **t2.xlarge** EC2 instance because it provides a cost-effective solution with 4 vCPUs and 16 GB of memory, sufficient for running training

scripts and handling moderate computational tasks. It is ideal for development and testing phases where GPU acceleration is not critical, ensuring a balance between performance and cost efficiency.
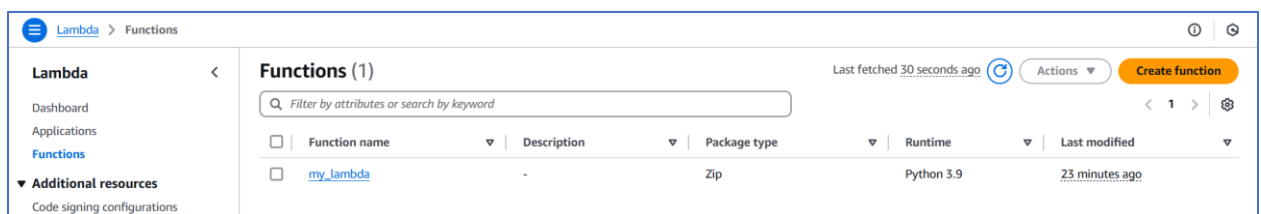


Training model on EC2, here is the terminal run train on vscode remote ssh to EC2.

- **Comparison of EC2 instances (ec2train1.py) with SageMaker Notebook Instances (train_and_deploy-solution.ipynb and hpo.py)**

  - ec2train1.py is suitable for simple, small-scale training setups, while the notebook and hpo.py are designed for production-grade pipelines on AWS.
  - SageMaker-based training supports distributed learning, making it ideal for large datasets or high-performance requirements.
  - The notebook setup integrates seamlessly with AWS tools, whereas ec2train1.py requires manual configuration.
  - EC2 instances are highly customizable and scalable, allowing users to tailor the setup to their specific needs, such as the choice of machine learning frameworks (e.g., PyTorch or TensorFlow), CPU and memory specifications, or GPU support. They are particularly suitable for high-performance computing tasks, enabling faster training for large-scale machine learning models. Additionally, EC2 instances can be optimized to match specific workloads, providing flexibility and control over computational resources.
  - SageMaker notebook instances provide a pre-configured environment with popular machine learning frameworks and libraries, making setup quick and user-friendly. These instances are well-integrated with other AWS services, such as SageMaker for training and deploying models, and offer collaboration tools for seamless teamwork. This integration with the AWS ecosystem simplifies the workflow, making it ideal for machine learning engineering and operations tasks.
  - In summary, EC2 instances are better suited for customized, resource-intensive training needs, while SageMaker notebook instances are designed for rapid setup and efficient integration with AWS machine learning tools.

## 3. Lambda function setup

- Lambda function invokes the student's deployed endpoint

## 4. Security and testing

- **Attach a security policy to the role associated with the Lambda endpoint**



- **Test the Lambda function.**
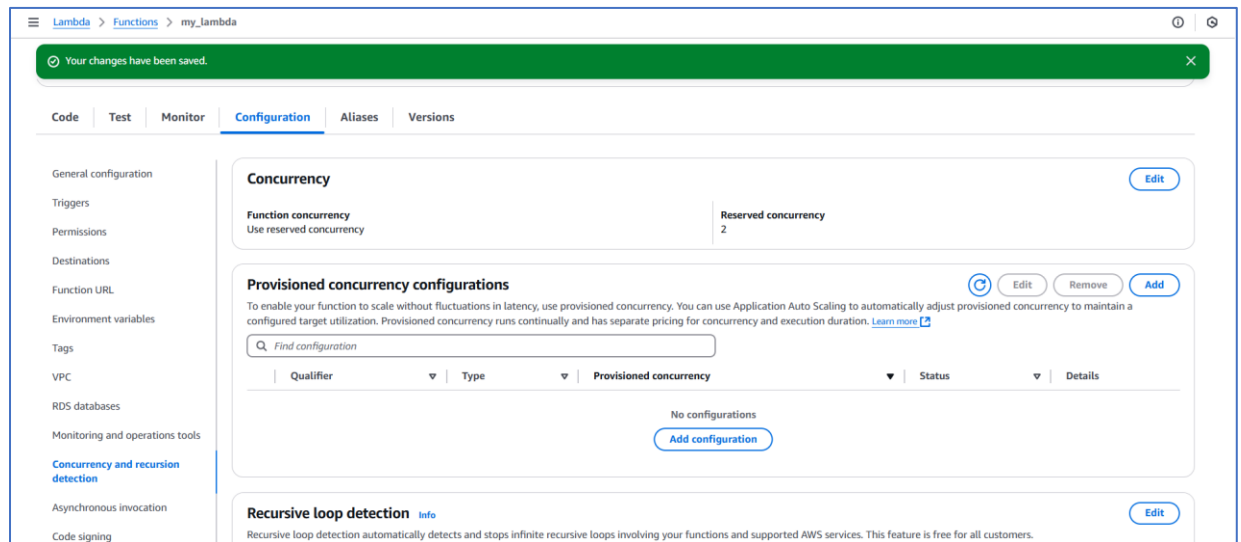
- **Identified Potential Vulnerabilities**

  - In the setup for the Lambda function, the IAM role assigned includes the **AmazonSageMakerFullAccess** policy to enable communication with the SageMaker inference endpoint. While this allows the function to execute the required tasks, it also introduces potential vulnerabilities that should be addressed for improved security.
  - The role attached to the Lambda function has AmazonSageMakerFullAccess, which provides broad permissions, including the ability to create, update, and delete SageMaker resources. If the Lambda function is compromised, an attacker could misuse these permissions to manipulate SageMaker resources, leading to data leaks, endpoint misuse, or increased costs.
  - The assigned role includes permissions for broader actions across all SageMaker resources. This violates the principle of least privilege, which states that roles should only have permissions necessary for their intended purpose.

## 5. Concurrency and auto-scaling

- **Lambda concurrency**

  In this project, the Lambda function was configured with Reserved Concurrency set to 2. By default, a Lambda function can handle only one request at a time. Configuring concurrency allows the function to handle multiple requests

simultaneously by reserving the specified number of execution environments. The chosen concurrency level of 2 is appropriate for a low-demand service or a development/test environment. For production or high-traffic scenarios, this setting should be revisited and potentially increased or replaced with an auto-scaling approach to ensure efficient resource utilization and cost management.



Considerations:

- o Traffic: This configuration enables the Lambda function to handle up to two simultaneous requests. While suitable for low-traffic scenarios, it may not be sufficient for high-demand services.
- o Cost: Provisioned concurrency incurs additional costs, as it pre-initializes environments to reduce cold start latency. By keeping the reserved concurrency value low (2), the additional cost is minimized.
- o Efficiency: Setting concurrency to 2 balances reduced latency for simultaneous requests while avoiding unnecessary expenses. However, for scalable applications with fluctuating traffic, higher concurrency or dynamic scaling might be required.

- **Auto-Scaling Sagemaker Endpoint**

To ensure the endpoint can efficiently handle fluctuating traffic, auto-scaling was configured for the SageMaker endpoint. This allows the system to dynamically adjust resources based on the number of incoming requests.

Considerations:

- o Traffic: This setup allows the endpoint to efficiently handle up to 15 simultaneous requests per instance, with a maximum of 3 instances ensuring scalability for moderate traffic.
- o Cost: Keeping the minimum instance count at 1 reduces idle resource costs, while the maximum of 3 ensures that scaling is limited to control expenses during peak demand.
- o Efficiency: By setting reasonable delays for scaling actions, the configuration minimizes rapid instance changes, improving stability and ensuring resources are used effectively.