

# **MLOps Laboratory Assignment 3**

## **Experiment Tracking and Versioning with MLFlow**

Machine Learning Operations

December 13, 2025

# Contents

1	Introduction .....	3
2	Repository and Deployment Links .....	3
2.1	GitHub Repositories .....	3
2.2	HuggingFace Spaces .....	3
3	Testing Strategy .....	3
3.1	Overview .....	3
3.2	Test Architecture .....	3
3.3	Test Categories .....	3
3.3.1	Unit Tests (test_logic.py) .....	3
3.3.2	CLI Integration Tests (test_cli.py) .....	4
3.3.3	API Integration Tests (test_api.py) .....	4
3.3.4	Artifact Existence Tests (test_artifacts.py) .....	5
3.4	Testing Fixtures and Utilities .....	5
4	Experiments Conducted .....	6
4.1	Experimental Setup .....	6
4.1.1	Dataset .....	6
4.1.2	Base Model Architecture .....	6
4.2	Hyperparameter Configurations .....	6
4.3	Logged Artifacts .....	7
4.3.1	Parameters Logged .....	7
4.3.2	Metrics Logged .....	7
4.3.3	Artifacts Logged .....	7
5	Results Analysis .....	8
5.1	MLFlow UI Analysis .....	8
5.1.1	Model Selection .....	8
5.2	Performance Analysis .....	9
5.2.1	Training Curves .....	9
5.2.2	Final Model Performance .....	9
5.3	Model Serialization .....	9
6	Implementation Details .....	9
6.1	Model Selection Script .....	9
6.2	Inference Script .....	10
7	Services Working .....	10
7.1	DockerHub .....	10
7.2	Render .....	10
7.3	HuggingFace .....	11
8	Challenges and Solutions .....	12
8.1	Docker .....	12

# 1 Introduction

This report presents the work completed for Laboratory Assignment 3 of the MLOps course, focusing on experiment tracking and model versioning using MLFlow. This assignment builds upon the previous laboratories (Lab 1 and Lab 2) by replacing the random prediction model with a deep learning classifier trained using transfer learning on the Oxford-IIIT Pet dataset.

The main objectives accomplished in this laboratory are:

- Implementation of transfer learning using lightweight deep learning models
- Experiment tracking and model versioning with MLFlow
- Model serialization in ONNX format for production deployment
- Integration of the serialized model into the inference API

## 2 Repository and Deployment Links

### 2.1 GitHub Repositories

- **Lab 1:** [github.com/MariaGoico/MLOps-Lab1](https://github.com/MariaGoico/MLOps-Lab1)
- **Lab 2:** [github.com/MariaGoico/MLOps-Lab2](https://github.com/MariaGoico/MLOps-Lab2)
- **Lab 3:** [github.com/MariaGoico/MLOps-Lab3](https://github.com/MariaGoico/MLOps-Lab3)

### 2.2 HuggingFace Spaces

- **Lab 2 Deployment:** [huggingface.co/spaces/spaces/MGoico/Lab-2-Mlops](https://huggingface.co/spaces/spaces/MGoico/Lab-2-Mlops)
- **Lab 3 Deployment:** [huggingface.co/spaces/MGoico/MLOps-Lab3](https://huggingface.co/spaces/MGoico/MLOps-Lab3)

## 3 Testing Strategy

### 3.1 Overview

The testing strategy for this project ensures the reliability and correctness of the MLOps pipeline, from data preprocessing to model inference.

### 3.2 Test Architecture

The test suite is organized into four main test modules:

1. **test\_logic.py** - Unit tests for core business logic and preprocessing functions
2. **test\_cli.py** - Integration tests for the command-line interface
3. **test\_api.py** - Integration tests for the FastAPI endpoints
4. **test\_artifacts.py** - Pre-deployment validation tests for model artifacts

### 3.3 Test Categories

#### 3.3.1 Unit Tests (test\_logic.py)

The unit tests validate the core functionality of the image processing and prediction pipeline:

##### Prediction Functions:

- **predict\_simple():** Tests both ONNX classifier usage and fallback to random prediction
  - Validates that predictions always return valid class labels from `class_labels.json`
  - Tests error handling when classifier is unavailable
  - Verifies graceful fallback on runtime errors
- **predict():** Tests the confidence-aware prediction function
  - Validates tuple return format (class, confidence)
  - Tests fallback behavior returning `None` confidence
  - Verifies ONNX classifier integration with mocked predictions

### Image Preprocessing Functions:

- **resize():** Tests image resizing with various parameter combinations
  - Specific width and height dimensions
  - Random dimensions when parameters not provided
  - Partial specification (width-only or height-only)
  - Error handling for invalid dimensions (negative, zero)
- **to\_grayscale():** Validates conversion to grayscale mode
- **normalize():** Tests autocontrast application and pixel value scaling to [0,1]
- **random\_rotate():** Verifies rotation within expected angle range (-20° to +20°)
- **random\_flip():** Tests conditional horizontal flipping based on probability
- **blur():** Validates Gaussian blur filter application
- **preprocess():** Tests complete preprocessing pipeline execution order

### Utility Functions:

- **ensure\_output\_dir():** Tests directory creation and idempotency

### Class Labels Validation:

- Validates 37 expected classes from Oxford-IIIT Pet dataset
- Confirms presence of known breed names

### 3.3.2 CLI Integration Tests (test\_cli.py)

Tests for the Click-based command-line interface covering two main command groups:

#### Classify Group:

- **predict command:** Tests image classification via CLI
  - Success cases with valid images
  - Error handling for non-existent files
  - Mocked prediction validation

#### Preprocess Group:

- **resize command:** Tests with explicit dimensions, random dimensions, and default output paths
- **grayscale command:** Validates grayscale conversion
- **rotate command:** Tests random rotation functionality
- **flip command:** Tests random horizontal flip
- **blur command:** Tests Gaussian blur application
- **pipeline command:** Tests complete preprocessing pipeline

#### Edge Cases:

- Partial dimension specification (width-only, height-only)
- Default output path handling
- Help message validation for all command groups

#### Testing Strategy:

- Uses Click's `CliRunner` for isolated command execution
- Temporary directories for output isolation
- Mocking with `unittest.mock` for unit-level CLI testing
- Fixtures for reusable test images and directory mocking

### 3.3.3 API Integration Tests (test\_api.py)

Tests for the FastAPI web service endpoints:

## Endpoints Tested:

1. **Home Page (GET /):**
  - Validates HTML response and correct content type
2. **Predict Endpoint (POST /predict):**
  - Tests with real images returning valid class predictions
  - Tests with mocked predictions for deterministic validation
  - Error handling for invalid file types
  - Validates JSON response structure (`predicted_class`, `filename`)
3. **Resize Endpoint (POST /resize):**
  - Fixed size resizing with explicit width/height
  - Random size resizing when dimensions not provided
  - Validates returned image dimensions
  - Error handling for invalid dimensions (negative values)
  - Mocked resize function testing
4. **Get Output File (GET /outputs/{filename}):**
  - Tests retrieval of existing files
  - Error handling for non-existent files

## Testing Approach:

- Uses FastAPI's `TestClient` for HTTP request simulation
- In-memory image generation with PIL for test isolation
- Temporary output directories to avoid filesystem pollution
- **Fixtures for expected class labels** with fallback for CI/CD environments
- Both real integration tests and mocked unit tests for comprehensive coverage

### 3.3.4 Artifact Existence Tests (`test_artifacts.py`)

Critical pre-deployment validation ensuring required artifacts exist:

```
def test_model_artifacts_exist():
    """Test that required model artifacts exist before deployment"""
    model_path = Path("model.onnx")
    labels_path = Path("class_labels.json")

    assert model_path.exists(), f"Model file not found at {model_path}"
    assert labels_path.exists(), f"Class labels file not found at {labels_path}"
```

These tests run before containerization to catch missing artifacts early, preventing deployment failures on platforms like Render or HuggingFace Spaces.

## 3.4 Testing Fixtures and Utilities

### Shared Fixtures Across Test Modules:

- **expected\_classes:** Session-scoped fixture that loads class labels from `class_labels.json` with fallback to hardcoded list for CI environments
- **test\_image / dummy\_image:** Creates temporary test images for validation
- **tmp\_outputs\_dir / mock\_outputs\_dir:** Provides isolated temporary directories for output files
- **runner:** CLI testing with Click's `CliRunner`

### Mocking Strategy:

The test suite uses `unittest.mock` to:

- Mock ONNX classifier availability and predictions
- Mock random number generation for deterministic testing
- Patch file system paths for isolated testing
- Mock image processing functions to test integration without full pipeline execution

## COVERAGE

```
tests/test_logic.py::test_preprocess_returns_image PASSED [ 91%]
tests/test_logic.py::test_preprocess_produces_grayscale PASSED [ 93%]
tests/test_logic.py::test_preprocess_call_order PASSED [ 94%]
tests/test_logic.py::test_ensure_output_dir PASSED [ 95%]
tests/test_logic.py::test_ensure_output_dir_idempotent PASSED [ 97%]
tests/test_logic.py::test_class_labels_loaded PASSED [ 98%]
tests/test_logic.py::test_class_labels_content PASSED [100%]

===== tests coverage =====
coverage: platform win32, python 3.11.13-final-0

Name                               Stmts  Miss  Cover
-----
api\api.py                          54      0  100%
cli\cli.py                          92      0  100%
logic\onnx_classifier.py            51      0  100%
logic\utilities.py                  72      0  100%
TOTAL                               269      0  100%
Coverage HTML written to dir htmlcov

===== 72 passed in 3.20s =====
```

## LINT

```
-----
Your code has been rated at 10.00/10 (previous run: 10.00/10, +0.00)
```

# 4 Experiments Conducted

## 4.1 Experimental Setup

### 4.1.1 Dataset

- **Dataset:** Oxford-IIIT Pet Dataset
- **Number of classes:** 37
- **Image preprocessing:** Resized to 256x256, take center crop of 224x224 and normalized using ImageNet statistics (mean and std)
- **Train/Validation split:** 80/20
- **Reproducibility:** Fixed random seeds (42) for dataset splitting and model initialization

### 4.1.2 Base Model Architecture

Describe the models you experimented with:

- **Primary model:** MobileNet\_v2 with IMAGENET1K\_V1 pretrained weights
- **Transfer learning approach:** Frozen feature extractor, modified classifier head
- **Optimizer:** Adam
- **Loss function:** Cross-entropy loss

## 4.2 Hyperparameter Configurations

Run Name	Model	Batch Size	Learning Rate	Epochs
mobilenet_v2_bs32_lr0.001	MobileNet_v2	32	0.001	3
mobilenet_v2_bs32_lr0.0001	MobileNet_v2	64	0.0001	3
mobilenet_v2_bs32_lr0.0005	MobileNet_v2	32	0.0005	5

Table 1: Experimental configurations tested

## 4.3 Logged Artifacts

### 4.3.1 Parameters Logged

For each experiment run, the following parameters were logged to MLFlow:

- Batch Size
- Dataset
- Epochs
- Image Size
- Learning Rate
- Loss Functions
- Model Name
- Number of classes
- Optimizer
- Seeds
- Train split

### 4.3.2 Metrics Logged

The following metrics were tracked throughout training:

- **Training accuracy** (per epoch)
- **Validation accuracy** (per epoch)
- **Training loss** (per epoch)
- **Validation loss** (per epoch)
- **Final training accuracy**
- **Final validation accuracy**

### 4.3.3 Artifacts Logged

Additional artifacts saved with each run:

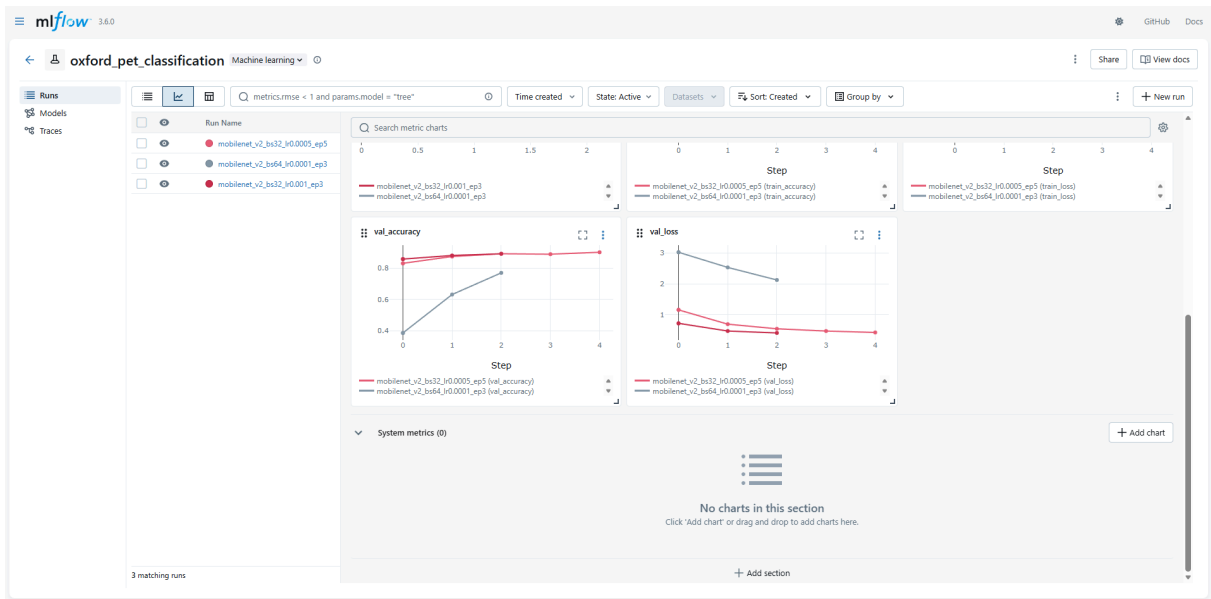
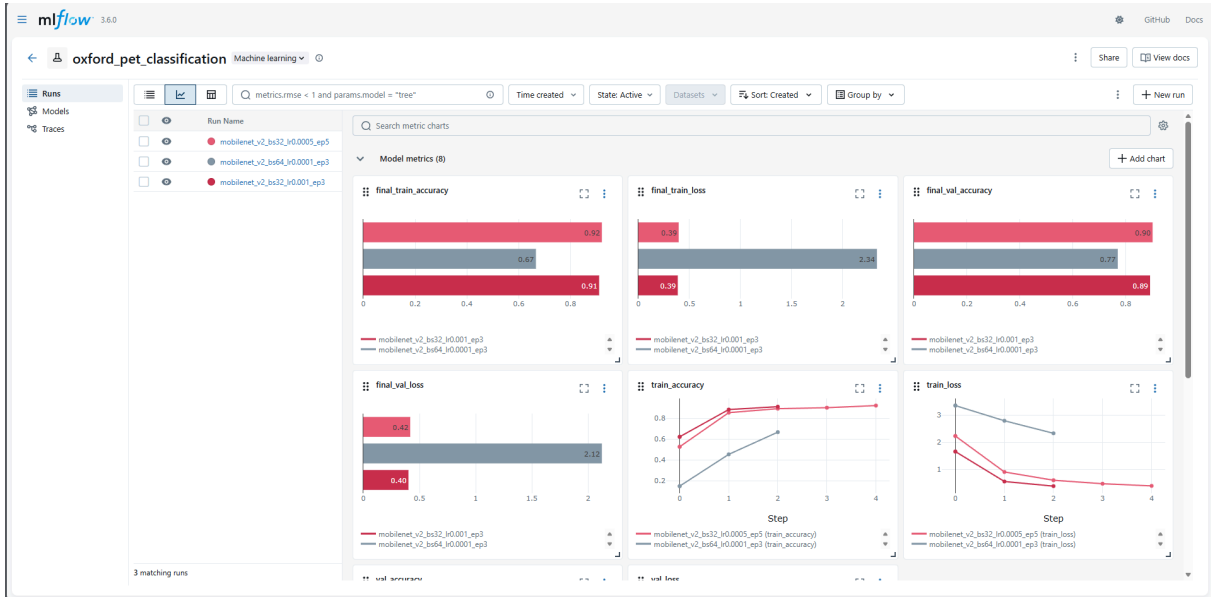
- **Training loss curve** (PNG image)
- **Validation loss curve** (PNG image)
- **Accuracy curves** (PNG image)
- **Class labels** (JSON file)
  - Example content:

```
[  
    "Abyssinian",  
    "american_bulldog",  
    "american_pit_bull_terrier",  
    ...  
]
```

- **Trained model** (PyTorch format, registered with MLFlow)

## 5 Results Analysis

### 5.1 MLFlow UI Analysis



#### 5.1.1 Model Selection

Based on the `final_val_accuracy` metric comparison, the models ranked as follows:

Metric	Best Model	Second Best	Third Best
Run Name	...bs32_lr0.0005_ep5	...bs32_lr0.001_ep3	...bs64_lr0.0001_ep3
Validation Accuracy	90%	89%	77%
Training Accuracy	91%	91%	67%
Final Val Loss	0.40	0.40	2.12
Final Train Loss	0.30	0.30	0.94

Table 2: Model performance comparison from MLFlow UI

**Selected Model for Production:**



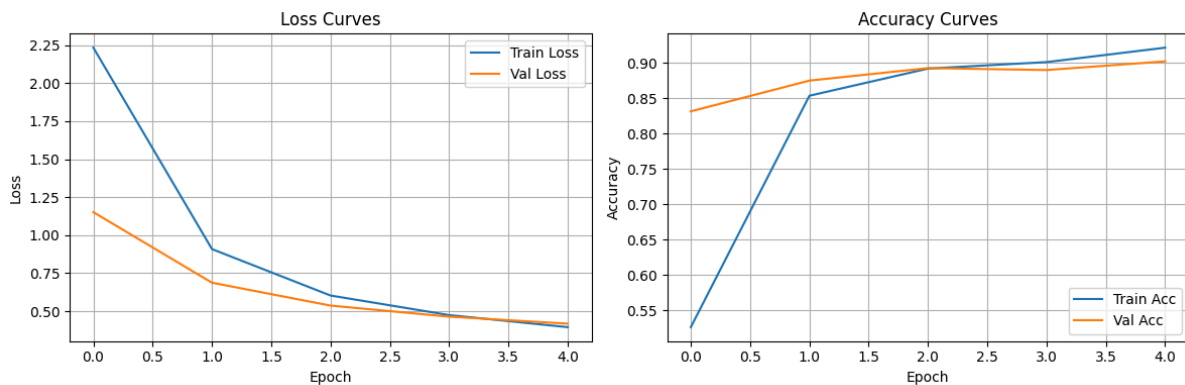
- **Run Name:** mobilenet\_v2\_bs32\_lr0.0005\_ep5
- **Configuration:** Batch size 32, Learning rate 0.0005, 5 epochs
- **Final Validation Accuracy:** 90%
- **Final Training Accuracy:** 91%

#### Justification for Selection:

Achieved the highest validation accuracy with strong generalization ( 1% gap), stable training curves, low validation loss ( 0.40), and well-balanced hyperparameters (lr=0.0005, batch size=32).

## 5.2 Performance Analysis

### 5.2.1 Training Curves



### 5.2.2 Final Model Performance

The selected model achieved:

- **Validation Accuracy:** 90%
- **Training Accuracy:** 91%
- **Generalization Gap:** 1% (good generalization)

## 5.3 Model Serialization

The best performing model was:

1. Loaded from MLFlow model registry
2. Converted to evaluation mode
3. Serialized to ONNX format (opset version 18)
4. Validated for inference compatibility

**ONNX model size:** 0.145 MB

**ONNX model data size:** 8.83 MB

# 6 Implementation Details

## 6.1 Model Selection Script

Automated selection process:

- Query all registered model versions
- Extract validation metrics for each version
- Select model with highest validation accuracy
- Serialize selected model to ONNX format
- Export class labels as JSON

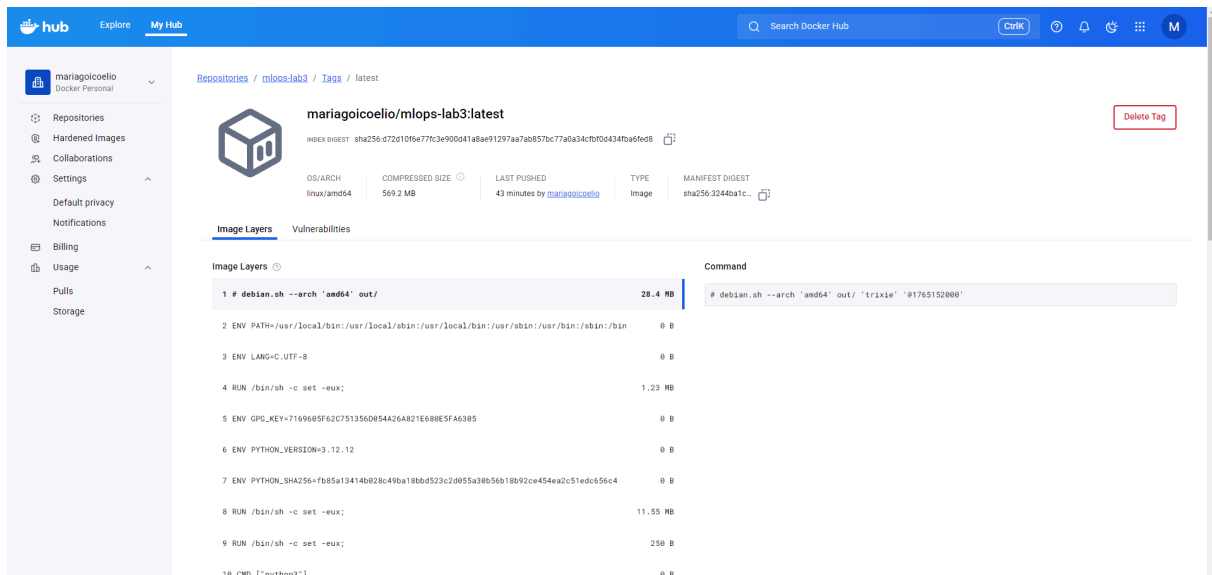
## 6.2 Inference Script

ONNX Runtime integration:

- InferenceSession with CPU execution provider
- Image preprocessing pipeline
- Prediction with class label mapping
- API integration for production deployment

## 7 Services Working

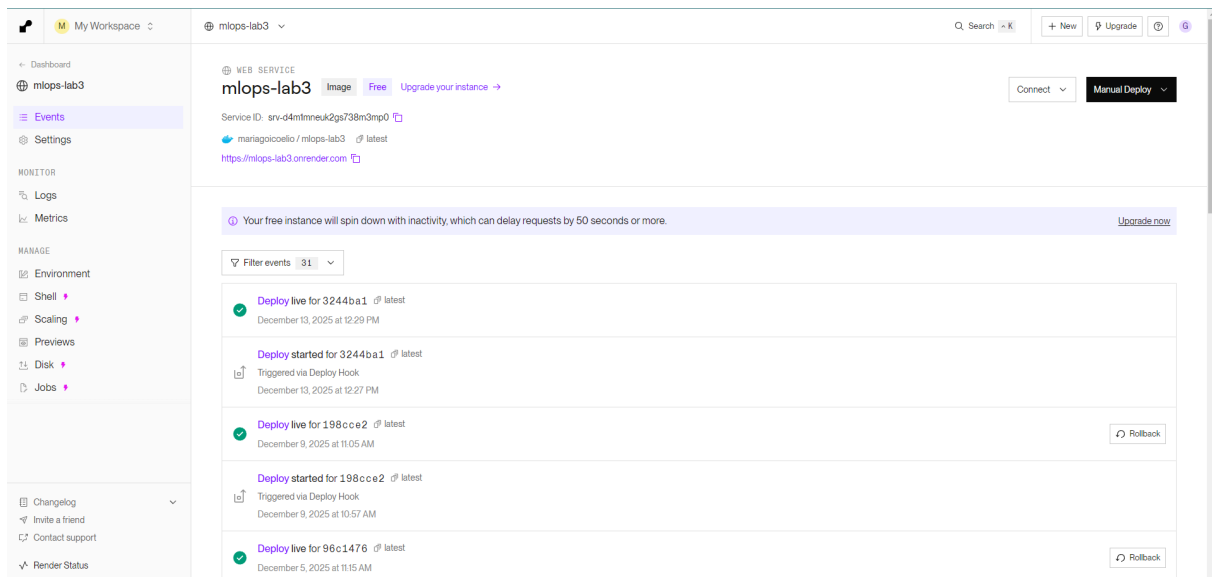
### 7.1 DockerHub



The screenshot shows the Docker Hub interface for the repository `mariagoicoelio/mlops-lab3:latest`. The page includes a sidebar with navigation options like Repositories, Hardened Images, and Settings. The main content area displays the repository details, including the image layers and the command to run the container. The image layers table lists 10 layers, with the first layer being the base image `debian.sh --arch 'amd64' out/` and the last layer being the command `python3`.

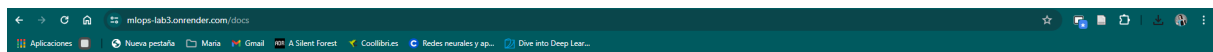
Image Layers	Command
1 # debian.sh --arch 'amd64' out/	# debian.sh --arch 'amd64' out/ 'trixie' '#1765152888'
2 ENV PATH=/usr/local/bin:/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/sbin:/bin	
3 ENV LANG=C.UTF-8	
4 RUN /bin/sh -c set -eux;	
5 ENV GPU_KEY=7169685f62c751356d85426a821e680e5fa6305	
6 ENV PYTHON_VERSION=3.12.12	
7 ENV PYTHON_SHA256=f85a13414b828c49ba18bd523c2d855a38b56b18992ce454ea2c51edc656c4	
8 RUN /bin/sh -c set -eux;	
9 RUN /bin/sh -c set -eux;	
10 CMD ["python3"]	

### 7.2 Render



The screenshot shows the Render dashboard for the service `mlops-lab3`. The dashboard includes a sidebar with navigation options like Dashboard, Events, Settings, and Monitor. The main content area displays the service details, including the service ID, the image used, and the deployment history. The deployment history table lists 4 deployments, with the first deployment being the latest one.

Deployment	Status	Timestamp	Action
Deploy live for 3244ba1	Success	December 13, 2025 at 12:29 PM	
Deploy started for 3244ba1	In Progress	December 13, 2025 at 12:27 PM	
Deploy live for 198cce2	Success	December 9, 2025 at 11:05 AM	Rollback
Deploy started for 198cce2	In Progress	December 9, 2025 at 10:57 AM	
Deploy live for 96c1476	Success	December 5, 2025 at 11:15 AM	Rollback



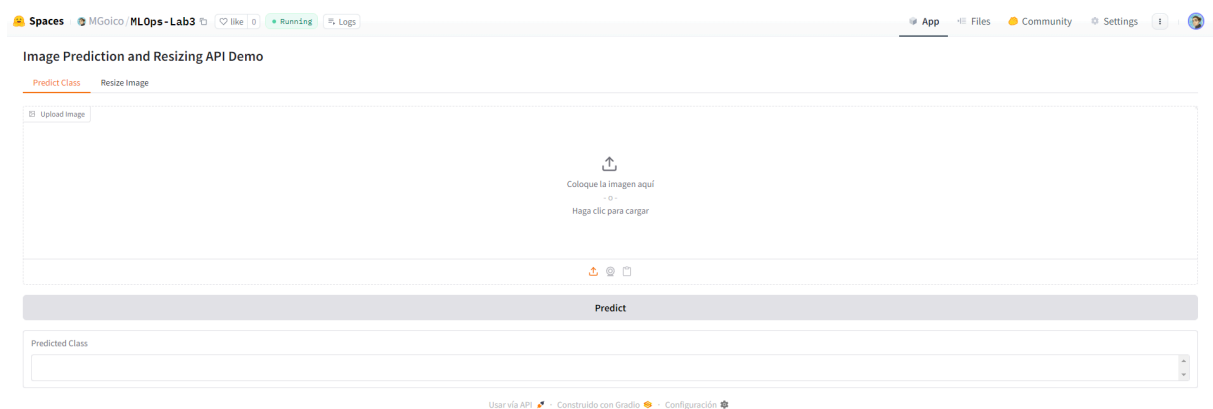
## API of the Lab 1 using FastAPI 1.0.0 OAS 3.1

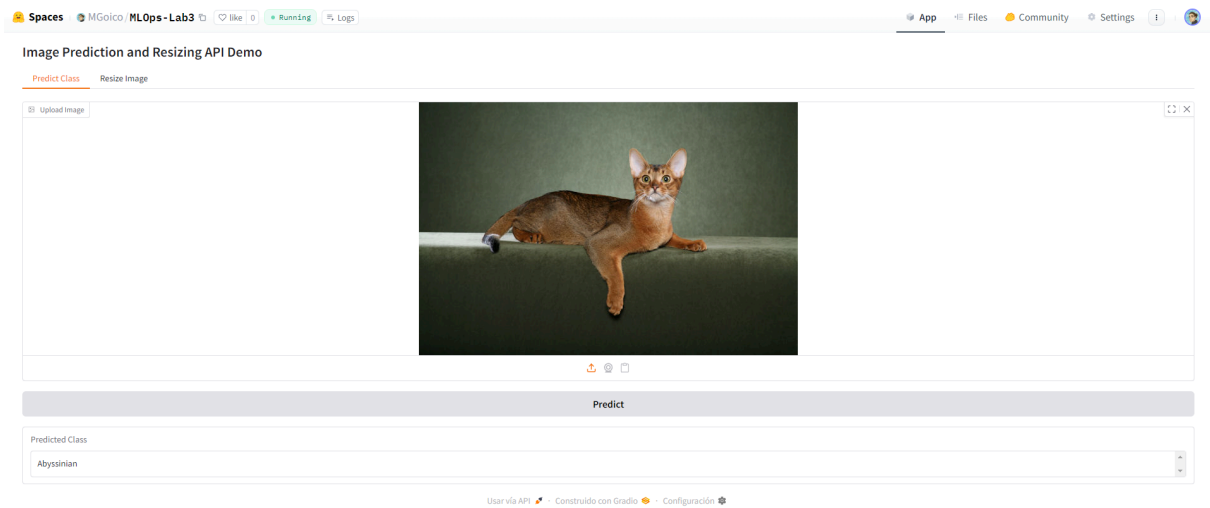
[OpenAPI JSON](#)

API to perform preprocessing on images

default		^
GET	/ Home	▼
POST	/predict Predict Class	▼
POST	/resize Resize Image	▼
GET	/outputs/{filename} Get Output File	▼
Schemas		^
Body_predict_class_predict_post > Expand all object		
Body_resize_image_resize_post > Expand all object		
HTTPValidationError > Expand all object		
ValidationError > Expand all object		

## 7.3 HuggingFace





## 8 Challenges and Solutions

### 8.1 Docker

**Challenge:** Took too long to build and deployment failed.

**Solution:** Modified the pytorch version so that it does not include the nvidia CUDA libraries since our model will only be used for inference and no further training will be performed.