

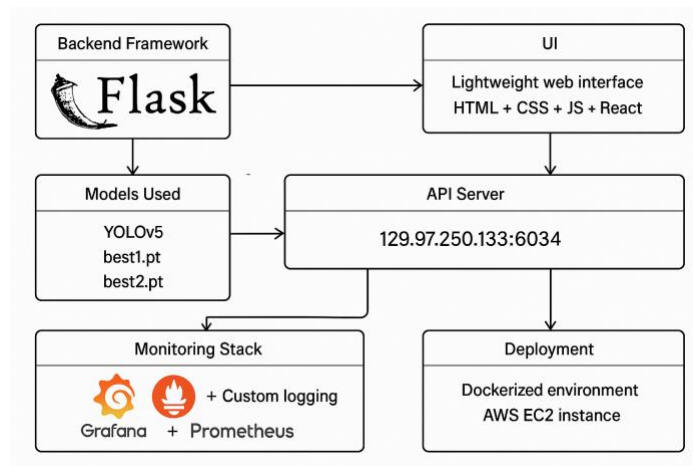
SYDE 770 – Assignment 4 Report

1. Introduction

This project focused on deploying trained YOLO object detection models using a RESTful API. The system was designed to handle inference requests, support multiple models, and provide monitoring capabilities using tools like Grafana and Docker. A front-end UI was also built for user interaction.

2. System Architecture

- Backend Framework: Flask
- Models Used: YOLOv5 (`best1.pt`, `best2.pt`)
- API Server: Hosted on `129.97.250.133:6034`
- Monitoring Stack: Grafana + Prometheus + Custom logging
- Deployment: Dockerized environment
- UI: Lightweight web interface using HTML + CSS + JS + React



3. REST API Implementation

Our object detection service is implemented as a Flask-based REST API with the following components:

Model Management: Handles loading and switching between YOLO models

Prediction Endpoint: Processes image inputs and returns detection results

Monitoring System: Tracks performance metrics and system health

Evaluation Framework: Assesses model performance on validation datasets

1. Complete REST API implementation

1.1 Prediction API

Endpoint: POST /predict

Description: Processes an image and returns detected objects

```
@app.route('/predict', methods=['POST'])
def predict():
    # Input validation
    if 'image' not in request.files:
        return jsonify({"error": "No image provided"}), 400

    # Model selection with A/B testing
    model_name = request.form.get('model', current_default_model)
    if 'model' not in request.form and random.random() < 0.5:
        model_name = "best2.pt" if current_default_model == "best1.pt" else "best1.pt"

    # Processing pipeline
    results = model(process_image(file))
    return jsonify(format_predictions(results, model_name))
```

POST {{base_url}}/predict Send

Params Authorization Headers (8) Body Scripts Tests Settings Cookies

☐ none ☒ form-data ☐ x-www-form-urlencoded ☐ raw ☐ binary ☐ GraphQL

Key	Value	Description	Bulk Edit
<input checked="" type="checkbox"/> image	File GenAI_img_3.1.jpeg		
Key	Text Value	Description	

Body Cookies Headers (5) Test Results (2/2) 200 OK 172 ms 294 B Save Response

JSON Preview Visualize

```
1 {
2   "model_used": "best1.pt",
3   "predictions": [
4     {
5       "bbox": [
6         382,
7         338,
8         849,
9         1135
10      ],
11      "confidence": 0.9739964604377747,
12      "label": "Tim Hortons Cup"
13    }
14  ]
15 }
```

Predict with specific model either model 1 or model 2.

YOLO Model API Testing / Single Image Prediction / Successful Prediction (Specific Model) Save Share

POST {{base_url}}/predict Send

Params Authorization Headers (8) Body Scripts Tests Settings Cookies

☐ none ☒ form-data ☐ x-www-form-urlencoded ☐ raw ☐ binary ☐ GraphQL

Key	Value	Description	Bulk Edit
<input checked="" type="checkbox"/> image	File GenAI_img_3.2.jpeg		
<input checked="" type="checkbox"/> model	Text best2.pt		
Key	Text Value	Description	

Body Cookies Headers (5) Test Results (2/2) 200 OK 175 ms 282 B Save Response

JSON Preview Visualize

```
1 {
2   "model_used": "best2.pt",
3   "predictions": [
4     {
5       "bbox": [
6         395,
7         356,
8         866,
9         1106
10      ],
11      "confidence": 0.9598243832588196,
12      "label": "Tim"
13    }
14  ]
15 }
```

1.2. Evaluation API

1.2.1 Batch Evaluation

Endpoint: POST /evaluate

Description: Evaluates model performance on a dataset

Parameters:

-model (required): Model name

-dataset_path (optional): Path to evaluation dataset

```
@app.route('/evaluate', methods=['POST'])
def evaluate_model():
    model_name = request.form['model']
    dataset_path = request.form.get('dataset_path', 'data/eval')

    # Evaluation logic
    results = model.val(data=dataset_path)
    return jsonify({
        "precision": results.results_dict['metrics/precision'],
        "recall": results.results_dict['metrics/recall']
    })
```

The screenshot shows a REST client interface with the following details:

- URL:** `POST {{base_url}}/evaluate`
- Params:** `model` (Text, Value: `best2.pt`) and `dataset_path` (Text, Value: `data/eval`).
- Body:** The response is a JSON object with the following structure:

```
{
  "class_metrics": {
    "0": {
      "precision": 0.0,
      "recall": 0
    },
    "1": {
      "precision": 1.0,
      "recall": 0.9902439024390244
    }
  },
  "f1_score": 0.9806763285024155,
  "map50": 0.999,
  "map50_95": 0.9915470129031091,
  "model": "best2.pt",
  "precision": 0.9712918660287081,
  "processed_files": 201,
  "recall": 0.9902439024390244
}
```
- Status:** 200 OK

1.2.2 Model comparison:

Endpoint: GET / metrics/compare?time_range=24

Description: Returns model statistics for performance comparison.

YOLO Model API Testing / Metrics & Monitoring / Compare Models

GET Send

Params Authorization Headers (6) Body Scripts Tests Settings Cookies

Query Params

<input checked="" type="checkbox"/> Key	Value	Description	Bulk Edit
<input checked="" type="checkbox"/> time_range	24		
Key	Value	Description	

Body Cookies Headers (5) Test Results (1/2) 200 OK - 4 ms - 692 B Save Response

{ } JSON Preview Visualize

```
1 {
2   "best1.pt": {
3     "avg_map50": 0.995,
4     "avg_map50_95": 0.9865461214701645,
5     "avg_precision": 0.966824644549763,
6     "avg_recall": 0.9951219512195122,
7     "latest_map50": 0.995,
8     "latest_map50_95": 0.9865461214701645,
9     "latest_precision": 0.966824644549763,
10    "latest_recall": 0.9951219512195122
11  },
12  "best2.pt": {
13    "avg_map50": 0.995,
14    "avg_map50_95": 0.9915470129031091,
15    "avg_precision": 0.9712918660287081,
16    "avg_recall": 0.9902439024390244,
17    "latest_map50": 0.995,
18    "latest_map50_95": 0.9915470129031091,
19    "latest_precision": 0.9712918660287081,
20    "latest_recall": 0.9902439024390244
21  }
22 }
```

1.3.1 System Metrics

Endpoint: GET /metrics

Description: Returns current system performance metrics

```
@app.route('/metrics', methods=['GET'])
def get_metrics():
    return jsonify({
        "request_rate": metrics['request_rate'],
        "latency": metrics['avg_latency']
    })
```

YOLO Model API Testing / Metrics & Monitoring / Get System Metrics

GET `{{base_url}}/metrics` Send

Params Authorization Headers (6) Body Scripts Tests Settings Cookies

Query Params

Key	Value	Description	Bulk Edit
Key	Value	Description	

Body Cookies Headers (5) Test Results (2/2) 200 OK 6 ms 403 B Save Response

JSON Preview Visualize

```
1 {
2   "system_metrics": {
3     "avg_latency_ms": 462.6,
4     "error_rate": 0.0,
5     "max_latency_ms": 754.9,
6     "min_latency_ms": 170.3,
7     "model_usage": {
8       "best1.pt": 1,
9       "best2.pt": 1
10    },
11     "request_rate_per_minute": 0,
12     "total_requests": 2
13   },
14   "timestamp": "2025-04-07T18:10:53.516083"
15 }
```

1.3.2 Model Metrics

Endpoint: GET /metrics/models/<model_name>

Description: Returns metrics for specific model

```
@app.route('/management/models', methods=['GET'])
def list_models():
    return jsonify({"models": list(MODEL_INFO.keys())})
```

YOLO Model API Testing / Metrics & Monitoring / Get Model Metrics

GET `{{base_url}}/metrics/models/best2.pt?time_range=24` Send

Params Authorization Headers (6) Body Scripts Tests Settings Cookies

Query Params

Key	Value	Description	Bulk Edit
<input checked="" type="checkbox"/> time_range	24		
Key	Value	Description	

Body Cookies Headers (5) Test Results (1/2) 200 OK 5 ms 841 B Save Response

JSON Preview Visualize

```
19 {
20   "1": {
21     "precision": [
22       {
23         "timestamp": "2025-04-07T18:06:05.863828",
24         "value": 1.0
25       }
26     ],
27     "recall": [
28       {
29         "timestamp": "2025-04-07T18:06:05.863828",
30         "value": 0.9902439024390244
31       }
32     ]
33   },
34   "map50": [
35     {
36       "timestamp": "2025-04-07T18:06:05.863828",
37       "value": 0.995
38     }
39   ],
40   "map50_95": [
41     {
42       "timestamp": "2025-04-07T18:06:05.863828",
43       "value": 0.9915470129031091
44     }
45   ]
46 }
```

1.4. Model Management API

1.4.1 List Models

Endpoint: GET /management/models

Description: Lists all available models

```
@app.route('/management/models', methods=['GET'])
def list_models():
    return jsonify({"models": list(MODEL_INFO.keys())})
```

Body Cookies Headers (5) Test Results (1/1) 200 OK • 17 ms • 210 B Save Response

{ } JSON Preview Visualize

```
1 {
2   "available_models": [
3     "best1.pt",
4     "best2.pt"
5   ]
6 }
```

1.4.2 Model Description

Endpoint: GET /management/models/{model}/describe

Description: Provides metadata for the selected model

YOLO Model API Testing / Model Management / Describe Model Save Share

GET {{base_url}}/management/models/best1.pt/describe Send

Params Authorization Headers (6) Body Scripts Tests Settings Cookies

Query Params

Key	Value	Description	Bulk Edit
Key	Value	Description	

Body Cookies Headers (5) Test Results (1/1) 200 OK • 10 ms • 433 B Save Response

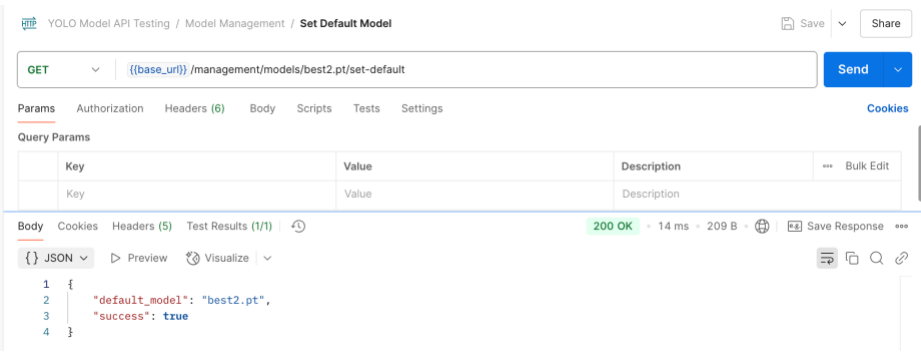
{ } JSON Preview Visualize

```
1 {
2   "config": {
3     "batch_size": 64,
4     "confidence_threshold": 0.4,
5     "input_size": [
6       416,
7       416
8     ]
9   },
10  "date_registered": "2025-03-13",
11  "metrics_summary": {
12    "accuracy": null,
13    "map50": 0.995,
14    "map50_95": 0.9865461214701645,
15    "precision": 0.966824644549763,
16    "recall": 0.9951219512195122
17  },
18  "model": "best1.pt"
19 }
```

1.4.3 Set Default Model

Endpoint: GET /management/models/{model}/set-default

Description: Updates default inference model.



1.5. Health Check API

1.5.1 System Health

Endpoint: GET /health-status

Description: Returns system health status

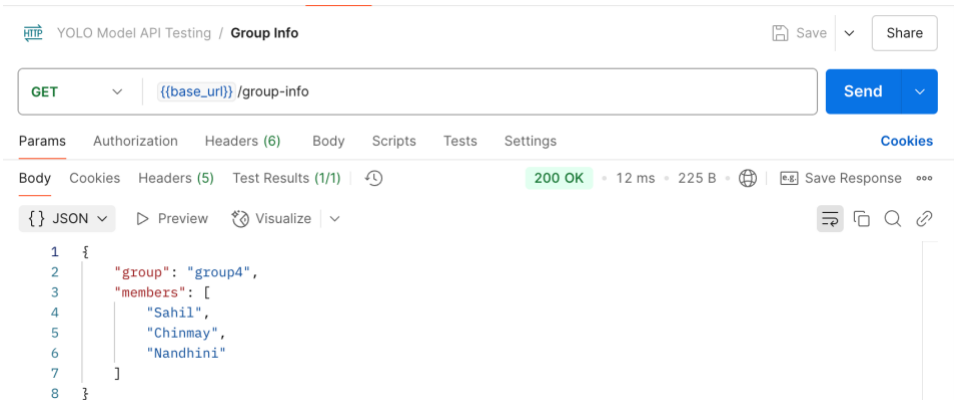
```

1 {
2   "server": "Flask",
3   "status": "Healthy",
4   "uptime": "1:36:30"
5 }
```

6.Group metadata

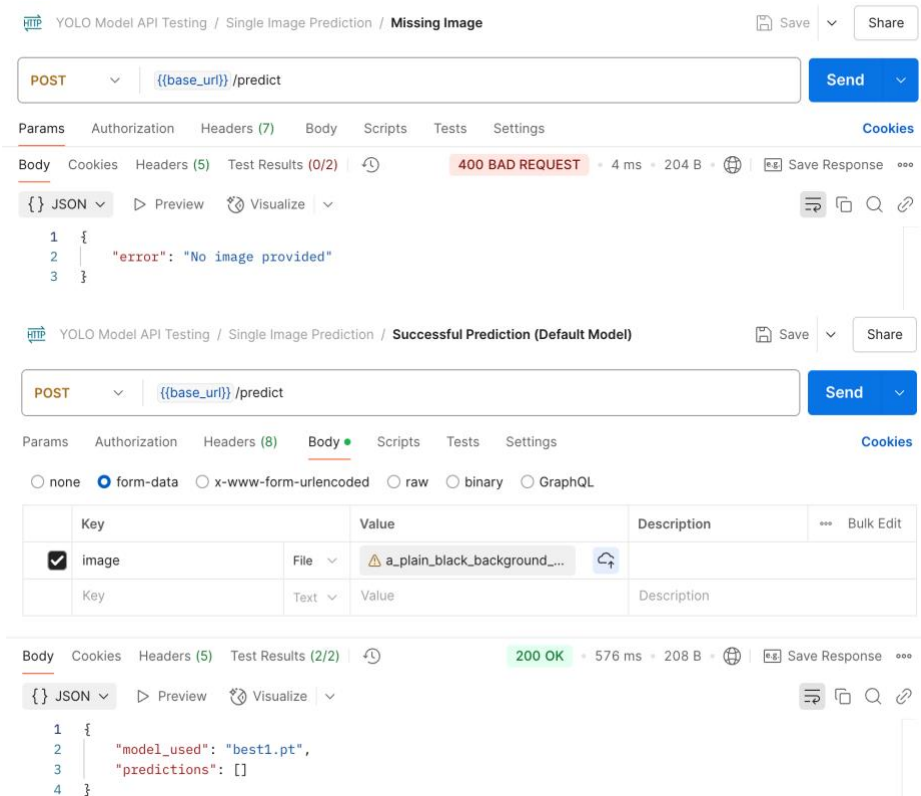
Endpoint: GET /group-info

Description: Reflects the group information.



API Testing Matrix

Endpoint	Test Cases Covered	Status Codes Tested
/predict	5 (valid, invalid, missing image)	200, 400, 500
/evaluate	3 (default path, custom path)	200, 400
/metrics	2 (normal, alert state)	200
Model Management	4 (list, describe, set default)	200, 404
Health Check	1	200



Non-Cup image with no detections.

Key Features

1. Model A/B Testing:

Randomly routes 50% of unspecified requests to alternate model

```
if 'model' not in request.form and random.random() < 0.5:
    model_name = "best2.pt" if current_default_model == "best1.pt" else "best1.pt"
```

2. Performance Monitoring:

Background thread tracks:

```
def monitor_metrics():
    while True:
        time.sleep(60)
        print(f"Current metrics: {calculate_metrics()}")
```

3. Alert System:

Detects performance degradation:

```
if avg_latency > 2.5: # 2.5 second threshold
    return "High latency alert"
```


4. Comprehensive Error Handling

```
try:
    image = process_image(file)
except Exception as e:
    metrics["errors"].append(str(e))
    return jsonify({"error": str(e)}), 500
```

5. Flexible Model Management

```
def load_model(model_name):
    if model_name not in models:
        models[model_name] = YOLO(os.path.join(MODELS_DIR, model_name))
```

2. User Interface

2.1 App Structure -

```
/src
├── /components
│   ├── ImageUploader.js
│   └── MetricsTable.jsx
├── /api
│   └── apiHelper.js
├── App.js
└── index.js
```

2.2 App Overview

The frontend user interface (UI) of the application is structured as a single-page React application, designed for ease of use, real-time interaction, and modularity. The UI primarily consists of two core components that facilitate model evaluation and image-based detection: the **Model Comparison Component** and the **Image Upload & Detection Component**. Each component is modular and built using React functional components with Material-UI for consistent design.

1. Model Comparison Component (MetricsTable)

This component is responsible for rendering a side-by-side comparison of two different YOLO-based detection models. It fetches model performance metrics from a backend API endpoint and displays them in a clean tabular format.

The goal of this component is to assist users (developers or evaluators) in quantitatively assessing the performance differences between two trained object detection models, allowing data-driven decision making when selecting a model for deployment.

2. Image Upload & Detection Component (ImageUploader)

This component provides an intuitive interface for uploading images and choosing which detection model to use for inference. After an image is uploaded, it is sent to the backend, which returns bounding box predictions, labels, and confidence scores.

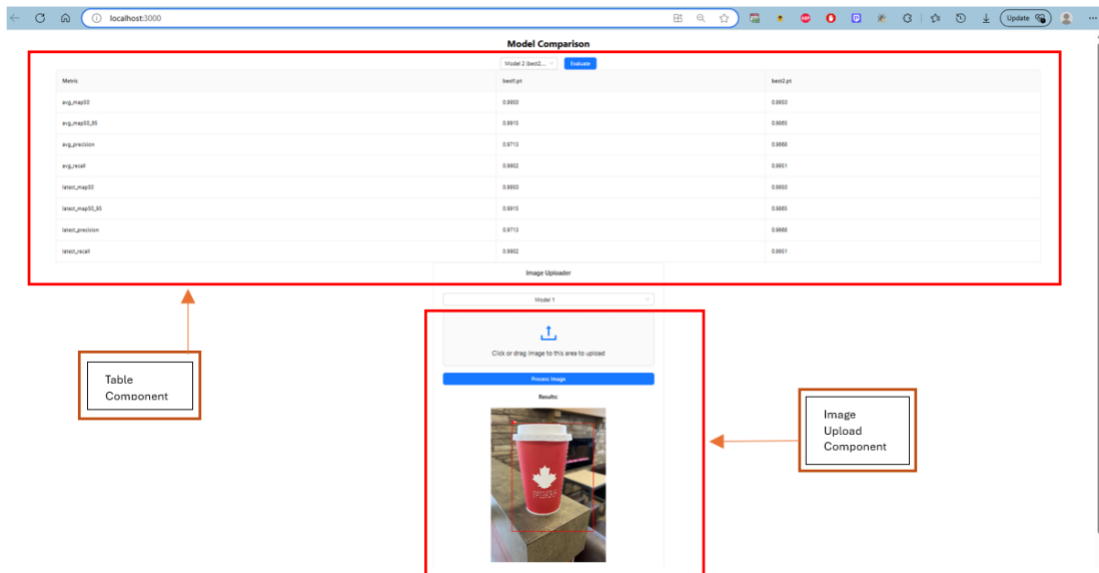
Technical Implementation:

- Uses React hooks (useState, useEffect) for managing local state (selected model, image, and status).
- API call handled via a centralized apiHelper.js configuration to ensure consistency and maintainability across different endpoints.

Integration Summary

Both components are rendered within the main App.js component and work together to offer a full detection evaluation flow:

1. Users **compare models** using the Metrics Table.
2. Users **select a model and upload an image** for real-time inference and validation



3. Docker Deployment

This system leverages a microservice architecture with monitoring and observability, deployed via Docker Compose.



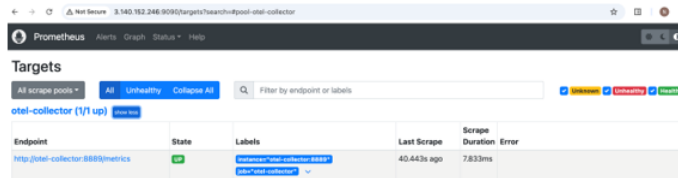
- Hosts the Flask-based inference API using YOLOv8 models (best1.pt, best2.pt), Exposed on port 6034
- Instrumented with OpenTelemetry for tracing and metrics

otel-collector

- Collects telemetry data from the Flask app
- Converts and forwards telemetry to Prometheus-compatible format
- Uses custom configuration from otel-collector-config.yaml
- Ports:
 - 4318: Receives OTLP traces/metrics from Flask
 - 8889: Exposes metrics in Prometheus format

prometheus

- Scrapes metrics from otel-collector (on port 8889)
- Configured via prometheus.yml



The screenshot shows the Prometheus web interface at the 'Targets' page. It displays a table with one target named 'otel-collector (1/1 up)'. The table has columns for Endpoint, State, Labels, Last Scrape, Scrape Duration, and Error. The 'otel-collector' target is in a 'UP' state, with a last scrape time of 40.643s ago and a duration of 7.833ms.

Endpoint	State	Labels	Last Scrape	Scrape Duration	Error
http://otel-collector:8889/metrics	UP	instance="otel-collector-8889" otel-collector="otel-collector"	40.643s ago	7.833ms	

Grafana

- Visualizes metrics from Prometheus
- Dashboards for application health, response times, system load
- Default port: 3000 (login: admin / admin by default)

react-frontend

- Lightweight UI with HTML/CSS/JS/React
- Built in folder ./react-frontend and exposed on port 3001

AWS

To securely and efficiently deploy a Dockerized YOLOv5 Flask API on an AWS EC2 instance, begin by choosing a suitable instance type (e.g., t2.medium)

Instances (1/1) Info Last updated 1 minute ago Connect Instance state Actions Launch instances

Find instance by attribute or tag (case-sensitive) All states

Instance state = running Clear filters

Name	Instance ID	Instance state	Instance type	Status check	Alarm status	Availability Zone
<input checked="" type="checkbox"/>	i-0e709dadb9a09bb7c	Running	t2.micro	2/2 checks passed	View alarms +	us-east-2a

i-0e709dadb9a09bb7c

Details | Status and alarms | Monitoring | Security | Networking | Storage | Tags

▼ Instance summary Info

Instance ID i-0e709dadb9a09bb7c	Public IPv4 address 3.140.152.246 open address	Private IPv4 addresses 172.31.13.171
IPv6 address -	Instance state Running	Public IPv4 DNS ec2-3-140-152-246.us-east-2.compute.amazonaws.com

Associating an Elastic IP for a stable external address.

EC2 > Elastic IP addresses > 3.140.152.246 Options

3.140.152.246 Actions Associate Elastic IP address

Summary

Allocated IPv4 address 3.140.152.246	Type Public IP	Allocation ID eipalloc-0748013a59ff74cfa	Reverse DNS record -
Association ID eipassoc-082884fc164a83c1f	Scope VPC	Associated instance ID i-0e709dadb9a09bb7c	Private IP address 172.31.13.171
Network interface ID eni-08a91b7f4b0e85ea6	Network interface owner account ID 375725685055	Public DNS ec2-3-140-152-246.us-east-2.compute.amazonaws.com	NAT Gateway ID -
Address pool Amazon			

Configure the Security Group to expose only necessary ports: 6034 (Flask API), 80 or 3001 (React frontend), 3000 (Grafana), 9090 (Prometheus), 4318 and 8889 (OpenTelemetry), and restrict SSH (port 22) access to your IP.

Instances (1/1) Info Last updated less than a minute ago Connect Instance state Actions Launch instances

Find instance by attribute or tag (case-sensitive) All states

Instance state = running Clear filters

Name	Instance ID	Instance state	Instance type	Status check	Alarm status	Availability Zone
<input checked="" type="checkbox"/>	i-0e709dadb9a09bb7c	Running	t2.micro	2/2 checks passed	View alarms +	us-east-2a

i-0e709dadb9a09bb7c

▼ Inbound rules

Filter rules

Name	Security group rule ID	Port range	Protocol	Source	Security group
-	sgr-0dba284a4bc6165a9	22	TCP	0.0.0.0/0	flask
-	sgr-0cae4d8cb1ddc718b	3000 - 9090	TCP	0.0.0.0/0	flask

Since model files are large and default volume of 8Gb is not sufficient to accomplish the models. The root volume is expanded to 20 Gb(/dev/xvda1) via the AWS console, then resize the partition and filesystem using growpart and resize2fs.



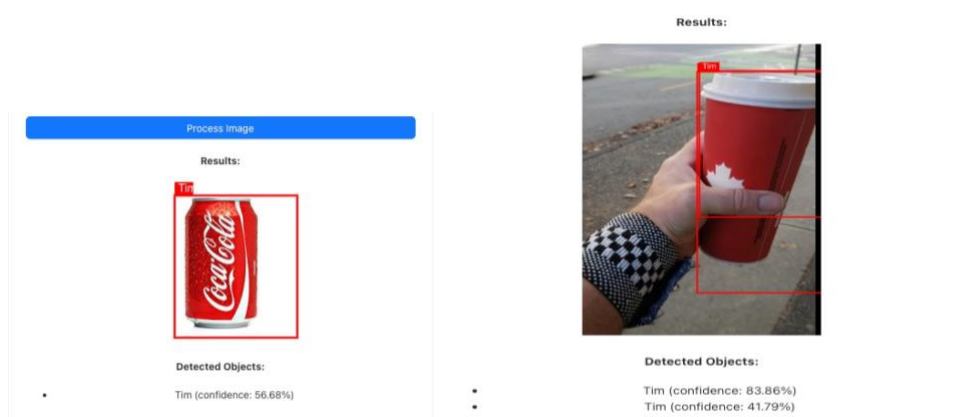
5. Evaluation Results

Using the provided evaluation script:

- All API requirements passed successfully
- Metrics returned correctly with standardized keys

```
syde770_group_4@SYDECLCLOUD:~/SYDE770$ curl -X POST -F "api_link=129.97.250.133:6034" http://129.97.250.133:7070/submit
{"message": "All tests passed\n"}syde770_group_4@SYDECLCLOUD:~/SYDE770$
```

During edge case testing, Model 2 demonstrated improved performance on edge cases compared to Model 1, which exhibited signs of overfitting on the edge-case test dataset. To enhance robustness, we plan to address some identified weaknesses—such as misclassification of lids-only portions, sleeve-bearing cups (e.g., Tim Hortons) where cup portion above and below sleeve or hand is treated separately or in an overlapping sense doubly detected, and red Coke with white text cans—by expanding the training dataset with additional edge-case examples. This refinement will improve generalization and ensure more reliable real-world deployment.



6. Individual Work

Sahil:

- Developed Flask-based backend REST APIs in accordance with the specified requirements.
- Designed and implemented APIs with Custom System and Model Monitoring metrics, alerts, Logging, Bulk Evaluation, Processing and Model Comparison.
- Configured and executed comprehensive Postman tests for all APIs to ensure functionality and reliability and tried but failed to integrate Evidently AI.

Chinmay:

- Developed a web-based application that enables users to compare the performance of Yolo models and run image-based inference using selected models.
- Curated and created targeted edge-case datasets (e.g., lids-only, sleeve-covered cups, red Coke cans) to stress-test model generalization and improve real-world robustness.
- Conducted comprehensive edge-case testing and final end-to-end evaluation, identifying key failure modes (e.g., overfitting in Model 1, overlapping detections), and informed iterative dataset refinement and model retraining.

Nandhini:

- Worked in Containerizing the application
- Used OTEL, Prometheus and Grafana for monitoring metrics and setup Dashboard & Panel.
- Deployed on AWS using a free-tier **t2.micro EC2 instance**. Despite limited resources, successfully expanded the root volume and configured a swap file to extend available RAM.

7. Conclusion

This project successfully demonstrates the deployment of deep learning models using a RESTful API, with integrated monitoring, visualization, and a user-friendly interface. The modular architecture ensures easy future extension, such as adding more models or using GPU-backed servers.