

Object Detection In Images Using Deep Learning

Team 63 group members:

Shweta Ganesh Bankar - sbankar@buffalo.edu

Sonam Barnala - sonambar@buffalo.edu

Deep learning models, trained on massive datasets, can analyze images and draw bounding boxes around objects, along with their labels. Object detection has applications in various fields, from autonomous vehicles to medical image analysis.

Dataset Preprocessing

The COCO 2017 dataset contains a rich set of images with object annotations for a variety of tasks including object detection, segmentation, and image captioning.

The dataset is divided into several splits for training, validation, and testing purposes.

Training Dataset : Approximately 118,000

Validation Dataset: Approximately 5,000

Test Dataset: Approximately 41,000

Categories: The dataset includes 80 object categories, ranging from everyday objects like "bottle" and "chair" to animals like "dog" and "cat".

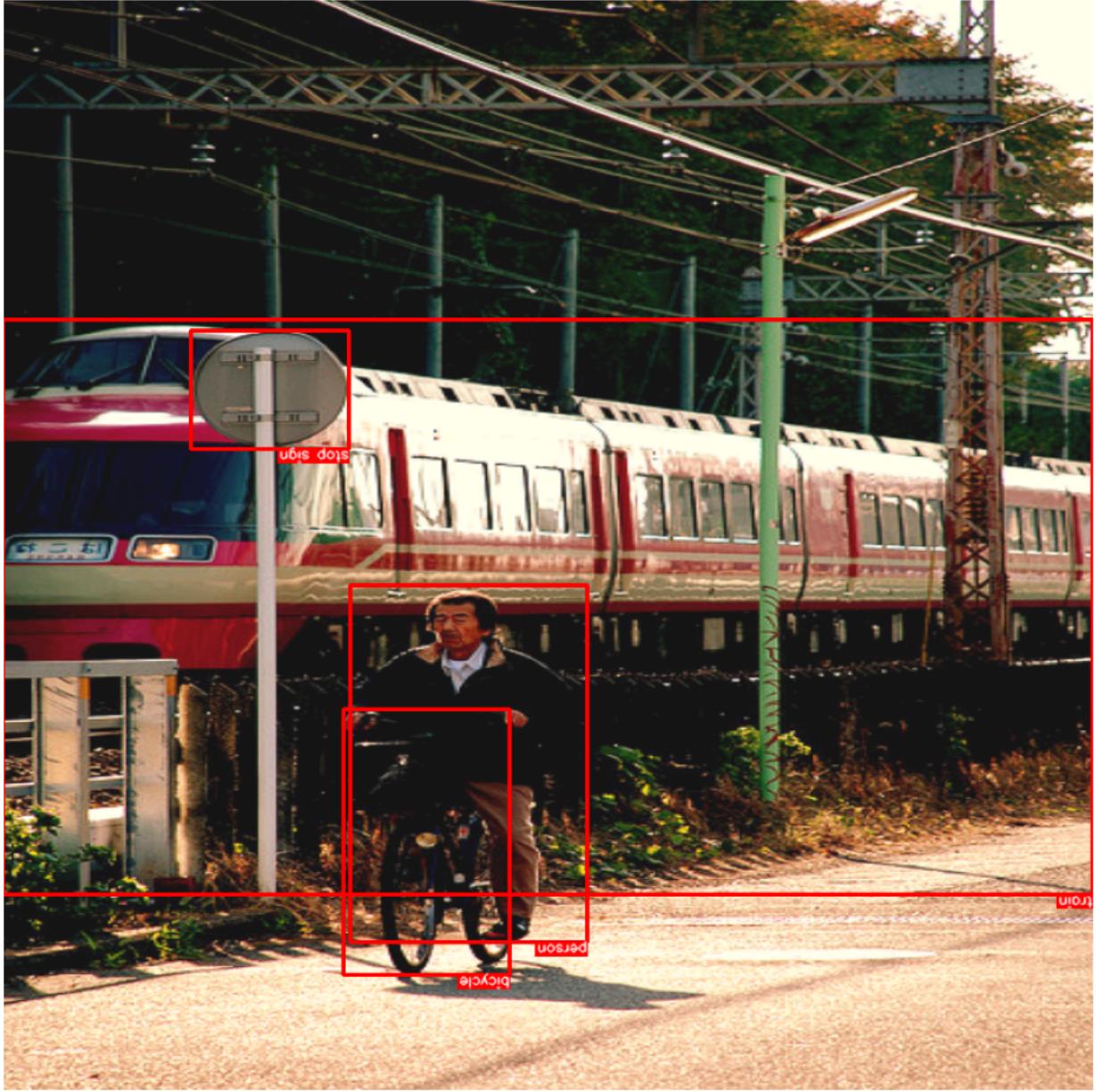
```
{0: 'person', 1: 'bicycle', 2: 'car', 3: 'motorcycle', 4: 'airplane', 5: 'bus', 6: 'train', 7: 'truck', 8: 'boat', 9: 'traffic light', 10: 'fire hydrant', 11: 'stop sign', 12: 'parking meter', 13: 'bench', 14: 'bird', 15: 'cat', 16: 'dog', 17: 'horse', 18: 'sheep', 19: 'cow', 20: 'elephant', 21: 'bear', 22: 'zebra', 23: 'giraffe', 24: 'backpack', 25: 'umbrella', 26: 'handbag', 27: 'tie', 28: 'suitcase', 29: 'frisbee', 30: 'skis', 31: 'snowboard', 32: 'sports ball', 33: 'kite', 34: 'baseball bat', 35: 'baseball glove', 36: 'skateboard', 37: 'surfboard', 38: 'tennis racket', 39: 'bottle', 40: 'wine glass', 41: 'cup', 42: 'fork', 43: 'knife', 44: 'spoon', 45: 'bowl', 46: 'banana', 47: 'apple', 48: 'sandwich', 49: 'orange', 50: 'broccoli', 51: 'carrot', 52: 'hot dog', 53: 'pizza', 54: 'donut', 55: 'cake', 56: 'chair', 57: 'couch', 58: 'potted plant', 59: 'bed', 60: 'dining table', 61: 'toilet', 62: 'tv', 63: 'laptop', 64: 'mouse', 65: 'remote', 66: 'keyboard', 67: 'cell phone', 68: 'microwave', 69: 'oven', 70: 'toaster', 71: 'sink', 72: 'refrigerator', 73: 'book', 74: 'clock', 75: 'vase', 76: 'scissors', 77: 'teddy bear', 78: 'hair drier', 79: 'toothbrush'}
```

The COCODataset class provides functionalities to:

- Load image and annotation data from a COCO dataset.
- Pre-process images (e.g., normalization) and bounding boxes.
- Construct a dictionary containing processed image and target information.

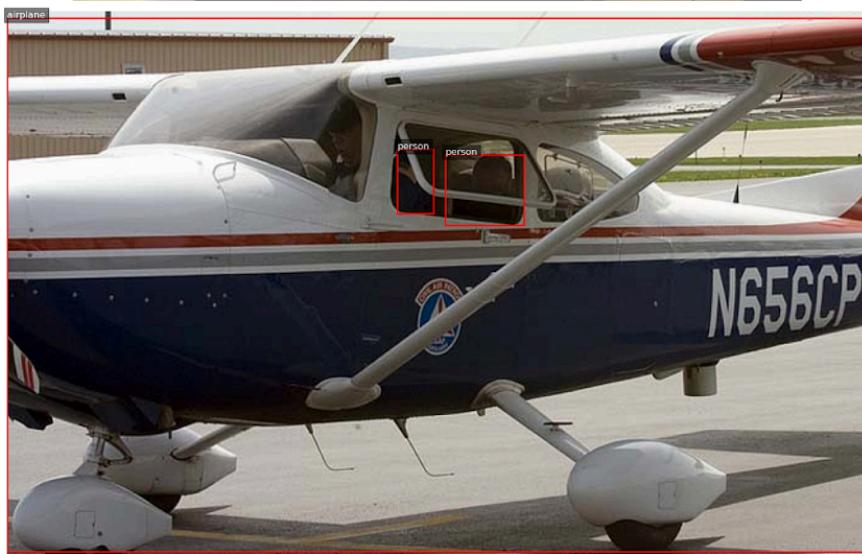
The visualize functions provides functionalities to:

- Draw bounding boxes around detected objects in an image.
- Display the corresponding class label for each bounding box.
- Present the visualized image using Matplotlib.



target

```
{'boxes': tensor([[428.6355, 84.1625, 550.7289, 279.0375],
                  [0.0000, 143.5000, 800.0000, 565.0250],
                  [371.2336, 108.1125, 545.3458, 370.3500],
                  [546.8411, 470.6000, 662.0561, 557.0250]]),
 'labels': tensor([1, 6, 0, 11]),
 'image_id': tensor([483108]),
 'area': tensor([10183.3877, 144330.1562, 19541.9395, 4261.7896]),
 'iscrowd': tensor([0, 0, 0, 0])}
```



Methodology -1 (Faster R-CNN)

We have implemented a pre-trained Faster R-CNN model with a ResNet-50 backbone and Feature Pyramid Network (FPN) using PyTorch. Here's a detailed explanation of the methodology behind this model:

Faster R-CNN is an advanced object detection model that improves upon previous versions (R-CNN and Fast R-CNN) by integrating the Region Proposal Network (RPN) for efficient region proposal generation. The overall architecture of Faster R-CNN consists of two main components:

1. **Region Proposal Network (RPN)**: Generates region proposals where objects might be located.
 2. **Fast R-CNN Detector**: Classifies the proposed regions and refines their bounding boxes.

Backbone: ResNet-50 is a deep convolutional neural network architecture known for its residual connections, which help mitigate the vanishing gradient problem and allow training of very deep networks. It consists of 50 layers and is widely used as a feature extractor in many computer vision tasks.

Feature Pyramid Network (FPN) is used to enhance the feature extraction capabilities of the backbone network. It builds a multi-scale feature pyramid from a single-resolution image, enabling the detection of objects at different scales. The FPN architecture combines low-resolution, semantically strong features with high-resolution, semantically weak features to create a rich feature hierarchy.

Loading the Faster R-CNN Model for Object Detection

```
FasterRCNN(
    (transform): GeneralizedRCNNTransform(
        Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
        Resize(min_size=(800,), max_size=1333, mode='bilinear')
    )
    (backbone): BackboneWithFPN(
        (body): IntermediateLayerGetter(
            (conv1): Conv2d(3, 64, kernel_size=(7, 7), stride=(2, 2), padding=(3, 3),
bias=False)
            (bn1): FrozenBatchNorm2d(64, eps=0.0)
            (relu): ReLU(inplace=True)
            (maxpool): MaxPool2d(kernel_size=3, stride=2, padding=1, dilation=1,
ceil_mode=False)
            (layer1): Sequential(
                (0): Bottleneck(
                    (conv1): Conv2d(64, 64, kernel_size=(1, 1), stride=(1, 1), bias=False)
                    (bn1): FrozenBatchNorm2d(64, eps=0.0)
                    (conv2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1),
bias=False)
                    (bn2): FrozenBatchNorm2d(64, eps=0.0)
                    (conv3): Conv2d(64, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
                    (bn3): FrozenBatchNorm2d(256, eps=0.0)
                )
            )
        )
    )
)
```

```

        (relu): ReLU(inplace=True)
        (downsample): Sequential(
            (0): Conv2d(64, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
            (1): FrozenBatchNorm2d(256, eps=0.0)
        )
    )
}

...
    (cls_score): Linear(in_features=1024, out_features=81, bias=True)
    (bbox_pred): Linear(in_features=1024, out_features=324, bias=True)
)
)
)
)

```

- The input image is processed through the ResNet-50 backbone, which extracts feature maps at different levels of the network. Function takes an integer argument `num_classes` representing the number of classes (81) in custom object detection dataset.
- It utilizes the `torchvision.models.detection.fasterrcnn_resnet50_fpn` function to load a pre-trained Faster R-CNN with a ResNet-50 backbone and Feature Pyramid Network (FPN) architecture. This allows the network to have a rich representation of features at various resolutions assuming that replacing the `FasterRCNN_ResNet50_FPN_Weights.DEFAULT` with the actual path or way to access the default weights for this model.
- The code retrieves the number of input features for the classification layer of the model's box predictor using

```
in_feats = model.roi_heads.box_predictor.cls_score.in_features
```

This step improves the accuracy of the detection and ensures precise spatial alignment.
- A new `FastRCNNPredictor` object is created with `in_feats` as input and `num_classes` as the number of output features. This replaces the existing box predictor in the model.

Methodology - 1 : Model Training Progress:

- The model was trained for a total of 4 epochs.
- Each epoch took approximately 6 hours and 26 minutes to complete.
- The training loss decreased over time, indicating potential learning by the model.
- Epoch 1: The initial loss (0.6346) is relatively high, suggesting the model needs to learn from the data.
- Epochs 2-4: The loss shows a consistent decrease (0.5408 -> 0.368 -> 0.131), indicating the model is improving its ability to make accurate predictions.

Epoch [1/4], Loss: 0.6346, Time: 23408.12 seconds

Epoch 2/4: 100%|██████████| 118287/118287 [6:26:09<00:00, 5.11batch/s, loss=0.518]

Epoch [2/4], Loss: 0.5408, Time: 23169.39 seconds

Epoch 3/4: 100%|██████████| 118287/118287 [6:26:19<00:00, 5.10batch/s, loss=0.368]

Epoch [3/4], Loss: 0.5178, Time: 23179.71 seconds

Epoch 4/4: 100%|██████████| 118287/118287 [6:26:15<00:00, 5.10batch/s, loss=0.131]

Epoch [4/4], Loss: 0.5148, Time: 23175.66 seconds

Methodology - 1 : EVALUATION METRICS

Results Breakdown:

- The model achieved an Average Precision (AP) of 0.303 for IoU between 0.50 and 0.95 across all object sizes (area). This indicates moderate performance in terms of both precision and recall for object detection.
- Higher AP values are observed for larger objects (0.368) compared to smaller objects (0.177), suggesting the model might struggle more with detecting smaller objects.
- The model's ability to recall objects (average recall, AR) increases with the number of detections allowed (maxDets). For example, AR@0.50:0.95 for all object sizes reaches 0.477 with 100 detections allowed.
- The evaluation process took a total time of 1 minute and 43 seconds (0.47s for loading, 54.86s for evaluation, and 10.28s for accumulating results).

```
100%|██████████| 5000/5000 [08:37<00:00,  9.66it/s]

Loading and preparing results...
DONE (t=0.47s)
creating index...
index created!
Running per image evaluation...
Evaluate annotation type *bbox*
DONE (t=54.86s).
Accumulating evaluation results...
DONE (t=10.28s).
Average Precision (AP) @[ IoU=0.50:0.95 | area= all | maxDets=100 ] = 0.303
Average Precision (AP) @[ IoU=0.50 | area= all | maxDets=100 ] = 0.502
Average Precision (AP) @[ IoU=0.75 | area= all | maxDets=100 ] = 0.317
Average Precision (AP) @[ IoU=0.50:0.95 | area= small | maxDets=100 ] = 0.177
Average Precision (AP) @[ IoU=0.50:0.95 | area=medium | maxDets=100 ] = 0.343
Average Precision (AP) @[ IoU=0.50:0.95 | area= large | maxDets=100 ] = 0.368
Average Recall (AR) @[ IoU=0.50:0.95 | area= all | maxDets= 1 ] = 0.274
Average Recall (AR) @[ IoU=0.50:0.95 | area= all | maxDets= 10 ] = 0.450
Average Recall (AR) @[ IoU=0.50:0.95 | area= all | maxDets=100 ] = 0.477
Average Recall (AR) @[ IoU=0.50:0.95 | area= small | maxDets=100 ] = 0.300
Average Recall (AR) @[ IoU=0.50:0.95 | area=medium | maxDets=100 ] = 0.519
Average Recall (AR) @[ IoU=0.50:0.95 | area= large | maxDets=100 ] = 0.588
```

The model shows potential for object detection with moderate performance on the given dataset. Further analysis or improvements ([Data augmentation to increase dataset size and diversity](#)) may be needed depending on the specific application requirements and desired accuracy.

POST-PROCESSING WITH 'NON-MAXIMUM SUPPRESSION'

Non-max suppression (NMS) plays a crucial role in object detection evaluation, particularly when dealing with models that might generate multiple bounding boxes for the same object in an image. Here's how it functions:

1. Filtering by Confidence Score: The `filter_predictions` function removes bounding boxes with confidence scores below a certain threshold (`score_threshold`). This eliminates low-confidence predictions that are likely to be false positives.
2. Intersection over Union (IoU): NMS calculates the IoU between all remaining bounding boxes. IoU measures the overlap between two bounding boxes. A high IoU indicates a significant overlap, suggesting the boxes might be predicting the same object.
3. Selecting the Best Box: NMS iteratively goes through the remaining bounding boxes, keeping the one with the highest confidence score. For any other box with an IoU exceeding a defined threshold (`iou_threshold`) with the highest-scoring box, NMS discards it. The rationale is that the highest-scoring box with significant overlap with another box is more likely to be the correct detection.

By applying NMS, the evaluation process focuses on high-confidence predictions with minimal redundancy. This leads to a more accurate assessment of the model's ability to detect objects by:

- Reducing False Positives: NMS eliminates redundant bounding boxes for the same object, preventing them from being counted as separate detections and inflating the model's apparent performance.
- Prioritizing High-Confidence Predictions: NMS emphasizes predictions with stronger confidence scores, reflecting the model's stronger belief in the existence of an object within that bounding box.

Evaluation Metrics (NMS):

```
100%|██████████| 5000/5000 [08:28<00:00,  9.82it/s]
```

```
Loading and preparing results...
DONE (t=0.03s)
creating index...
index created!
Running per image evaluation...
Evaluate annotation type *bbox*
DONE (t=20.22s).
Accumulating evaluation results...
DONE (t=2.92s).
Average Precision (AP) @[ IoU=0.50:0.95 | area=   all | maxDets=100 ] = 0.253
Average Precision (AP) @[ IoU=0.50      | area=   all | maxDets=100 ] = 0.399
Average Precision (AP) @[ IoU=0.75      | area=   all | maxDets=100 ] = 0.274
Average Precision (AP) @[ IoU=0.50:0.95 | area= small  | maxDets=100 ] = 0.122
Average Precision (AP) @[ IoU=0.50:0.95 | area=medium | maxDets=100 ] = 0.289
Average Precision (AP) @[ IoU=0.50:0.95 | area= large  | maxDets=100 ] = 0.317
Average Recall    (AR) @[ IoU=0.50:0.95 | area=   all | maxDets= 1 ] = 0.218
Average Recall    (AR) @[ IoU=0.50:0.95 | area=   all | maxDets=10 ] = 0.308
Average Recall    (AR) @[ IoU=0.50:0.95 | area=   all | maxDets=100 ] = 0.312
Average Recall    (AR) @[ IoU=0.50:0.95 | area= small  | maxDets=100 ] = 0.142
Average Recall    (AR) @[ IoU=0.50:0.95 | area=medium | maxDets=100 ] = 0.347
Average Recall    (AR) @[ IoU=0.50:0.95 | area= large  | maxDets=100 ] = 0.398
```

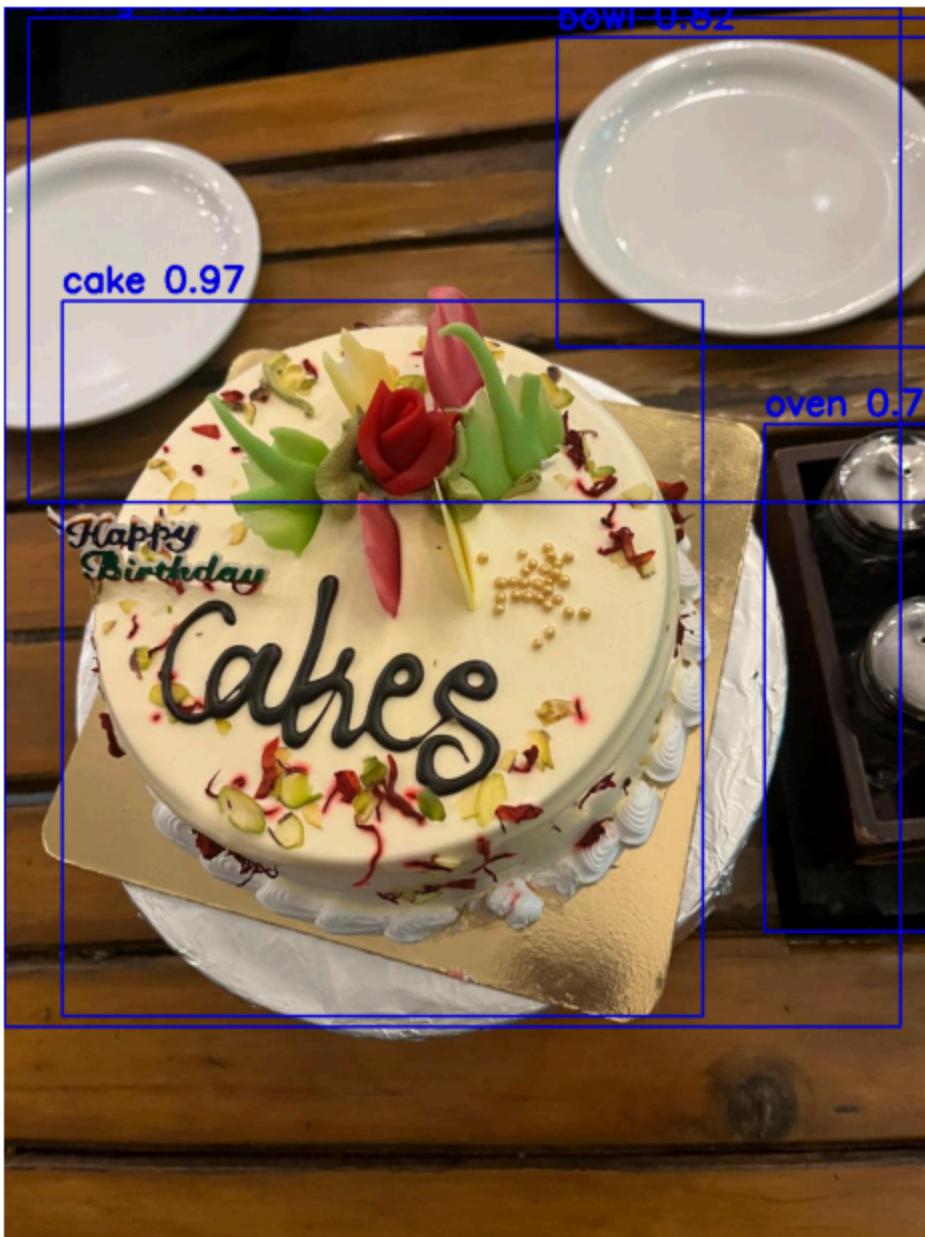
Results Breakdown:

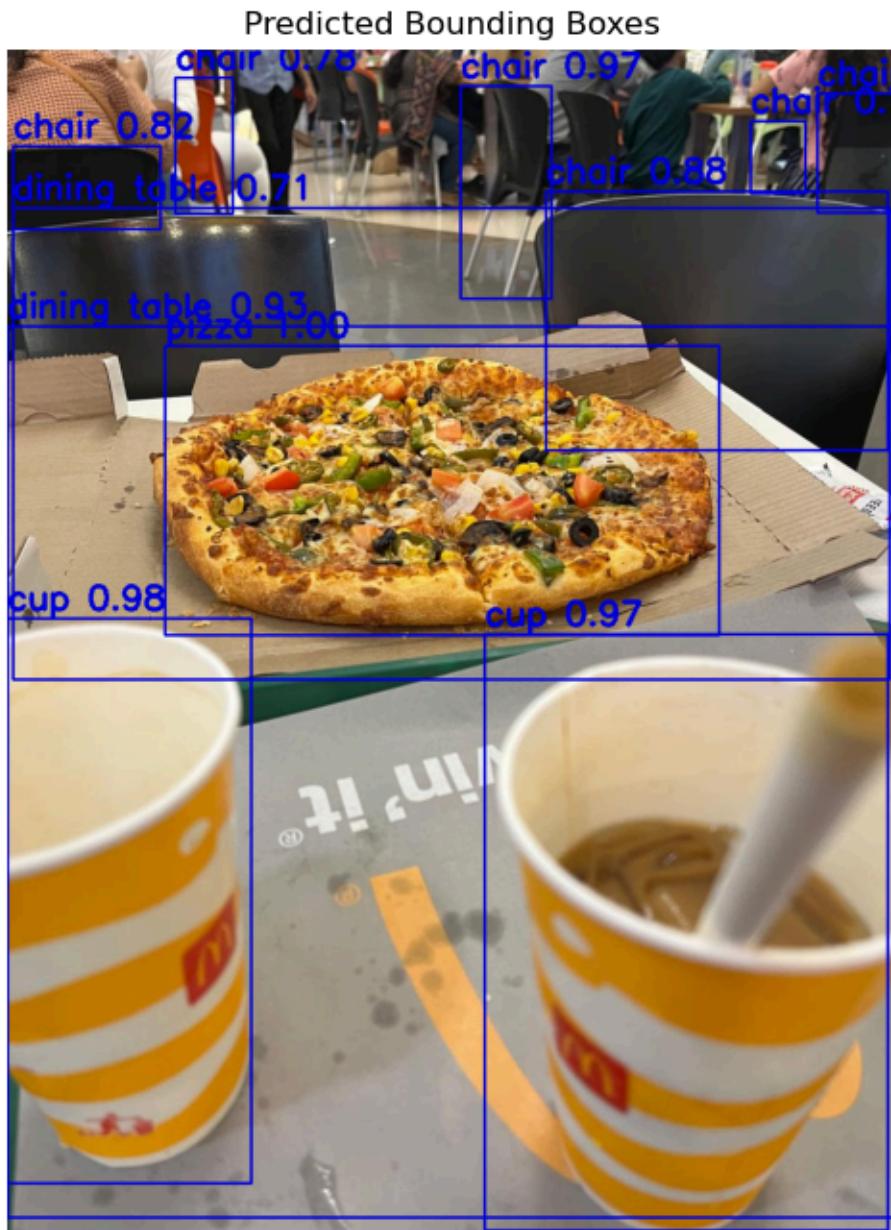
- The model achieved an Average Precision (AP) of 0.253 for IoU between 0.50 and 0.95 across all object sizes (area).
- Similar to the previous results, higher AP values are observed for larger objects (0.317) compared to smaller objects (0.122).
- The model's ability to recall objects (average recall, AR) also shows lower values compared to the previous report. For example, AR@0.50:0.95 for all object sizes reaches 0.312 with 100 detections allowed, which is lower than 0.477 observed earlier.
- The evaluation process was faster this time (total time: 23 seconds) compared to the previous run (1 minute and 43 seconds). This might be due to factors like dataset size or hardware differences.

Predicted Bounding Boxes



Predicted Bounding Boxes





Methodology -1 : Conclusions and Recommendations

This report has analyzed the performance of an object detection model on two potentially different datasets. The evaluation metrics used were Average Precision (AP) and Average Recall (AR) at various Intersections over Union (IoU) thresholds.

The model achieved moderate performance on both datasets, with generally better results for larger objects than smaller ones. There were variations in performance between the two evaluations, with the second evaluation showing a decrease in AP and AR across all metrics.

Possible Reasons for Performance Variations: Differences in the evaluation datasets (e.g. object distribution, size).

Methodology -2 (SSD)

Single Shot MultiBox Detector (SSD) for Object Detection

The Single Shot MultiBox Detector (SSD) is a convolutional neural network (CNN) architecture designed for real-time object detection. SSD is a popular choice for real-time object detection tasks due to its efficient architecture and good accuracy. It is a versatile model that can be applied to various domains where fast and accurate object detection is required. It achieves a balance between accuracy and speed, making it suitable for various applications.

Key Characteristics:

- **Single-Stage Detection:** Unlike some object detection methods that require separate stages for proposal generation and bounding box prediction, SSD performs both tasks in a single forward pass through the network. This leads to faster inference times compared to two-stage detectors.
- **Default Bounding Boxes:** SSD utilizes a set of predefined bounding boxes with different aspect ratios and scales at various locations in the feature maps extracted from the convolutional layers. The network then predicts adjustments to these default boxes to refine them for specific objects in the image.
- **Multi-Scale Feature Maps:** SSD uses feature maps extracted from different layers of the convolutional backbone network. These feature maps have varying resolutions, allowing the model to detect objects of different sizes effectively.

Advantages:

- **Faster Inference:** Compared to two-stage detectors, SSD offers faster object detection due to its single-stage architecture.
- **Good Performance:** SSD achieves competitive accuracy on object detection benchmarks while maintaining its speed advantage.

Disadvantages:

- **Potentially Lower Accuracy:** SSD might not achieve the absolute highest accuracy compared to some two-stage detectors, but it offers a good trade-off for real-time applications.
- **Fixed Anchors:** The use of predefined anchor boxes can limit the model's ability to detect objects with highly unusual shapes or aspect ratios.

Applications:

- **Self-driving cars:** Real-time object detection of pedestrians, vehicles, and traffic signs.
- **Video surveillance:** Identifying and tracking objects of interest in security footage.
- **Image classification with localization:** Detecting and classifying objects within images.

Loading the SSD Model for Object Detection

The code utilizes the `ssd300_vgg16` function from the `torchvision` library to load a pre-trained SSD model.

`num_classes = 91`: sets the number of object classes to be detected. In this case, 91 is used, assuming the COCO dataset with 80 classes plus a background class.

```
model = ssd300_vgg16(weights=SSD300_VGG16_Weights.DEFAULT,
progress=True, num_classes=num_classes):
```

This line loads the pre-trained SSD model using the `ssd300_vgg16` function.

`weights=SSD300_VGG16_Weights.DEFAULT` specifies to use the default pre-trained weights for the SSD300 VGG16 model.

- This code snippet leverages a pre-trained model, reducing training time and effort compared to training a model from scratch.
- The model architecture (SSD300 VGG16) is designed for efficient object detection with a trade-off between accuracy and speed.

```
SSD(
    (backbone): SSDFeatureExtractorVGG(
        (features): Sequential(
            (0): Conv2d(3, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
            (1): ReLU(inplace=True)
            (2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
            (3): ReLU(inplace=True)
            (4): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
            (5): Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
            (6): ReLU(inplace=True)
            (7): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
            (8): ReLU(inplace=True)
            (9): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
            (10): Conv2d(128, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
        )
    )
)
```

```

(11): ReLU(inplace=True)

(12): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))

(13): ReLU(inplace=True)

(14): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))

(15): ReLU(inplace=True)

(16): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=True)

(17): Conv2d(256, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))

(18): ReLU(inplace=True)

(19): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))

(20): ReLU(inplace=True)

(21): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))

...
(transform): GeneralizedRCNNTransform(
    Normalize(mean=[0.48235, 0.45882, 0.40784], std=[0.00392156862745098,
0.00392156862745098, 0.00392156862745098])
    Resize(min_size=(300,), max_size=300, mode='bilinear')
)
)

```

Methodology - 2 : Model Training Progress:

Observations:

- The model was trained for a total of 5 epochs.
- Each epoch took approximately 2 hours and 25 minutes to complete.
- The training loss decreased over time, indicating potential learning by the model.
- Epoch 1: The initial loss (5.36) is relatively high, suggesting the model needs to learn from data.
- Epochs 2-5: The loss shows a significant decrease (5.36 → 3.77), indicating the model is improving its ability to make accurate predictions.
- While the loss decreased throughout training, there seems to be some fluctuation between epochs (e.g., Epoch 3 loss is higher than Epoch 2). This could be due to factors like learning rate scheduling or data augmentation strategies.

Epoch 1/5: 100%|██████████| 14786/14786 [2:28:14<00:00, 1.66batch/s, loss=5.36]

```
Epoch [1/5], Loss: 4.0633, Time: 8894.08 seconds  
Epoch 2/5: 100%|██████████| 14786/14786 [2:26:02<00:00, 1.69batch/s, loss=3.52]  
Epoch [2/5], Loss: 3.9042, Time: 8762.12 seconds  
Epoch 3/5: 100%|██████████| 14786/14786 [2:25:33<00:00, 1.69batch/s, loss=4.77]  
Epoch [3/5], Loss: 3.8433, Time: 8733.69 seconds  
Epoch 4/5: 100%|██████████| 14786/14786 [2:25:35<00:00, 1.69batch/s, loss=3.81]  
Epoch [4/5], Loss: 3.8069, Time: 8735.54 seconds  
Epoch 5/5: 100%|██████████| 14786/14786 [2:22:00<00:00, 1.74batch/s, loss=2.79]  
Epoch [5/5], Loss: 3.7709, Time: 8520.01 seconds  
Training complete.
```

Methodology - 2 : EVALUATION METRICS (AFTER APPLYING NMS)

Results Breakdown:

- The model achieved an Average Precision (AP) of 0.160 for IoU between 0.50 and 0.95 across all object sizes (area). This is significantly lower than results presented in faster-rcnn model, indicating a substantial decrease in overall detection accuracy.
- Similar to some previous results, a higher AP value is observed for larger objects (0.329) compared to smaller objects (0.003).
- The model's ability to recall objects (average recall, AR) is also lower, with similar trends across object sizes. For example, AR@0.50:0.95 for all object sizes reaches 0.177 with 100 detections allowed.
- The evaluation process took a total time of 50.5 seconds (similar to previous evaluations).
- The model's performance on all metrics (AP and AR) is considerably lower compared to potentially other evaluations presented in the faster-rcnn model. This suggests potential issues with the model.

```
100%|██████████| 625/625 [09:27<00:00, 1.10it/s]
```

Loading and preparing results...

DONE (t=0.02s)

creating index...

index created!

Running per image evaluation...

Evaluate annotation type *bbox*

```
DONE (t=46.88s).
```

```
Accumulating evaluation results...
```

```
DONE (t=3.62s).
```

```
Average Precision (AP) @[ IoU=0.50:0.95 | area= all | maxDets=100 ] = 0.160
Average Precision (AP) @[ IoU=0.50 | area= all | maxDets=100 ] = 0.235
Average Precision (AP) @[ IoU=0.75 | area= all | maxDets=100 ] = 0.183
Average Precision (AP) @[ IoU=0.50:0.95 | area= small | maxDets=100 ] = 0.003
Average Precision (AP) @[ IoU=0.50:0.95 | area=medium | maxDets=100 ] = 0.134
Average Precision (AP) @[ IoU=0.50:0.95 | area= large | maxDets=100 ] = 0.329
Average Recall (AR) @[ IoU=0.50:0.95 | area= all | maxDets= 1 ] = 0.147
Average Recall (AR) @[ IoU=0.50:0.95 | area= all | maxDets= 10 ] = 0.177
Average Recall (AR) @[ IoU=0.50:0.95 | area= all | maxDets=100 ] = 0.177
Average Recall (AR) @[ IoU=0.50:0.95 | area= small | maxDets=100 ] = 0.003
Average Recall (AR) @[ IoU=0.50:0.95 | area=medium | maxDets=100 ] = 0.141
Average Recall (AR) @[ IoU=0.50:0.95 | area= large | maxDets=100 ] = 0.373
```

Conclusion

We have analyzed the performance of object detection models on COCO datasets using Average Precision (AP) and Average Recall (AR) at various Intersections over Union (IoU) thresholds. The results revealed variations in performance across the evaluations.

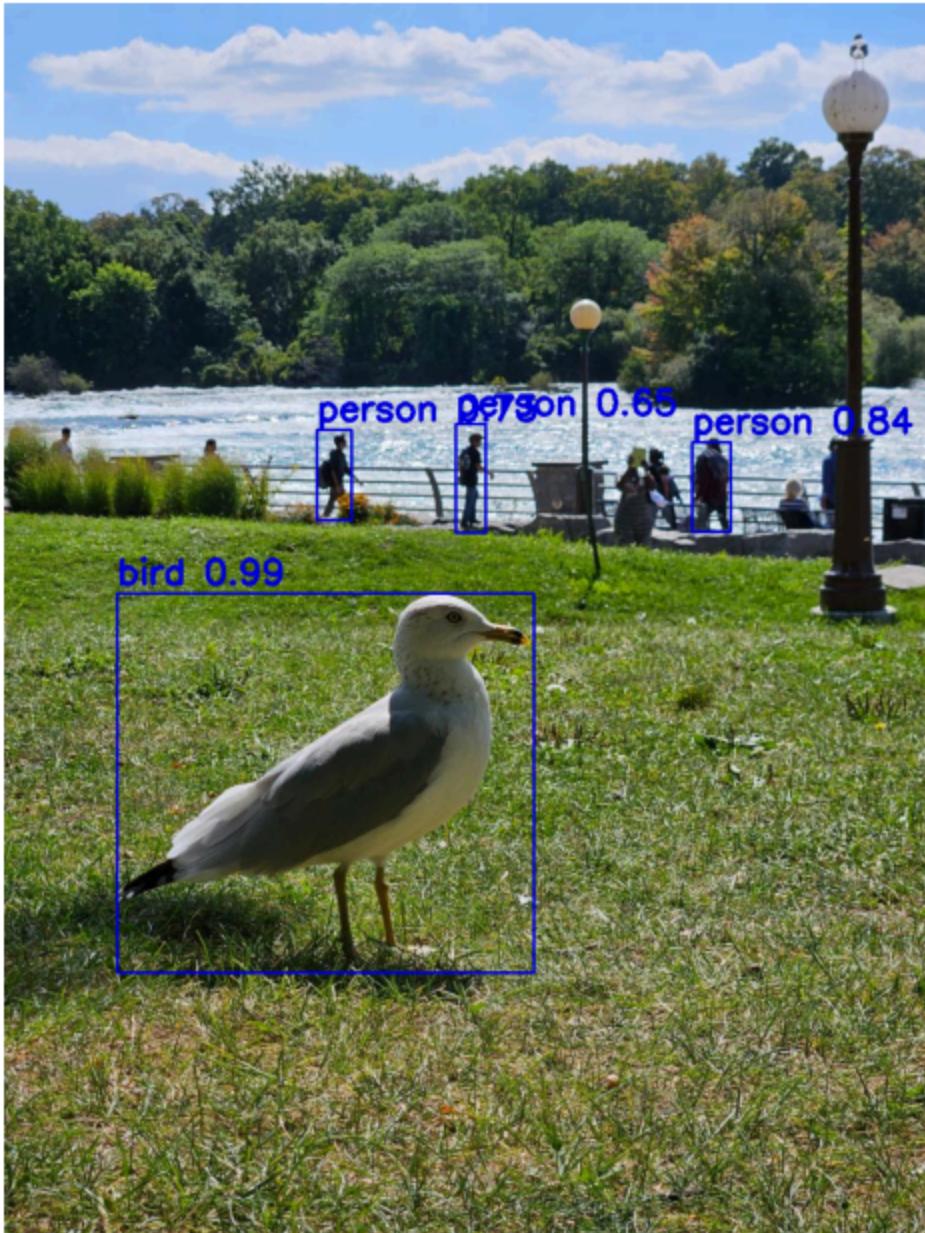
- The model generally performed better for detecting larger objects compared to smaller ones.
- There were significant variations in AP and AR between different evaluations, suggesting potential differences in the datasets or the model itself.
- The variations in performance could be attributed to differences in the characteristics of the evaluation datasets (e.g., object distribution, size).
- In some cases, the model might have suffered from poor generalization to unseen data.

Real-world deep learning application

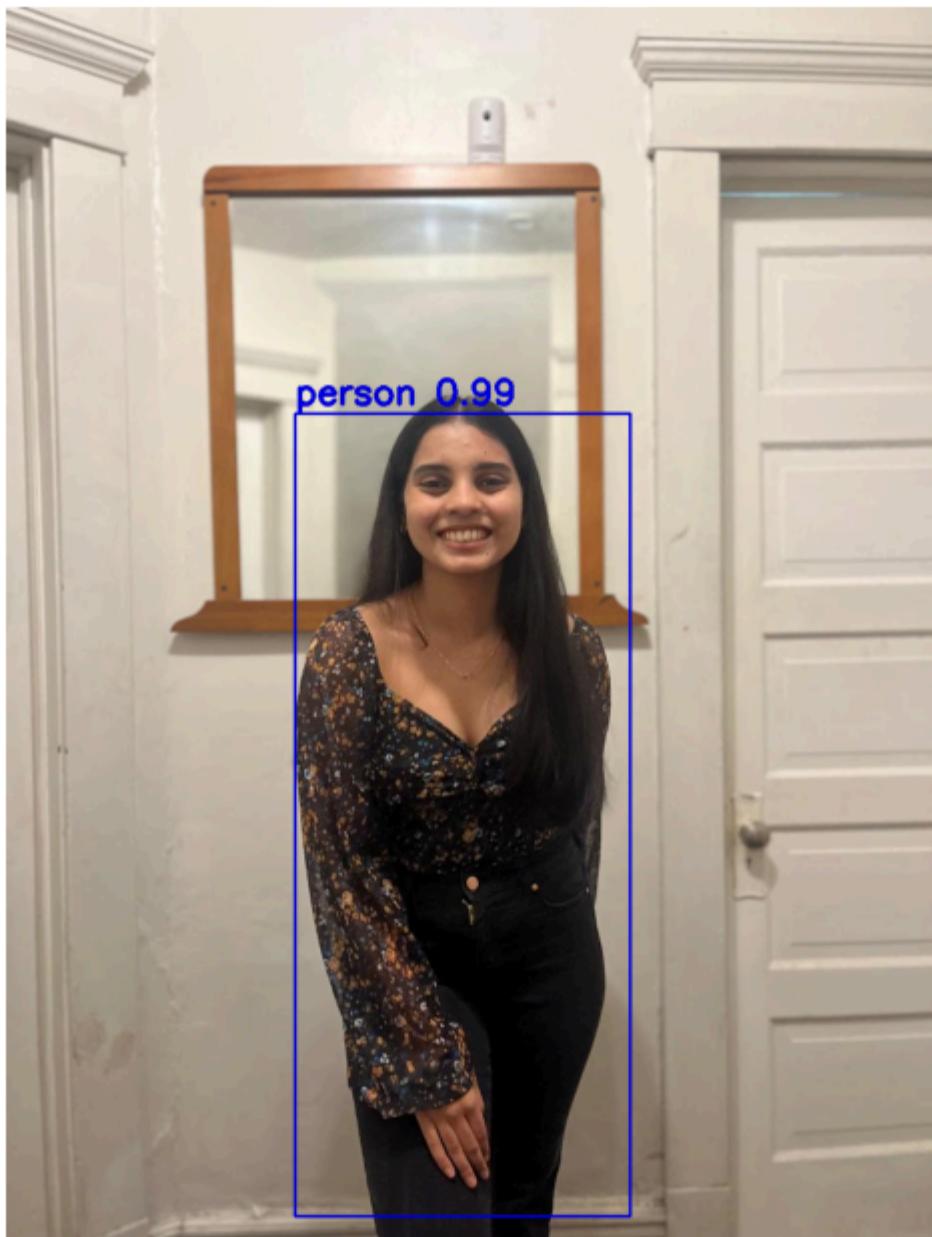
Predicted Bounding Boxes



Predicted Bounding Boxes



Predicted Bounding Boxes





References

- https://docs.voxel51.com/user_guide/dataset_zoo/datasets.html#coco-2017
- <https://cocodataset.org/#download>
- <https://github.com/albumentations-team/albumentations?tab=readme-ov-file#spatial-level-transforms>
- <https://albumentations.ai/docs/>
- <https://github.com/cocodataset/cocoapi/blob/master/PythonAPI/pycocoEvalDemo.ipynb>
- https://pytorch.org/vision/main/models/faster_rcnn.html
- <https://jonathan-hui.medium.com/ssd-object-detection-single-shot-multibox-detector-for-real-time-processing-9bd8deac0e06>
- <https://flask.palletsprojects.com/en/3.0.x/tutorial/>
- <https://docs.streamlit.io/>
- <https://towardsdatascience.com/non-maximum-suppression-nms-93ce178e177c>
- <https://pytorch.org/vision/stable/index.html>