

Object Detection

Alan Casey and Alex Racapé

October 2023

1 Introduction

In this lab, we tested the Faster R-CNN and YOLO models for Object Detection. This is a central problem in computer vision that is concerned with not only discerning the class of an object, but also discerning where the objects are located in an image. The primary objective of our experiment was to both qualitatively and quantitatively evaluate the detection performance of the Faster R-CNN and YOLO models, which are two prominent neural network architectures for this task. Leveraging these popular neural network models, we conducted studies to evaluate their overall performance based on their object and boundary box classifications, as well as their total inference times.

2 Methodology

Our studies aimed to investigate the performance of the Faster R-CNN and YOLO models on a curated set of images and a larger testing dataset. We used the official PyTorch release of Faster R-CNN with a Resnet backbone and PyTorch hub's implementation of YOLOv5.

To initiate our study, we selected five images showcasing a variety of scenes and objects to assess potential variations between the two models. The images we chose can be found in the appendix.

There are a few pictures with cats and people as the subjects, and we chose an image of a shed with numerous objects. We also chose an image with some cats and cars mixed in. The 5 images contain a range of scene and object complexities, and we used them to make qualitative observations about the models and their different results. We ran all 5 images on both models, and plotted the images with the resulting bounding boxes and class predictions.

Afterwards, moving beyond qualitative observations, we shifted towards a quantitative analysis, using a kaggle data set with images of buses and trucks to test detection accuracy and inference time. The dataset includes images from everyday scenes with buses in trucks in the foreground and background. To simplify our analysis, we focused on a smaller random subset of the data. We randomly selected 50 images that included only one labeled bus or truck. We made this decision to simplify our code, so we would only need to check our

predictions against a single ground truth. This analysis of both Faster R-CNN and YOLOv5 models provides us with some numeric data to compare with our previous qualitative judgements.

We measured the models’ detection accuracy using the models’ classification confidence levels and IoU (Intersection over Union). We established different threshold values for confidence levels (30% confidence) as well as for IoU (40%), which determined whether an inference was successful or not. We chose these thresholds based on observations from our first experiment where most correct observations were greater than 30% confidence, and detections below this were mostly wrong. In our analysis, we looked at the detections outputted by the model that were greater than our confidence threshold. Then we calculated the accuracy by looking at how many of these confident detections were adequately overlapping with a truck’s true bounding box. A detection could be wrong if there was no bus in the scene or the bounding box was too far from our true bounding box (IoU below our threshold). This logic can be found in the Lines 189 and 276 of our code in the appendix. In order to calculate the IoU between two bounding boxes, we wrote a function called *compute-IoU* (Line 123) that computes the overlap between two bounding boxes—specifically the training box and the predicted box—to determine the accuracy of the object’s bounding box. We measured the inference time by adding a timer to our code. After we calculated both the confidence levels and IoU for all images (using both models), as well as the inference time, we proceeded to compute the overall detection accuracy, and compared the results for both models in a following table. The quantitative analysis code can be found on line 148 for the Faster-RCNN model and 241 for the YoloV5 model.

3 Results and Discussion

3.1 Qualitative Observations

For our qualitative analysis, we observed the results of running 5 of our images through both models. We saw distinct characteristics for these two models. The Faster R-CNN model outputted many more proposed detections, and the model appeared to be more confident in its predictions, often producing predictions with higher confidence scores. One reason for this could be that the Faster R-CNN model has a lower confidence threshold for its predictions than YOLOv5, or maybe the region proposal network is able to find better proposals that lead to higher confidences. Faster R-CNN displayed a wider range of accurate and erroneous predictions. Moreover, we observed a considerable degree of overlapping boundary boxes in the results generated by Faster R-CNN.

This confidence from Faster R-CNN, however, did not necessarily correlate with the accuracy of the predictions. For example, in the image of Alex’s shed, it was 99% sure that a sander was a clock. All of our results can be seen in the Appendix. In contrast, YOLOv5 exhibited a more reserved approach, generating fewer bounding boxes and predictions with lower confidence scores,



leading to fewer false positives. However, in one of the images with the cats lying on the clay roof it did not make any prediction at all, and it misclassified cats as dogs as illustrated in the figure.

Furthermore, it is noteworthy that both models better able to detect objects on the road, especially with respect to cars and pedestrians. These detection may be explained by the fact that both Faster R-CNN and YOLOv5 have been designed and optimized for tasks closely related to autonomous driving, making them particularly adept at recognizing objects typically encountered in such scenarios.

3.2 Quantitative Observations

Table 1: Model Performance Comparison

Model	Detection Accuracy	Inference Time
Faster R-CNN	75.0%	6m 39s
YOLO	87.0%	18s

In our experiment with bus images, the YOLO model outperformed faster r-cnn in both speed and detection accuracy. We expected YOLO to be faster because of its fully convolutional architecture, but we were surprised by the magnitude of this difference. One contributing factor could be that we had to look through more proposed detections from the faster r-cnn model when calculating detection accuracy. Additionally, there are some differences in the way data is loaded for each models' input, though each image is read in both cases. In terms of detection accuracy, YOLO performed better once again with an accuracy of 87%. For each model we displayed the images that were incorrectly detected, and they were understandably difficult. Both models struggled with images of unique looking busses, and both had instances where it labeled a truck as a bus. There were two particular images that led to failures in both models, and they included a bus-like car in the background and a large obstructing object with a bus peeking out from behind.

In some research online we found that faster r-cnn is considered by some to be slower but more accurate. Perhaps this is because of its convolutional method for proposing regions of interest that is less restrictive than the grid used by YOLO. This reasoning fits with our qualitative findings where YOLO seemed less accurate with less confident detections, but when calculating detection accuracy for the buses, faster r-cnn was actually less accurate. While it is unclear why, this could be due to the nature of the dataset since buses are generally large rectangular objects that might be better suited for YOLO.

4 Conclusion

In this study, we conducted a comprehensive evaluation of the Faster R-CNN and YOLO object detection models, both qualitatively and quantitatively. Object detection is a pivotal task in computer vision, as it involves not only recognizing object classes but also accurately localizing objects within images. Our qualitative analysis of five diverse images revealed distinct characteristics of these models. Faster R-CNN demonstrated a high prediction rate and a tendency to provide a lot of bounding boxes with high confidence scores. However, this confidence did not always translate to accurate predictions, leading to a considerable number of false positives and overlapping bounding boxes. In contrast, YOLOv5 exhibited a more cautious approach with fewer bounding boxes and lower confidence scores, resulting in a higher precision rate.

Our quantitative results revealed that YOLOv5 significantly outperformed Faster R-CNN in terms of both speed and detection accuracy. YOLOv5's superior speed can be primarily attributed to its architectural design, which avoids the computationally intensive fully connected layers used by Faster R-CNN. This advantage is particularly pronounced when dealing with a substantial number of proposed detections, highlighting YOLOv5's efficiency in processing images.

Overall, our study underscores the significance of both qualitative and quantitative assessments when selecting object detection models for specific applications. Furthermore, the results highlight the importance of considering the

intended purpose and training strategies of these models in understanding their performance variations, which may vary depending on the specific requirements of the task at hand.

5 Appendix

```
1  # -*- coding: utf-8 -*-
2  """lab5.ipynb
3
4  Automatically generated by Colaboratory.
5
6  Original file is located at
7      https://colab.research.google.com/drive/11fMwQaswak4hwoPZji-eUxGjmZu5NhXf
8
9  > Indented block Yooooo
10 """
11
12 import pandas as pd
13 import torch
14 import torch.nn as nn
15 from torch.optim import Adam
16 from torchsummary import summary
17 import numpy as np
18 import time
19 import os
20 import random
21 from google.colab import drive
22
23 # Get Faster R-CNN models
24 from torchvision.models.detection import fasterrcnn_resnet50_fpn, FasterRCNN_ResNet50_FPN_Weights
25
26 # Imports for Data
27 from torch.utils.data import Dataset, DataLoader
28 from torchvision.transforms import transforms, Resize, ToTensor
29 from torchvision.transforms.functional import to_pil_image, convert_image_dtype
30 from torchvision.io import read_image
31 from torchvision.utils import draw_bounding_boxes
32
33 # Set up drive storage and device
34 drive.mount('/content/drive')
35 device = 'cuda' if torch.cuda.is_available() else 'cpu'
36 print(device)
37
38 # Get models set up
39 classes = FasterRCNN_ResNet50_FPN_Weights.DEFAULT.meta["categories"]
40 fast_rcnn = fasterrcnn_resnet50_fpn(weights=FasterRCNN_ResNet50_FPN_Weights.DEFAULT)
41 fast_rcnn.eval()
```

```

42
43 yolo = torch.hub.load('ultralytics/yolov5', 'yolov5s', pretrained=True)
44
45 # Experiment 1: Testing models on random images
46
47 # Images
48 imgs = ["alan.jpg", "cats2.jpg", "cats.jpg", "shed.jpg", "market.jpg"]
49 img_tensors = []
50 for path in imgs:
51     image = read_image(path)
52     image = image / 255.0 # Normalise the image to [0, 1]
53     img_tensors.append(image)
54
55 # Inference
56 fast_rcnn_results = fast_rcnn(img_tensors)
57 yolo_results = yolo(imgs)
58
59 # Display Results with bounding boxes and captions
60 def show_bb(img, boxes, captions):
61     img = (img * 255).to(torch.uint8)
62     img = draw_bounding_boxes(img, boxes, captions, width=1)
63     img = img.detach()
64     img = to_pil_image(img)
65     display(img)
66
67 # Display the results for each image from faster r-cnn
68 yolo_results.print()
69 for i, result in enumerate(fast_rcnn_results):
70
71     # Get captions, labels, confidences
72     confs = result["scores"]
73     labels = result["labels"]
74     captions = []
75     for label, conf in zip(labels, confs):
76         captions.append(f"{classes[label]} {(conf * 100).round()}")
77
78     # Display
79     boxes = result["boxes"]
80     show_bb(img_tensors[i], boxes, captions)
81
82 # Display the YOLO results
83 yolo_results.show()
84
85 # Experiment 2: Testing models on bus dataset
86
87 # Get Data for R-CNN
88 class BusDataset(Dataset):
89     def __init__(self, data_root):
90         self.root = data_root
91         self.image_root = data_root + "/images/images"

```

```

92
93     # Read CSV data from path
94     data = pd.read_csv(data_root + "/subset.csv")
95
96     # Set up x and y from CSV stuff
97     x = data["ImageID"]
98     y = data[["LabelName", "XMin", "XMax", "YMin", "YMax"]]
99     self.x, self.y = x, y
100
101     def __getitem__(self, ix):
102
103         # Retrieve image using ID
104         img_id = self.x.iloc[ix]
105         path = self.image_root + f"/{img_id}.jpg"
106         img = read_image(path)
107         img = img / 255.0
108         label = self.y.iloc[ix].values
109         return img, label
110
111     def __len__(self):
112         return len(self.x)
113
114 test_dataset = BusDataset("drive/MyDrive/Data/bus_data")
115
116 def area_from_points(min_p, max_p):
117     """Calculates the area of a square defined by its min and max point"""
118     width = max_p[0] - min_p[0]
119     height = max_p[1] - min_p[1]
120     return width * height
121
122
123 def compute_IoU(box1, box2):
124     """Computes Intersection over Union for two bounding boxes"""
125
126     min1, max1 = (box1[0], box1[2]), (box1[1], box1[3])
127     min2, max2 = (box2[0], box2[2]), (box2[1], box2[3])
128
129     i_min = max(min1, min2)
130     i_max = min(max1, max2)
131     if i_min[0] > i_max[0] or i_min[1] > i_max[1]:
132         return 0
133     intersection = area_from_points(i_min, i_max)
134
135     area1 = area_from_points(min1, max1)
136     area2 = area_from_points(min2, max2)
137     union = area1 + area2 - intersection
138     return intersection / union
139
140 def convert_boxes(boxes):
141     """Converts bounding boxes from xxyy to xyxy format"""

```

```

142     new_boxes = []
143     for box in boxes:
144         new_box = [box[0], box[2], box[1], box[3]]
145         new_boxes.append(new_box)
146     return new_boxes
147
148 def evaluate_fast_rcnn(model, dataset):
149     correct = 0
150     total = 0
151     confidence_thresh = .3
152     iou_thresh = .4
153
154     # Set the model to evaluation mode
155     model.eval()
156
157     # Disable gradient calculation
158     start = time.time()
159     for i in range(50):
160         print(f"STARTING IMAGE {i}...")
161
162         # Get image, label, and prediction
163         img, label = dataset[i]
164         prediction = model([img])[0]
165         true_class = label[0]
166         true_box = label[1:5]
167         height, width = img.size()[1], img.size()[2]
168
169         # Find which predicted boxes are buses, with conf > thresh
170         bus_indexes = []
171         pred_classes = [classes[label] for label in prediction["labels"]]
172         for j in range(len(pred_classes)):
173             if pred_classes[j] == "bus" and prediction["scores"][j] > confidence_thresh:
174                 bus_indexes.append(j)
175
176         # Adjust predicted bounding box to size of image
177         true_box = [val * width if idx < 2 else val * height for idx, val in enumerate(true_box)]
178
179         # For each confident detection, check if it was correct
180         for bus_index in bus_indexes:
181
182             # Get our predicted bounding box for a bus
183             predict_box = prediction["boxes"][bus_index].tolist()
184             predict_box = convert_boxes([predict_box])[0]
185
186             # Compute predicted IoU with true bounding box
187             print(true_box, predict_box)
188             iou = compute_IoU(true_box, predict_box)
189             if iou > iou_thresh and true_class == "Bus":
190                 correct += 1
191

```



```

192         # Display our wrong guess, if there is no bus
193         elif true_class != "Bus":
194             boxes = convert_boxes([predict_box])
195             show_bb(img, torch.tensor(boxes), ["Predicted Bus"])
196
197         # Display the bad guess, if we are confident and missed
198         elif true_class == "Bus" and iou < iou_thresh:
199             print([true_box, predict_box])
200             boxes = convert_boxes([true_box, predict_box])
201             print(boxes)
202             show_bb(img, torch.tensor(boxes), ["True Bus", "Predicted Bus"])
203
204         print(f"Predicted Box: {predict_box}, IoU: {iou}")
205         total += 1
206
207     # Calculate accuracy
208     accuracy = (correct / total) * 100.0
209     testing_time = time.time() - start
210     return accuracy, testing_time
211
212 # Run through our tests
213 evaluate_fast_rcnn(fast_rcnn, test_dataset)
214
215 class YoloDataset(Dataset):
216     def __init__(self, data_root):
217         self.root = data_root
218         self.image_root = data_root + "/images/images"
219
220         # Read CSV data from path
221         data = pd.read_csv(data_root + "/subset.csv")
222
223         # Set up x and y from CSV stuff
224         x = data["ImageID"]
225         y = data[["LabelName", "XMin", "XMax", "YMin", "YMax"]]
226         self.x, self.y = x, y
227
228     def __getitem__(self, ix):
229
230         # Retrieve image using ID
231         img_id = self.x.iloc[ix]
232         path = self.image_root + f"/{img_id}.jpg"
233         img_size = read_image(path).size()
234         label = self.y.iloc[ix].values
235         return path, (img_size[1], img_size[2]), label
236
237     def __len__(self):
238         return len(self.x)
239
240
241 def evaluate_yoloV5(model, dataset):

```

```

242
243     accurate_predictions = 0
244     total_bus_predictions = 0
245     confidence_thresh = 0.3
246     iou_thresh = 0.4
247
248     start = time.time()
249     for i in range(50):
250
251         img, size, true_label = dataset[i]
252         height, width = size
253
254         # Get the true class and true box dimensions
255         true_class = true_label[0]
256         true_box = true_label[1:5]
257         # Adjust predicted bounding box to size of image
258         true_box = [val * width if idx < 2 else val * height for idx, val in enumerate(true_box)]
259
260         # Compute and store the results for the yolo model
261         results = model(img)
262         label_table = results.pandas().xyxy[0]
263
264         """
265         Count the total number of true predictions
266         A true prediction consists of:
267             1) A confident "bus" when the image contains a bus, with
268             2) An accurate bounding box
269         """
270         for index, row in label_table.iterrows():
271             # Only loop through the predictions that are a confident "bus"
272             if row["name"] == "bus" and row["confidence"] > confidence_thresh:
273                 total_bus_predictions += 1
274                 pred_box = [row["xmin"], row["xmax"], row["ymin"], row["ymax"]]
275                 iou = compute_IoU(pred_box, true_box)
276
277                 if iou > iou_thresh and true_class == "Bus":
278                     # Increase accuracy if IoU passes our threshold
279                     accurate_predictions += 1
280                 elif true_class != "Bus":
281                     # Display the incorrect bus prediction box
282                     draw_box = convert_boxes([pred_box])
283                     img = read_image(img)
284                     show_bb(img, torch.tensor(draw_box), ["Predicted Bus"])
285                 else:
286                     # Display the true bus bounding box and our incorrect bus prediction box
287                     draw_boxes = convert_boxes([true_box, pred_box])
288                     img = read_image(img)
289                     show_bb(img, torch.tensor(draw_boxes), ["True Bus", "Predicted Bus"])
290
291

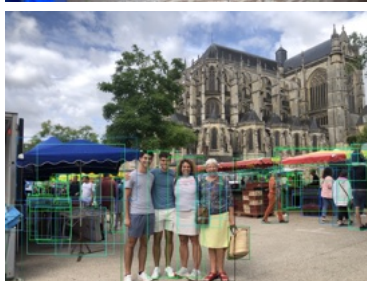
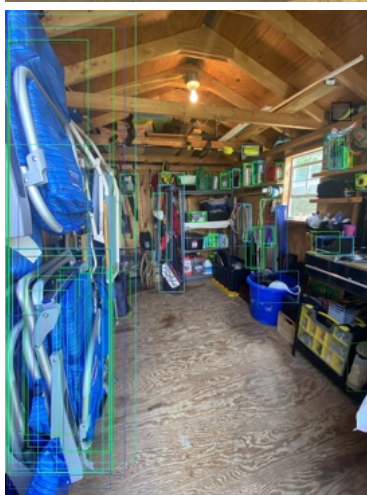
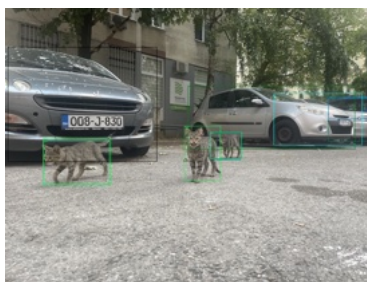
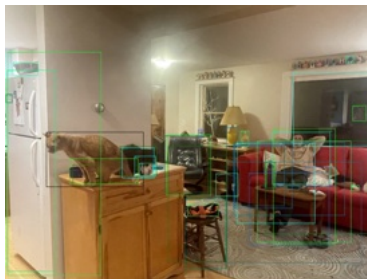
```

```

292
293     accuracy = (accurate_predictions / total_bus_predictions) * 100.0
294     testing_time = time.time() - start
295     return accuracy, testing_time
296
297     # What the results format should look like:
298     #      xmin      ymin      xmax      ymax  confidence  class   name
299     # 0  749.50   43.50  1148.0   704.5    0.874023      0  person
300     # 2  114.75  195.75  1095.0   708.0    0.624512      0  person
301     # 3  986.00  304.00  1028.0   420.0    0.286865     27   tie
302
303 yolo_dataset = YoloDataset("drive/MyDrive/Data/bus_data")
304 evaluate_yoloV5(yolo, yolo_dataset)
305
306 # Create a list of images we were able to download
307 downloaded_files = os.listdir("drive/MyDrive/Data/bus_data/images/images")
308
309 # Create a new DataFrame with only the downloaded images (12,308)
310 original_df = pd.read_csv("drive/MyDrive/Data/bus_data/df.csv")
311 original_df['ImageID'] = original_df['ImageID'] + '.jpg' # Add file extension so we can compare
312 downloaded_df = original_df[original_df['ImageID'].isin(downloaded_files)]
313 downloaded_df['ImageID'] = downloaded_df['ImageID'].str.rstrip('.jpg')
314 downloaded_df = downloaded_df.drop_duplicates(subset=['ImageID'], keep=False)
315
316 # Filter down to a subset
317 size_of_subset = 50
318 print(downloaded_df["ImageID"].nunique())
319 subset = random.sample(downloaded_df['ImageID'].tolist(), size_of_subset)
320 subset_df = downloaded_df[downloaded_df["ImageID"].isin(subset)]
321 print(subset_df)
322
323 # Download to drive
324 downloaded_df.to_csv('drive/MyDrive/Data/bus_data/subset.csv', index=False)

```

Faster R-CNN



YOLOd

