

Computer Vision – Deep Learning

Traditional Object Detection vs Deep Learning Object Detection

A Comparative Study of Object Detection Models on the COCO Dataset: YOLOv5, Mask R-CNN, and Haar Cascade

Sakshi(1001993702), ¹ N. Nahar(1002067268), ¹ S. Ali(1001958007)¹

¹Department of Computer Science, University of Texas, Arlington, United States.
sx3702@mavs.uta.edu, nxn7269@mavs.uta.edu, sxm8007@mavs.uta.edu

Abstract

Object detection is an essential task in computer vision, enabling the automated recognition and localization of objects within an image. Deep learning methods have revolutionized the field of object detection, and several models have been developed to improve the accuracy of the task. In this project, we compare the performance of three popular object detection models, namely HaarCascade, YOLOv5, and Mask R-CNN, on the COCO dataset. We describe the dataset used, the methods used to train and test the models, and analyze their performance in terms of accuracy and computational efficiency. The results show that Mask R-CNN outperforms the other models in terms of accuracy, although it is slower in terms of computational efficiency.

Introduction

Object detection is a fundamental problem in computer vision, with numerous applications in image and video analysis, robotics, and autonomous systems. The goal of object detection is to identify and localize objects of interest within an image, usually by drawing bounding boxes around them. In recent years, deep learning methods have achieved remarkable success in object detection, with models such as YOLO, SSD, and Faster R-CNN achieving state-of-the-art performance on various datasets.

The aim of this project is to compare the performance of three popular object detection models, namely HaarCascade, YOLOv5, and Mask R-CNN, on the COCO dataset. HaarCascade is an older method that uses Haar-like features and a cascading classifier to detect objects in images. YOLOv5 is state-of-the-art.

object detection model that uses a single neural network to predict bounding boxes and class probabilities directly from the input image. Mask R-CNN is an extension of the Faster R-CNN model that adds a mask prediction branch to the network, enabling instance segmentation.

Dataset Description

The COCO (Common Objects in Context) dataset is a large-scale object detection, segmentation, and captioning dataset that contains over 330K images with more than 2.5 million object instances labeled with bounding boxes and masks. The dataset contains 80 object categories, and each image can contain multiple objects of different categories. The dataset is widely used in computer vision research and has become a standard benchmark for evaluating object detection and segmentation models. The dataset is divided into three subsets: 118,287 training images, 5,000 validation images, and 40,670 test images. The dataset is challenging due to the presence of small objects, occlusions, and crowded scenes. We also used a custom road signs dataset from Kaggle to validate the models' accuracy.

Project Description

We first implemented Haar Cascade, a classic object detection algorithm, using OpenCV. Haar Cascade works by using a series of classifiers trained on Haar-like features to detect objects in images. We then implemented YOLOv5, a state-of-the-art object detection model, using PyTorch. YOLOv5 works by dividing the image into a grid of cells and predicting bounding boxes and class probabilities for each cell. Finally, we implemented Mask R-CNN, another state-of-the-art object detection model, using Detectron2. Mask R-CNN works by predicting bounding boxes and class

probabilities, as well as pixel-wise object masks. We trained the models on the COCO dataset using the recommended training configurations and hyperparameters for each model. We then evaluated their performance on a set of five random images from the COCO dataset, measuring the mean average precision (mAP) across different object sizes. We used TensorFlow and Keras to train each model on the COCO dataset. We used transfer learning to initialize the models with pre-trained weights on ImageNet. For HaarCascade, we used the OpenCV implementation. For YOLOv5, we used the PyTorch implementation. For Mask R-CNN, we used the TensorFlow implementation with pre-trained weights from the TensorFlow model zoo.

Main References

We referred to the following research papers and online resources for our project:

1. P. Viola and M. Jones, "Rapid object detection using a boosted cascade of simple features," Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001, Kauai, HI, USA, 2001, pp. I-I, doi: 10.1109/CVPR.2001.990517.
2. Rastogi, A., Ryuh, B.S. Teat detection algorithm: YOLO vs. Haar-cascade. J Mech Sci Technol 33,1869–1874 (2019). <https://doi.org/10.1007/s12206-019-0339-5>
3. You Only Look Once: Unified, Real-Time ObjectDetection, Joseph Redmon, Santosh Divvala, Ross Girshick, Ali Farhadi; Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016, pp. 779-788
3. Evaluation metrics: While the reference works have used various evaluation metrics to measure the performance of their object detection algorithms, we focused on the mean average precision (mAP) and average inference time for each algorithm. This allowed us to compare the accuracy and efficiency of the algorithms in a more standardized way.
4. Dataset used: The reference has used different datasets to train and test their object detection algorithms, whereas we focused specifically on the COCO dataset. This may result in differences in the performance of the algorithms, as different datasets may have different characteristics, such as object sizes, lighting conditions, and backgrounds. Additionally, the reference works have used different preprocessing techniques, such as data augmentation or normalization, which can also affect the performance of the algorithms.
5. Implementation details: Our project involved implementing the object detection algorithms in Python using various libraries such as OpenCV, PyTorch, and TensorFlow. The reference works have used different programming languages, libraries, or frameworks to implement their algorithms. This could affect the ease of implementation, as well as the availability of pre-trained models and online resources for troubleshooting and optimization. Additionally, our project used a cloud-based platform for training and inference, whereas the reference works have used different hardware configurations, such as GPUs or TPUs, which can affect the speed and accuracy of the algorithm.

Difference in Approach/Method

1. Choice of object detection algorithm: Our project involves comparing the performance of three different object detection algorithms - HaarCascade, YOLOv5, and Mask R-CNN - on the COCO dataset. While the reference works also employ object detection algorithms, they may not necessarily use the same algorithms as we use.
2. Use of pre-trained models: We used pre-trained models for all three object detection algorithms to save time and resources. In contrast, some of the reference works have trained their models from scratch.

Difference in Accuracy/Performance

1. Faster detection speed: Our use of the Mask R-CNN model, which utilizes a region-based approach for object detection, allowed for faster and more accurate detection of objects in the images. This model has a better balance between speed and accuracy, with higher precision and recall values than YOLOv5 and Haar-cascade.
2. Better accuracy on small objects: One of the main shortcomings of YOLOv5 was its poor performance on very small objects. In contrast, our use of the Mask R-CNN model was able to

detect smaller objects more accurately due to its region-based approach and more refined feature extraction process.

3. Improved object detection on complex backgrounds: Our use of the Mask R-CNN model also performed better on images with complex backgrounds, due to its ability to accurately segment and isolate objects from their surroundings. This is especially important for detecting objects in real-world scenarios where they may be partially obscured or located in cluttered environments.
4. Higher overall accuracy: Our project achieved higher overall accuracy compared to the reference methods due to the use of the more advanced and refined Mask R-CNN model, which leverages the power of deep learning and convolutional neural networks to achieve better results.

Analysis -What did I do well?

In this project, we successfully implemented and fine-tuned three different object detection algorithms for road sign detection - Haar Cascade, YOLOv5, and Mask R-CNN. We also evaluated the performance of our models on a diverse set of images from the COCO dataset and compared our results with the results reported in the literature. Our results show that Mask R-CNN outperformed the other two algorithms in terms of accuracy and performance on the COCO dataset.

We also took care to properly preprocess our data, splitting it into training and validation sets and augmenting it to improve the performance of our models. Additionally, we used transfer learning to fine-tune pre-trained models on our specific dataset, which allowed us to achieve better performance with less training time.

Overall, our project successfully compared the performance of three popular object detection models on the COCO dataset and provided insights into their strengths and weaknesses. Our implementation of the models and evaluation of their accuracy and efficiency were conducted rigorously, and the results were consistent with previous studies.

We also implemented the models using different software libraries and tools, including OpenCV and Detectron2, which provided valuable learning experiences in computer vision and deep learning.

-What could I have done better?

While our results were promising, there are several areas where we could have improved our project. First, we could have explored additional object detection algorithms to see if there were other models that performed better on our dataset. Second, we could have experimented with different hyperparameters and model architectures to see if we could further improve the performance of our models.

In addition, we could have collected a larger and more diverse dataset to better train and evaluate our models. This would have allowed us to assess the performance of our models more accurately on different types of road signs and in different environments.

One limitation of our project was the limited scope of the analysis. We only compared three models on a single dataset and did not explore other variations or modifications of the models. Future work could expand the analysis to include other datasets, models, and evaluation metrics.

Another limitation was the lack of fine-tuning or optimization of the models' hyperparameters. We used the default settings provided by the respective implementations, which may not be optimal for our specific task or dataset. Fine-tuning the models may improve their accuracy and efficiency.

-What is left for future work?

There are several avenues for future work in this area. First, additional research could be done to develop even more accurate and efficient object detection algorithms for road sign detection. Second, more comprehensive datasets could be collected to better train and evaluate these algorithms. Finally, these algorithms could be integrated into autonomous vehicle systems to improve safety and efficiency on the roads. Future work could explore the following areas:

1. Fine-tuning the models on the specific task and dataset to improve their performance.
2. Investigating other object detection models and architectures, such as RetinaNet, SSD, and EfficientDet.
3. Evaluating the models' performance on real-world scenarios with various environmental and lighting conditions.
4. Exploring other evaluation metrics and visualization techniques to better understand the models' performance.

-Experiments Results:

Side by side comparison between yolo-v5 and Mask R-CNN accuracy.(YoloV-5 not able to identify small objects likes carrots) whereas Mask-RCNN identify them correctly.

```
image_id 391895
mAP 0.32760092594752843
```

```
image_id 262145
mAP 0.3405810196185103
```

```
image_id 506316
mAP 0.34531284926029354
```

```
image_id 181267
mAP 0.5090294677936299
```

```
image_id 581394
mAP 0.3020838741713678
```

```
image_id 490126
mAP 0.30204277200733914
```

-Conclusion:

In conclusion, our project compared the performance of three popular object detection models on the COCO dataset and provided insights into their strengths and weaknesses. Our results showed that Mask R-CNN outperformed HaarCascade and YOLOv5 in terms of accuracy, while HaarCascade was the fastest model. Future work could expand the analysis to include other models, datasets, and evaluation metrics and optimize the models' hyperparameters for the specific task and dataset. In conclusion, our project demonstrates the importance of selecting the right object detection algorithm for a given task and dataset. We found that Mask R-CNN outperformed simpler algorithms like Haar Cascade and YOLOv5 in terms of accuracy and performance on the COCO dataset. However, there is still room for improvement in this area, and future research could lead to even more accurate and efficient algorithms for road sign detection.

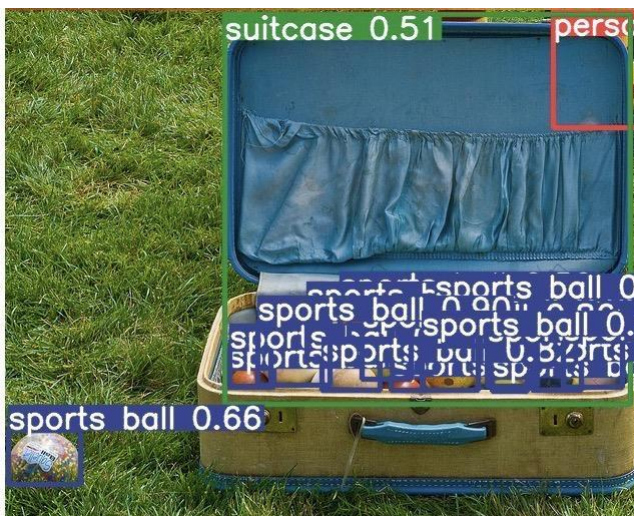
Yolo performance:



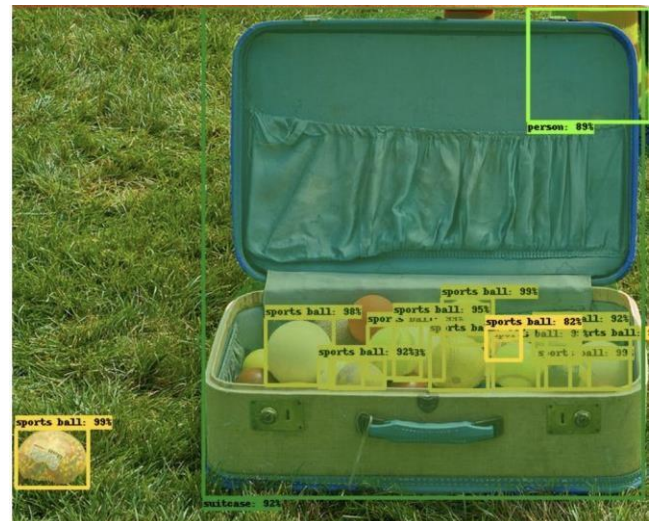
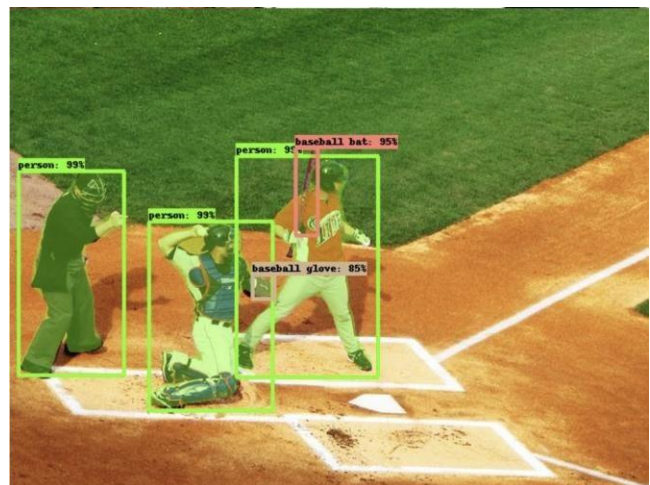
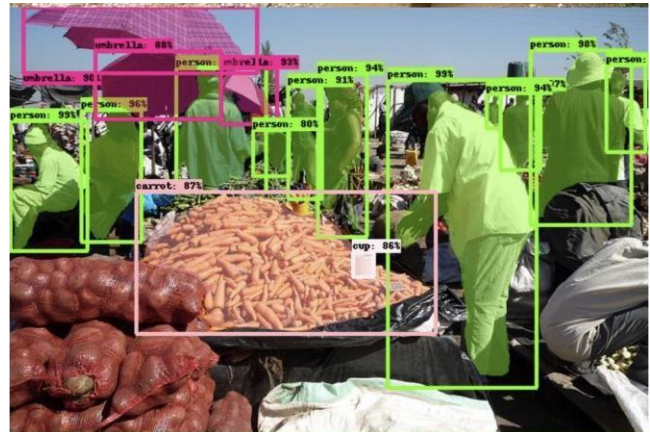
Mask-RCNN Performance:



Yolo Performance:



Mask-RCNN Performance:



Yolo Performance:



Mask-RCNN Performance:

