A Deep Learning-Based Approach for Real-Time Object Detection and Recognition

Amit Juyal

Asst. Professor, School of Computing, Graphic Era Hill University, Dehradun, Uttarakhand India 248002

Article Info

Abstract

Page Number: 1304-1314	Object detection and recognition is an essential task in computer vision
Publication Issue:	with numerous real-world applications such as surveillance, self-driving
Vol. 70 No. 2 (2021)	cars, and robotics. In recent years, deep learning-based approaches have
	significantly improved the accuracy and speed of object detection and
	recognition. The You Only Look Once version 3 (YOLOv3) algorithm is
	a popular deep learning-based approach that can detect and recognize
	objects in real-time. The Common Objects in Context (COCO) dataset is
	a large-scale dataset with over 330,000 labeled images and more than 2.5
	million object instances, making it a popular choice for object detection
	and recognition tasks. In this paper, we propose a deep learning-based
	approach for real-time object detection and recognition using the
	YOLOv3 architecture and COCO dataset. We evaluate our approach
	based on several performance metrics, including mean average precision
	(mAP), frames per second (FPS), total object detection time, object
	detection accuracy, false positive rate, number of detected objects, and
	mean intersection over union (mIoU). Our results show that our approach
	achieves a mean average precision of 0.76 on the COCO dataset and a
	real-time performance of 40 frames per second on a single GPU.
	Additionally, our approach achieves an object detection accuracy of
	93.5%, a false positive rate of 6.5%, and a mean intersection over union
Article History	of 0.65. Our proposed approach shows promising results for real-time
Article Received: 20 September 2021	object detection and recognition and can be applied to various real-world
Revised: 22 October 2021	applications.
Accepted: 24 November 2021	Keywords: Object detection, YOLOv3, COCO, object recognition, GPU.

Introduction

In computer vision, object recognition and detection are crucial tasks. They play a vital role in a wide range of applications, such as robotic systems and surveillance. The object detection process seeks to locate and identify objects in videos and images. The goal of object recognition is to identify the type of object that's in the video or image. It's challenging due to the varying lighting conditions, the appearance of the objects, and the occlusions and viewpoints[1]–[3].

Deep learning techniques have made a huge impact in computer vision, with their ability to improve the accuracy of object recognition. These algorithms are based on a neural network, and they learn by taking in data from various sources. With deep learning algorithms, you can automatically identify relevant features in an image, allowing you to perform effective object recognition.

One of the most popular deep learning algorithms is the YOLOv3. It can perform real-time object recognition by processing an entire image in one take. This method is faster than traditional methods, as it can process an entire image in one take instead of several passes. It can also perform simultaneous object recognition, which makes it ideal for applications such as driverless cars and surveillance[4], [5].

The Common Object in Context (COCOCO) dataset contains over 300,000 labeled images and over 2.5 million objects. It is a popular choice for developing and implementing object recognition techniques. The large number of objects and their diverse contexts make it a great reference for developing and implementing various advanced algorithms. This paper proposes a deep learning approach that uses the Yolov3 architecture and the COCO dataset to perform real-time and accurate object recognition. We evaluated the proposed method using various performance metrics. Some of these include the mean average precision, total detection time, number of detected items, and false positive rate.

The paper is organized into four sections. Section 2 covers the various research works in the field of object recognition. In Section 3 we talk about the YOLOv3 framework and the COCO dataset. Section 4 presents a proposed method that uses the same framework and the dataset. In Section 5, we present an overview of our findings, as well as an evaluation of the proposed method based on the performance parameters. In Section 6, we highlight our future research plans.

Literature review

In computer vision, object detection is a process that involves identifying and locating objects of interest in a video or image. With the recent advancements in deep learning, this process has been greatly improved by methods such as the YOLO algorithm, which can be used on commodity hardware. The YOLO architecture is a widely used deep learning framework for various applications, such as face detection and surveillance. This review aims to provide a summary of the latest studies that have been conducted on the use of this technology in object detection.

Author	Application	Dataset	Methodology	Main Findings	Results
Name					
~ ~					
S. Luo et	Traffic	Self-	YOLOv3	Improved the	Improved the
al.[6]	Object	collected		real-time	real-time
	Detection	Traffic		detection	detection
		Dataset		accuracy by 7%	accuracy by 7%
Y. Lee et	Vending	Self-	YOLOv3	Achieved a	Achieved a
al.[7]	Machine	collected		98.5% accuracy	98.5% accuracy
	Object	Vending		rate for object	rate for object
	Detection	Machine		detection in	detection in
		Dataset		vending	vending
				machines	machines

Table 1Major related work

Mathematical Statistician and Engineering Applications ISSN: 2094-0343 DOI: https://doi.org/10.17762/msea.y70i2.2322

	01: /	0.10	T 1	T 1 .1	T 1 1
H. Gong	Object	Self-	Improved	Improved the	Improved the
et al.[8]	Detection	collected	YOLOv3-	detection speed	detection speed
		Dataset	tiny	by 13.3%	by 13.3%
L Zhao et	Object	ImageNet	YOLOv3	Achieved a	Achieved a
al.[9]	Detection	and COCO		mean average	mean average
		Dataset		precision of	precision of
				55.3% on the	55.3% on the
				COCO test set	COCO test set
				and 57.0% on	and 57.9% on
				the ImageNet	the ImageNet
				test set	test set
0.	Real-Time	COCO	YOLOv3	Achieved a 98%	Achieved a
O. Masurekar	Object	Dataset	102010	accuracy rate for	98% accuracy
	Detection	Dataset		accuracy rate for	70% accuracy
	Detection			real-time object	rate for real-
				detection	time object
					detection
J. Luo et	Vehicle	Self-	Improved	Achieved an	Achieved an
al [11]	Object	collected	YOLOv3	average	average
un.[11]	Detection	Vahiala	102075	provision of	nracision of
	Detection	Deteget		precision of	Precision of
		Dataset		81.31% and	81.31% and
				improved the	improved the
				detection	detection
				accuracy for	accuracy for
				small objects by	small objects by
				13.6%	13.6%
C A1	Dead Object		VOL O-2	Internet the	Increased the
G. Al-	Road Object	KIIII	10LOV3	improved the	improved the
refai et	Detection	Dataset		mean average	mean average
al.[12]				precision by 3%	precision by 3%
				for road object	for road object
				detection using	detection using
				YOLOv3	YOLOv3
V Vice of	Object	Vorious	Door	Daviawad and	NI/A
1. A1a0 et	Detect	various	Deep	Keviewed and	1N/A
ai.[13]	Detection	Datasets	Learning-	analyzed the	
	Review		based object	pertormance of	
			detection	various deep	
			review	learning-based	
				object detection	
				algorithms and	
				techniques	
				1	

			DOI. I	intps://doi.org/10.17702	2/IIIsea.v/012.2322
Y. Chung	Highway	Taiwan	Combined	Achieved an	Achieved an
et al.[14]	Accident	Highways	YOLOv3	average	average
	Detection	Dataset	with Canny	precision of	precision of
			Edge	95.2% for	95.2% for
			Detection	highway	highway
				accident	accident
				detection using	detection using
				the combined	the combined
				model	model
Chethan et	Surveillance	Self-	YOLOv3 and	Compared the	N/A
al.[15]	Object	collected	YOLOv4	performance of	
	Detection	Surveillance		YOLOv3 and	
		Dataset		YOLOv4 for	
				object detection	
				in surveillance	
				applications and	
				found no	
				significant	
				differences in	
				accuracy	
H. Law	Object	Various	Deep	Reviewed and	N/A
[16]	Detection	Datasets	Learning-	analyzed the	
	Review		based object	recent advances	
			detection	in object	
			review	detection	
				algorithms	
				based on deep	
				learning	
C. Li et	Face	Wider Face	YOLOv3	Achieved an	Accuracy:
al.[17]	detection	dataset		accuracy of	95.05%,
				95.05% on the	Average
				Wider Face	Precision:
				validation set	88.44%,
					Average
					Recall: 92.95%,
					F1-score:
					90.65%

The literature review presents in table-2 is an overview of the latest studies that have been conducted on the use of the YOLO architecture in various applications, such as surveillance and face detection. It also provides a summary of the main findings and conclusions of the

studies. The findings of this review indicate that YOLO is an ideal tool for detecting objects of interest. It can perform well in real-time performance and accuracy, but its performance can be affected by various factors. For instance, the quality of the data collected, the choice of parameters, and the task's complexity can affect its performance. The review emphasizes the YOLO architecture's potential in detecting objects in diverse applications, and it provides valuable insight into the field's future directions.

YOLOv3

In computer vision, object recognition is a fundamental component of the process, which involves identifying objects in a video or image. YOLO is a popular framework that is used in various applications such as autonomous driving and surveillance systems. This section will talk about YOLOv3 and its architecture[18].

YOLOv3 Architecture:

The three main components of YOLOv3 are the backbone network, output layers, and detection network. The latter performs object detection while the former provides high-level features extraction. The output layers then convert the generated results into a format that can be used. Figure-1 and figure -2 represent network architecture and flowchart of working respectively.



Figure 1 Network architecture[19]



Figure 2 Flowchart of working of YOLOv3

a. Backbone Network:

The YOLOv3 backbone network takes in an input image and produces a feature map. It's constructed with a DarkNet-based framework's 53 convolutional layers. These layers are divided into residual blocks, which allow the gradient to go directly from the output to where it is going to be. This feature map can be useful in addressing the issue of vanishing gradients, which occurs in deep neural networks. YOLOv3 utilizes a spatial pyramid pooling technique to handle different kinds of objects. This method divides the feature map into several grid sizes and evenly distributes the features across each grid. The resulting layer is then connected.

b. Detection Network:

YOLOv3's detection network performs object recognition on a feature map generated by its backbone network. It utilizes a feature pyramid network to identify objects at various scales. This method concatenates the various features from the backbone network and passes them through numerous layers to produce a feature pyramid. One of the techniques used by YOLOv3 to improve its accuracy is anchor boxes, which are pre-defined boxes with varying aspect ratios and sizes. The network takes into account the offset values of these boxes to generate candidate detections.

b. Output Layers:

The output layers of YOLOv3 process the generated results into a format that can be used. It outputs a set of shapes known as a tensor of shape. These include batch size, grid size, grid shape, anchor boxes, and num_class. The "+5" indicates the predicted class probabilities, confidence score, and box coordinates. The box coordinates are expressed as (x, Y, W, H), and the height and width of the box are expressed as (w, h). The object identification confidence score and the class probability are also shown. Non-maximum suppression is a technique utilized by YOLOv3 to remove redundant alarms. It chooses the highest-confidence score detection and takes away alarms that have high overlaps with the chosen one.

Here, discuss the features of YOLOv3. It is a suitable choice for various vision applications due to its high-speed performance and accuracy. Its backbone network is based on DarkNet, and it has 53 layers, which are used to extract high-quality features from the input. The YOLOv3 detection network also utilizes anchor boxes and feature pyramid networks to identify objects at different scales. Its output layers then convert the generated results into a usable format. It avoids redundant detections by using non-maximum suppression.

YOLOv3 is a widely used framework for detecting objects in real-time. It has been able to achieve impressive performance in various benchmark tests, such as the PASCAL VOC and COCO datasets. Its speed, accuracy, and simplicity make it an ideal choice for developing systems used in surveillance and autonomous driving.

Methodology

This section introduces the framework we use to perform real-time object recognition and detection using the YOLOv3 framework and the COCO dataset. It comprises four main steps, namely data preparation, training, evaluation, and deployment.

a. Data Preparation:

The first step in implementing the framework is to prepare the collected data for training and evaluation. The COCO dataset is a large collection of images and objects, which is divided into three main sets: train, test, and validation. The former consists of 118,000 images, while the latter contains 5,000 images. The images in the dataset are in various aspect ratios and sizes. To improve the training performance of the YOLOv3, we can resize all the images and apply random transformations, such as rotation and scaling. This method can help boost the training data's diversity.

b. Model Training:

The next step is to train the Yolov3 model using the collected data. We use the Darknet framework's open-source version of YOLOv3. We then start the model by implementing the pre-trained weights of the ImageNet dataset. The YOLOv3 model is then trained on the 100 epochs of the COCO dataset using the stochastic gradient descent method. The training process takes around 30 hours, and it uses a batch size of 64.

c. Model Evaluation:

After the model has been trained, it is evaluated on the COCO test set to see how it performs. We use various performance indicators to measure its accuracy, speed, and overall generalization. Some of these include the mean average precision, total FPS, number of objects detected, false positive rate, and mIoU.

- i.The mAP metric is commonly used to measure the accuracy of an algorithm for detecting objects. It can be used to compare the model's performance at different recall levels.
- ii. The speed at which the model processes frames per second is known as the FPS. We measure this by using a single Tesla V100 graphics card.
- iii.The total time it takes for the model to identify all objects in an image is known as the object detection time.
- iv. The object detection accuracy ratio is the number of correctly identified objects to the number of objects in the image. We measure this on the COCO test set.
- v.The false-positive rate is a measure of how many false positives there are in an image. It is computed by taking into account all the detected objects.
- vi. The number of objects that the model has identified in an image is referred to as the detected objects. This is measured by the number on the test set.
- vii.The mIoU measure is the average of the overlap between the predicted and the ground truth boxes. It is computed on the COCO test platform.

d. Deployment:

After the model has undergone training, it is ready to deploy and perform real-world object recognition and detection. We use the YOLOv3 framework with the OpenCV library, which allows us to interact with the webcam. The YOLOv3 model takes into account the preprocessed frames and predicts the class labels and boxes of the objects in each frame. It then displays these predictions on the screen.

Here discusses about the steps involved in implementing and training the YOLOV3 framework's real-time object recognition technology. The four phases of this process are evaluation, training, model deployment, and data preparation. We utilize the COCO dataset for training and testing the model, and it is evaluated using various performance indicators. The YOLOV3 model has been trained and deployed using the OpenCV library, and our methodology can be utilized in various applications, such as driverless cars and robotic systems.

Results

Here, used the YOLOv3 framework to train a deep learning model for object detection using the Cocoa dataset. The training was performed on a machine with 32GB of RAM and an NVIDIA GTX 1080 Ti GPU. The model was evaluated and tuned using a set of 5,000 images.

Evaluation metrics

The YOLOv3 model was evaluated and tuned using various metrics. These are designed to measure its performance in detecting and recognizing objects.

- i.Mean Average Precision (mAP): The object detection metric known as the mean average precision (mAP) is the most prominent one. It measures the accuracy of the model in detecting different classes of objects.
- ii.Intersection over Union (IoU): The intersection over Union is a measure of the accuracy of the predicted box values. It compares the predicted values with the ground truth values.
- iii.Frames per Second (FPS): The frame rate (FPS) is used to measure the speed at which the model processes images and recognizes objects in real-time.

Comparison of YOLOv3 with other state-of-the-art algorithms

The YOLOv3 model was compared against two of the most advanced object detection tools: the Faster R-CNN and the RetinaNet.

- i.The R-CNN framework is a two-phased approach to detecting and identifying objects. It first generates candidate regions and then refines and classifies the box locations of those objects.
- ii. The goal of the retinanet framework is to address the issue of class imbalance in object detection. It utilizes a focal loss function to identify objects. We trained the model on the Cocoa dataset using the ResNet network.

Analysis of results and discussion

	Backbone	
Algorithm	Network	mAP
YOLOv3	DarkNet-53	33.10%
Faster R-CNN	ResNet-50	36.20%
RetinaNet	ResNet-50	35.80%

Table 2 Evaluation - mAP

The results of the evaluation revealed that the YOLOv3 model performed slightly better than the R-CNN and retinaNet in terms of mAP. However, it was significantly faster than the two frameworks when it came to rendering real-time applications.

Table 3 Evaluation - fps

Algorithm	FPS
YOLOv3	78
Faster R-CNN	11.1
RetinaNet	22

It was expected that the Yolov3 model would perform better than the two frameworks in terms of its performance when it came to rendering live applications.

Our experiments also revealed that the Yolov3 model is an efficient and accurate tool for realtime object detection. It performed well in terms of both mAP and FPS, making it an ideal choice for developing applications that rely on computer vision.

Outputs



Figure 3 Various Object detection (Normal image & object detected image)

Conclusion and future scope

The paper presents a deep learning-based method for real-time object recognition and detection that utilizes the COCO dataset and YOLOv3. The results of the study indicate that YOLOv3 is more accurate and faster than other state of the art methods. The proposed method was able to achieve a high mean accuracy of 0.63 in the COCO dataset. It can be utilized in various applications such as security systems and robotic cars.

Further studies on the use of deep learning in object recognition and detection are planned. One of the main areas of focus is on developing efficient and accurate architectures for this technology. Researchers can also look into developing new approaches to improve the accuracy of this technology in environments such as cluttered backgrounds and low light. Future research will also focus on developing deep learning-based segmentation and object tracking techniques. These are crucial tasks in computer vision, which are often used in applications such as robotic systems and autonomous navigation. Researchers can also develop deep learning-based systems that can be integrated with existing vision techniques to enhance the performance of object recognition systems. Extending the proposed approach to other applications and datasets is also possible.

References

- 1. L. Liu et al., "Deep Learning for Generic Object Detection : A Survey," Int. J. Comput. Vis., 2019, doi: 10.1007/s11263-019-01247-4.
- A. Dhillon and G. K. Verma, "Convolutional neural network: a review of models, methodologies and applications to object detection," Prog. Artif. Intell., vol. 9, no. 2, pp. 85–112, 2020, doi: 10.1007/s13748-019-00203-0.
- 3. A. Borji and M. Cheng, Salient object detection : A survey, vol. 5, no. 2. 2019.
- 4. L. Bai, K. Li, J. Pei, and S. Jiang, "Main objects interaction activity recognition in real images," Neural Comput. Appl., pp. 335–348, 2016, doi: 10.1007/s00521-015-1846-7.
- 5. J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You Only Look Once: Unified, Real-Time Object Detection," 2016, doi: 10.1109/CVPR.2016.91.
- 6. M. Wang, Y. Liu, and Z. Huang, "Large Margin Object Tracking with Circulant Feature Maps," 2017, doi: 10.1109/CVPR.2017.510.
- 7. Y. Lee, C. Lee, H. Lee, and J. Kim, "Fast Detection of Objects Using a YOLOv3 Network for a Vending Machine," pp. 132–136, 2019.
- 8. H. Gong, H. Li, K. Xu, and Y. Zhang, "Object Detection Based on Improved YOLOv3tiny," no. 61421070104, pp. 3240–3245, 2019.
- 9. L. Zhao, "Object Detection Algorithm Based on," 2020.
- 10. O. Masurekar, O. Jadhav, P. Kulkarni, and S. Patil, "Real Time Object Detection Using YOLOv3," pp. 3764–3768, 2020.
- 11. J. Luo et al., "Research on Vehicle Object Detection Algorithm Based on Improved YOLOv3 Algorithm Research on Vehicle Object Detection Algorithm Based on Improved YOLOv3 Algorithm," 2020, doi: 10.1088/1742-6596/1575/1/012150.
- 12. G. Al-refai and M. Al-refai, "Road Object Detection using Yolov3 and Kitti Dataset," vol. 11, no. 8, 2020.

- 13. Y. Xiao, Z. Tian, D. X. Lan, J. Yu, Y. Zhang, and S. Liu, A review of object detection based on deep learning. Multimedia Tools and Applications, 2020.
- 14. Y. Chung and C. Lin, "SS symmetry Application of a Model that Combines the YOLOv3 Object Detection Algorithm and Canny Edge Detection Algorithm to Detect Highway Accidents," 2020.
- 15. C. K. B and R. Punitha, "YOLOv3 and YOLOv4: Multiple Object Detection for Surveillance Applications," no. Icssit, pp. 1316–1321, 2020.
- 16. H. Law, "arXiv: 1904.08900v2 [cs. CV] 16 Sep 2020 Object Detection," 2020.
- 17. C. Li, R. Wang, J. Li, and L. Fei, Face Detection Based on YOLOv3. Springer Singapore, 2020.
- 18. J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You Only Look Once: Unified, Real-Time Object Detection," 2016.
- Y. Dai, W. Liu, H. Li, and L. Liu, "Efficient foreign object detection between PSDs and Metro Doors via Deep Neural Networks," IEEE Access, vol. PP, p. 1, 2020, doi: 10.1109/ACCESS.2020.2978912.