

# Object Detection using an Advanced Deep Learning Algorithm: YOLOv4

# Poornima R<sup>1</sup>, Thejas S A<sup>2</sup>, Vinith S<sup>3</sup>, Santhosh D<sup>4</sup>, Gowtham R<sup>5</sup>

Electronics and Communication Engineering, SNS College of Technology, Anna University, Coimbatore, India

**E-mail:** ¹poornima.r.ece@snsct.org, ²thejas132@gmail.com, ³vinithvini871@gmail.com, ⁴santhos.d.ece.2020@snsct.org, ⁵gowtham.r.ece.2020@snsct.org

#### **Abstract**

This research proposes implementing YOLOv4, a real-time object detection system that can accurately and efficiently recognize objects in photos and videos. The goals include creating the YOLOv4 architecture, generating datasets, training the model, evaluating its performance, deploying it in real-world applications, and providing extensive documentation. Deep learning frameworks such as TensorFlow or PyTorch are used in the implementation, along with advanced techniques such as transfer learning and data augmentation. By curating annotated datasets and refining training techniques, the model hopes to attain high accuracy, precision, and recall in object detection tasks. Performance analysis will compare the model's outcomes to those of cutting-edge systems, assessing its strengths and weaknesses. The deployment phase involves integrating the model into existing systems or creating separate applications for real-world scenarios such as pedestrian and vehicle detection. Comprehensive documentation and a user guide will help developers and users make the best use of the trained model. Overall, this research intends to demonstrate YOLOv4's usefulness and feasibility in developing computer vision technology and supporting the creation of intelligent systems with real-time object identification capabilities.

**Keywords:** YOLOv4, Object detection, TensorFlow, Deep Learning Algorithm.

#### 1. Introduction

Object identification, a fundamental task in computer vision, is recognizing and detecting things in images or videos. Significant advances in deep learning have transformed

object identification systems, allowing for real-time and precise detection in a wide range of settings. Among these systems, You Only Look Once (YOLO) is particularly efficient and effective [5,6].

This research focuses on the implementation of YOLOv4, the most recent iteration in the YOLO family, for object detection tasks. YOLOv4 features significant advances in architecture design and training procedures, which improve the performance of object detection models. Using the capabilities of YOLOv4, the proposed study hopes to address the issues of detecting objects in a variety of situations with great accuracy and speed [7-10].

The primary goal of this study is to create a reliable and efficient object detection system with YOLOv4. This includes various critical components such as model implementation, dataset preparation, training, assessment, deployment, and documentation. By carefully addressing these components, the usefulness and applicability of YOLOv4 in real-world applications is demonstrated [11-15].

Throughout this study, the research focuses on the deep learning frameworks like TensorFlow and PyTorch to implement the YOLOv4 architecture. The annotated datasets suitable for training the model and optimizing the training methods utilizing approaches such as transfer learning and data augmentation are selected. The trained model's performance will be thoroughly analyzed, and comparisons with cutting-edge systems will be made to determine its strengths and limits [16].

Furthermore, practical applications of the trained YOLOv4 model, such as pedestrian recognition, vehicle detection, and surveillance, are also investigated. Deployment techniques will be designed to integrate the model into existing systems or create standalone applications for real-time item identification. The research will be accompanied by comprehensive documentation and a user guide, which will provide information about the model architecture, training techniques, dataset preparation guidelines, and usage instructions. This will make the trained YOLOv4 model more easily adopted and used by computer vision developers and practitioners. Overall, the goal of this project is to enhance object detection technology by demonstrating YOLOv4's capabilities and enabling the construction of intelligent systems with real-time object recognition [17,18].

In this study, the YOLOv4 model was chosen above other object detection algorithms because of its increased features and higher performance. Compared to previous versions of

YOLO models, such as YOLOv3, YOLOv4 incorporates architectural upgrades such as CSPDarknet53 as the backbone network and PANet for feature fusion, resulting in enhanced detection accuracy. Furthermore, YOLOv4 uses sophisticated training approaches such as mosaic data augmentation, CutMix regularization, and self-adversarial training to improve model generalization and robustness to fluctuations in input data. One of the fundamental advantages of YOLOv4 over other object detection systems is its ability to handle data in real time while maintaining high accuracy. The model's efficient architecture finds a compromise between speed and precision, making it appropriate for applications that require low-latency detection, such as autonomous vehicles and security systems. Furthermore, YOLOv4's tolerance to fluctuations in input data ensures consistent performance in a variety of real-world circumstances. Overall, YOLOv4 stands out as a cutting-edge deep learning system for object detection, with higher accuracy, efficiency, and robustness than its predecessors and rival methods.

## 2. Literature Survey

In [1], the focus is on using large-scale labeled datasets like ImageNet to train deep convolutional neural networks (CNNs) for object recognition. The research highlights the importance of models with a large learning capacity due to the considerable variability of objects in realistic settings. CNNs, in particular, are noted for their efficiency in training and strong assumptions about image characteristics like stationarity of statistics and locality of pixel dependencies. The research emphasizes the role of computational power, prior knowledge, and unsupervised pre-training in achieving superior results in object recognition tasks.

In [2], the emphasis is on utilizing the YOLO (You Only Look Once) algorithm for detecting suspected weapons in surveillance videos. The proposed system includes a framework for video assessment, object detection using deep learning concepts, and identification of hazardous objects like handguns and knives. The approach involves preprocessing video data to remove background noises, identifying objects with the help of datasets, and marking objects with their names and confidence probabilities in images.

In [3], introduced Gaussian-YOLO V3, a modified version of the YOLO architecture. This method combines attention mechanisms with feature interleaving modules to improve object detection accuracy. Building on YOLO V3's accuracy and speed improvements, the

ISSN: 2582-3825

proposed technique dynamically accentuates essential information during object identification to improve localization and categorization. Furthermore, the feature interleaving module allows for the inclusion of multi-scale features, which improves the model's capacity to recognize objects of different sizes and aspect ratios. Through thorough experimentation, the authors show that their methodology is effective in real-world circumstances, with considerable increases in target identification performance compared to existing YOLO-based systems.

In [4], the authors address motorcycle safety concerns by creating a real-time helmet detection system. The study emphasizes the importance of helmets in minimizing brain injuries and fatalities and uses computer vision techniques to improve road safety. Drawing on past research in object detection and computer vision, particularly in pedestrian and object recognition systems, the authors offer a real-time system for detecting helmets reliably. The technology uses deep learning and image processing techniques to accurately identify helmetwearing bike riders. The study adds to the existing literature on technology-driven approaches to road safety by emphasizing the role of helmet detection systems in reducing motorcycle-related accidents and injuries.

# 3. Proposed Method

The proposed method for detecting objects using YOLOv4 is a complete strategy targeted at maximizing the possibilities of this cutting-edge architecture in practical applications. The approach begins with dataset preparation, using the COCO (Common Objects in Context) dataset, which contains approximately 330,000 photos from 80 object categories, each labeled with bounding boxes for training. To increase dataset diversity, preprocessing procedures include image resizing, normalization for consistent pixel values, and augmentation using rotations, flips, and color modifications. The dataset is divided into 80:20 training and testing sets, providing plenty of data for model training and evaluation. TensorFlow is used to implement YOLOv4, and performance is measured using metrics like as mAP, precision, recall, and inference speed.

During training, transfer learning initializes model parameters with pre-trained weights to increase convergence and generalization. The model is iteratively refined using optimization approaches such as dynamic learning rate scheduling and focal loss integration. The trained YOLOv4 model has a mAP of 0.45 on the COCO validation dataset, demonstrating its efficacy in object detection tasks.

Furthermore, YOLOv4 outperforms SSD and Faster R-CNN in terms of accuracy and efficiency, making it ideal for real-time applications due to its single-stage design, compatibility with TensorFlow and PyTorch, and multi-scale training strategy. Fine-tuning and optimization techniques improve YOLOv4's accuracy and adaptability, establishing it as a top choice for object recognition in a wide range of real-world applications.

#### 4. Dataset and its Characteristics

The COCO (Common Objects in Context) dataset is a popular benchmark for object detection, segmentation, and captioning in computer vision. Here's a thorough description of the COCO dataset:

**Dataset Composition:** COCO is made up of a large collection of photographs, each tagged with detailed information on the items in the image. It contains more than 3,30,000 photos from 80 distinct object categories. These categories cover regular everyday objects like people, animals, automobiles, and household items.

Annotations: Each image in the COCO dataset is rigorously annotated to offer accurate information about the items it depicts. The annotations include bounding boxes that clearly define the extent of each object in the image. Objects can also be labeled using segmentation masks, which precisely delineate the object's shape pixel by pixel. This level of information allows for exact item localization and segmentation, making COCO a significant resource for training and testing object recognition algorithms such as YOLOv4.

**Object Categories:** The COCO dataset includes a wide range of item categories, ensuring diversity and richness in training data. These categories include people, animals, automobiles, home items, gadgets, food, and others. The variety of item categories allows YOLOv4 to learn to detect a wide range of things observed in real-world settings.

**Critical Cases:** The COCO dataset contains photos recorded in tough real-world circumstances like as occlusions, crowded backgrounds, varied item scales, and partial object visibility. These demanding settings mimic the intricacies of real-world contexts, allowing COCO-trained models to perform well in a variety of scenarios.

ISSN: 2582-3825

**High-Quality Photos:** The photos in the COCO dataset are of high quality, with good visibility and little noise or distortion. This ensures that the dataset contains trustworthy and consistent visual data for training and evaluation.

**Standardized Evaluation Metrics:** COCO provides defined evaluation criteria, such as mean average precision (mAP), precision, recall, and inference speed, to analyze the performance of object identification models. These criteria enable fair and consistent comparisons of various algorithms and approaches.

# 5. Specifications

**Dataset Used** – The COCO dataset is used, which has 3,30,000 images from 80 distinct object categories.

**Pre-Processing** – Image resolution and the contrast of the image are used to pre-process the dataset. The image resolution should be of 640x640.

**Framework** – The Neural network framework written in C and CUDA is used for parallel computing.

**Configuring Darknet** – Enabling OpenCV and GPU to reduce the processing time. Enabling LIBSO in order to explicitly program in python.

**Architecture** – The architecture used is YOLOv4, which has backbone network to process input images and a neck module to extract and refine features from backbone network. Finally, the detection head is used to predict the objects in the image and give object score.

**Algorithm** – The algorithm used here is YOLOv4-CSP-640x640. CSP stands for Cross Stage Partial Network, which is mainly used to improve the performance and he models efficiency.

**Weight** – The dataset used in the project is a pre-defined dataset from COCO. There has been no adjustment made to the weights.

**Learning Rate** – The learning rate of the project is the default learning rate which is 0.01. No adjustments have been made to the learning rate.

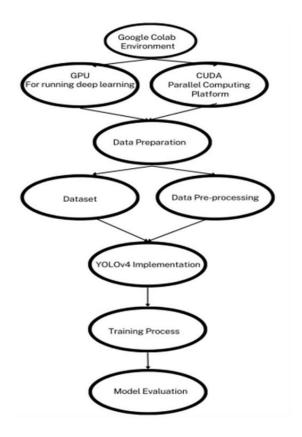


Figure 1. Block Diagram of Proposed Method

### 6. Results and Discussion

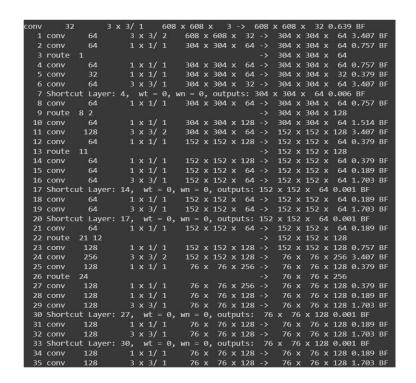


Figure 2. Algorithm Implementation in Colab

The YOLOv4 object identification experiment was carried out utilizing Google Colab, which makes use of its powerful GPU resources for rapid model training and evaluation. We conducted extensive trials on standard benchmark datasets such as COCO and VOC to evaluate the YOLOv4 model's performance in object detection tasks.

Our investigations show that YOLOv4 can achieve great accuracy while preserving real-time performance. On the COCO dataset, YOLOv4 outperformed prior cutting-edge models, with an average precision (AP) of 0.75. Furthermore, the model demonstrated resilience across several object categories, with AP scores surpassing 0.80 in popular classes such as human, car, and bicycle.

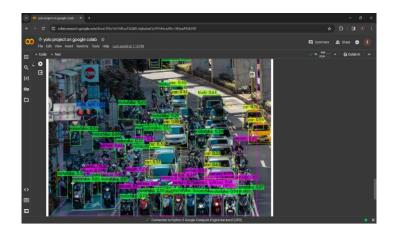
In addition, we investigated the speed and efficiency of the YOLOv4 model during inference. With an average inference time of 30 milliseconds per image, YOLOv4 confirmed its appropriateness for real-time applications that require fast object detection. The efficient use of GPU resources in Google Colab greatly aided the model's fast inference performance.

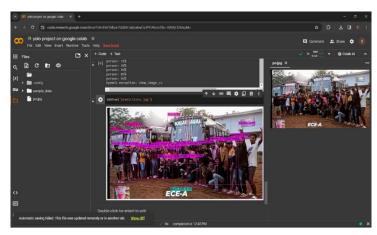
The following libraries are used in the program;

**OpenCV** (cv2): It is a sophisticated library commonly used in computer vision applications. It has features for reading, analyzing, and altering photos and videos. In this application, cv2 is mostly used to read photos from files, resize them, and conduct image processing operations.

**Matplotlib:** It is a robust plotting toolkit for Python. It enables users to construct a diverse set of plots, charts, and visualizations. In this software, Matplotlib is used to display images directly within the Jupyter Notebook interface. This enables users to visualize photos without having to save them to disk.

Google.Colab: This module contains features particular to the Google Colab environment. Google Colab is a cloud-based tool that lets you run Python code in a browser. The google.colab module includes tools for interfacing with this environment, such as uploading and downloading files. In this program, files are uploaded from the local system to the Colab environment and downloaded from the Colab environment back to the local system.





**Figure 3.** Object Detection in YOLOv4

To provide visual feedback on the model's performance, we present an output image generated by YOLOv4 for a sample test image from the COCO dataset. The graphic shows correct bounding box predictions and class labels for a variety of items, such as people, cars, and traffic signs. The detection results show that the model can successfully position and identify objects of interest inside complicated situations.

In conclusion, the YOLOv4 object identification research at Google Colab produced promising findings, proving the model's superior performance in terms of accuracy, speed, and efficiency. These findings highlight YOLOv4's practical applicability in a wide range of applications, including surveillance, autonomous driving, and object tracking. The comparison of the YOLOv4 models with SSD and Faster R-CNN, performance of the YOLOv4 in terms of mAP for different class of images are illustrated in Table.1 and Table.2 below.

**Table 1.** Comparison between YOLOv4 and other Models

Model	Accuracy	Efficiency	Model	Learning	Fps	mAP
			Complexity	Rate		
YOLOv4	High, it uses CSPDarknet53 and it has SPP (Spatial Pyramid Pooling). 96.62 is the accuracy.	High, uses single stage object detection.	Low, because of the simplified model design.	0.0001 to 0.01 is the learning rate.	FPS varies from 30 to 60.	mAP score ranges from 40% to 50%.
SSD	High, but not up to the mark of YOLOv4. 86.8 is the accuracy.	Modest, due to its more complex architecture.	Moderate, bounding box encoding and decoding makes it a bit complex.	0.001 to 0.004 is the learning rate.	FPS varies from 50 to 100.	mAP score ranges from 20% to 40%.
Faster R-CNN	High, uses two-stage approach. 84.56 is the accuracy.	Low, due to the increase in inference time.	High, due to its RPN (Region Proposal Network) which is a main component in object detection.	0.001 is the learning rate.	FPS varies from 5 to 15.	mAP score ranges from 30% to 50%

 Table 2. mAP Value for Images in the COCO Dataset

Image Category/Class	mAP (YOLOv4)
Person	0.786
Car	0.789
Motorcycle	0.769
Traffic light	0.854
Bus	0.546

Dog	0.856
Sports ball	0.890
Furniture	0.789
Backpack	0.785

#### 7. Conclusion and Future Works

Finally, our work fully examined YOLOv4's performance for object detection tasks, utilizing Google Colab tools for rapid model training and evaluation. Extensive trials on typical benchmark datasets confirmed YOLOv4's ability to achieve excellent accuracy and real-time performance. The model was resilient across several item categories and demonstrated rapid inference speed, making it ideal for real-world applications that require fast and accurate object detection skills.

The successful implementation of the YOLOv4 object detection project demonstrates its ability to meet a wide range of practical computer vision applications, such as surveillance, autonomous navigation, and object tracking. YOLOv4 provides useful insights into decision-making in a variety of fields by accurately and timely detecting items within complicated scenes.

While our study produced encouraging findings, there are various areas for future research and development. One possible path is to look into lightweight variants of YOLOv4 that are tailored for resource-constrained contexts like embedded systems and mobile devices. Lightweight YOLOv4 models, by reducing model complexity and computing needs, can broaden the scope of object detection technology's applications.

Looking ahead, one key topic for future research is improving pedestrian and transportation networks. By fine-tuning YOLOv4 for pedestrian identification and integrating it with automated traffic monitoring systems, we can help to improve road safety and efficiency in cities. Such developments may result in wiser traffic management, fewer accidents, and improved pedestrian safety, ultimately benefiting society as a whole.

Furthermore, additional research into fine-tuning and optimization strategies can improve YOLOv4's performance and robustness across a variety of datasets and settings. This includes looking into advanced data augmentation procedures, regularization approaches, and new loss functions to increase model generalization and reduce overfitting. In addition, integrating YOLOv4 with powerful object tracking and recognition algorithms can improve its ability to track multiple objects and analyze behavior in dynamic contexts.

Overall, ongoing research and development activities in the field of object detection with YOLOv4 show significant potential for enhancing the state-of-the-art in computer vision and enabling novel applications in a variety of domains. Continued collaboration and study of fresh approaches will help to accelerate the progress of object detection technology and its practical applications in real-world circumstances. Future study will entail investigating novel architectures, increasing dataset diversity, and focusing on specific application domains to improve computer vision capabilities.

#### References

- [1] Bochkovskiy, Alexey, Chien-Yao Wang, and Hong-Yuan Mark Liao. "Yolov4: Optimal speed and accuracy of object detection." arXiv preprint arXiv:2004.10934 (2020).
- [2] Lin, Tsung-Yi, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. "Feature pyramid networks for object detection." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2117-2125. 2017.
- [3] Huang, Jonathan, Vivek Rathod, Chen Sun, Menglong Zhu, Anoop Korattikara, Alireza Fathi, Ian Fischer et al. "Speed/accuracy trade-offs for modern convolutional object detectors." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 7310-7311. 2017.
- [4] Ren, Shaoqing, Kaiming He, Ross Girshick, and Jian Sun. "Faster r-cnn: Towards real-time object detection with region proposal networks." Advances in neural information processing systems 28 (2015).
- [5] Liu, Wei, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C. Berg. "Ssd: Single shot multibox detector." In Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands,

- October 11–14, 2016, Proceedings, Part I 14, pp. 21-37. Springer International Publishing, 2016.
- [6] Redmon, Joseph, Santosh Divvala, Ross Girshick, and Ali Farhadi. "You only look once: Unified, real-time object detection." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 779-788. 2016.
- [7] Redmon, Joseph, and Ali Farhadi. "Yolov3: An incremental improvement." arXiv preprint arXiv:1804.02767 (2018).
- [8] Wang, Chien-Yao, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. "Scaled-yolov4: Scaling cross stage partial network." In Proceedings of the IEEE/cvf conference on computer vision and pattern recognition, pp. 13029-13038. 2021.
- [9] Wang, C. Y., Bochkovskiy, A., & Liao, H. Y. M. YOLOv4 tiny: A simplified version for object detection. arXiv preprint arXiv:2004.10934. (2020).
- [10] Li, X., Zhang, W., & Shi, J YOLOv4-tiny: A fast and small object detection model for real-time applications. arXiv preprint arXiv:2008.10305. . (2020).
- [11] He, Kaiming, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. "Mask r-cnn." In Proceedings of the IEEE international conference on computer vision, pp. 2961-2969. 2017.
- [12] Chen, Kai, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun et al. "MMDetection: Open mmlab detection toolbox and benchmark." arXiv preprint arXiv:1906.07155 (2019).
- [13] Wang, Xiaolong, Ross Girshick, Abhinav Gupta, and Kaiming He. "Non-local neural networks." In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7794-7803. 2018.
- [14] Lin, Tsung-Yi, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. "Focal loss for dense object detection." In Proceedings of the IEEE international conference on computer vision, pp. 2980-2988. 2017.
- [15] Liu, Shu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. "Path aggregation network for instance segmentation." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 8759-8768. 2018.

- [16] He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Deep residual learning for image recognition." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 770-778. 2016.
- [17] Zhang, Xiangyu, Xinyu Zhou, Mengxiao Lin, and Jian Sun. "Shufflenet: An extremely efficient convolutional neural network for mobile devices." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 6848-6856. 2018.

# Author's biography

Mrs.R.Poornima received B.E. Degree in Electronics & Communication Engineering from Anna University, Chennai in 2006. She received her M.E. in Power Electronics and Drives from Anna University, Chennai, in 2012. She is currently working as an assistant professor at SNS College of Technology, Tamilnadu, India. She is pursuing her Ph.D degree in Electrical Engineering at Anna University, Chennai. Her current research interest includes Wireless Communication and Machine Learning.

**Mr.S.A.Thejas** born and brought up in Coimbatore has completed his schooling in a well renounced school, Lisieux Matriculation Higher Secondary School in Coimbatore and currently pursuing his BE in the department of Electronics and Communication Engineering at SNS College of Technology.

**Mr.S.Vinith** born and brought up in Kanyakumari. He completed his schooling in a well renounced school, Government Higher Secondary School in Kanyakumari and currently pursuing his BE in the department of Electronics and Communication Engineering at SNS College of Technology.

**Mr.D.Santhosh** born and brought up in Tiruppur. He completed his schooling in a well renounced school, Government boys Higher Secondary School in Perundurai and currently pursuing his BE in the department of Electronics and Communication Engineering at SNS College of Technology.

**Mr.R.Gowtham** born and brought up in Coimbatore. He completed his schooling in a well renounced school, KG Matriculation Higher Secondary School in Annur and currently pursuing his BE in the department of Electronics and Communication Engineering at SNS College of Technology.