# Object Detection and Deep Learning

## Shra Fatima[1], Maliha Fatima[2], and Wafa Zaidi[3]

[1] Assistant Professor, Department of Computer Science and Engineering, Integral University, Lucknow, India
[2, 3] B. Tech Scholar, Department of Computer Science and Engineering, Integral University, Lucknow, India

Correspondence should be addressed to Shra Fatima    rizvimaliha7@gmail.com

**ABSTRACT-** In this research, the process of implementing deep learning models for real time object detection using YOLOv4 and SSD are explored. The models were trained and evaluated with dataset from publicly available datasets, transfer learning, preprocessing techniques, metrics like mean Average Precision (mAP) and Intersection over Union (IoU). In the case of YOLOv4 the detection accuracy was better than and the speed was faster than SSD, which was optimal for simpler, resource constrained environments. The models were confirmed to be practical via real time testing with webcam and performance was robust under different conditions. The study also explains the pros and cons of each model and makes recommendations on further improvement of efficiency in surveillance, healthcare, and smart systems such as object tracking, edge deployment, and adversarial robustness.

**KEYWORDS-** Object Detection, Deep Learning, YOLOv4, SSD, Real-Time Inference, mAP, IoU, Computer Vision

## I. INTRODUCTION

Object detection is a key task in modern computer vision that treats the task of identification and localization of objects in an image or frame of the video. Object detection, on the other hand, is not like basic image recognition, where the job is only to label the objects present in the image, it also generates bounding boxes and finds several objects at the same time. Such applications require real time response and hence it is very essential.

With a change, the object detection accuracy and efficiency have skyrocketed with deep learning. Haar Cascades or HOG (Histogram of Oriented Gradients) have been the traditional approaches, which were relying heavily on handcrafted features and had poor generalisation to different environments. On the other hand, Convolutional Neural Networks (CNNs) have changed object detection ever since they allowed to learn hierarchical representations from raw pixel data. The state-of-the-art techniques to detect objects in real time with high precision include the YOLO (You Only Look Once) and the SSD (Single Shot Multi-Box Detector) and Faster R-CNN.

Each of the architectures in the deep learning models is unique and correspond to different trade-offs. Yolo makes inference speed trade-off for speed and accuracy by dividing an image into grids and predicting bounding box and class probability in one pass, whereas SSD balances speed and accuracy by using multiple feature maps. However, while faster R-CNN is very precise at proposing regions of interest for classification, it requires more computational demand.

With object detection systems needed to run intelligent automation and AI powered decision making, real time, efficient, and robust algorithms are needed. In this paper, we focused on deep learning algorithms optimization for object detection considering occlusion, scale variation, lighting changes etc that the object may encounter as it appears to us. By the implementation and the comparison of SSD and YOLO, this study shows how these methods could be practically usable and the limits and developments that they can bring into real systems.

## II. LITERATURE REVIEW

In the past couple of years, object detection as a part of computer vision has had an immense role in the growth of deep learning. The early approaches used the hand-crafted features and classifiers such as Viola-Jones using Haar features [8] that worked pretty well in the constrained environment but they failed under the real-world conditions like occlusions and light variable.

The switch to deep learning brought in more of the more adaptable and scalable solutions. In 2014, this was shifted by the Region based Convolutional Neural Network (R-CNN) by [6] that proposed regions of interest (ROIs) and used CNNs for feature extraction. Nevertheless, it was computationally expensive due to its multi stage pipeline. Fast R-CNN [6] and then Faster R-CNN [6] resolved this limitation by integrating the work of a Region Proposal Network (RPN) for improving the efficiency of generating detection candidates and their accuracy.

Object detection utilizing YOLO [5] treating object detection as a regression problem is predicting bounding boxes and class probabilities simultaneously on an image in a single pass. YOLO's real time performance suited it very well for time sensitive applications. [4] further increased the accuracy as well as detection speed of SSD (Single Shot Multi-Box Detector) by utilising multi scale object detection using feature maps of different resolutions.

In recent work speed has not been sacrificed at the expense of accuracy. YOLOv4 and YOLOv5 improves backbone architecture, data augmentation and training strategies to achieve the state of art level on the COCO and Pascal VOC benchmarks [1].

While progress has been made in detection, dense scenes present great challenges as well as improving in generalization from domain to domain. Therefore, researchers are moving to the hybrid models, transfer learning, and lightweight architectures for deployment on the edge devices[7]

## III. OBJECTIVES

The following are the key objectives on which this research is based:

- This is to accurately identify and locate key points on multiple objects in an image.
- This is to build and train deep learning models for real time object detections using the SSD and YOLO architectures.
- A study in order to evaluate model performance using precision-recall, IoU, and mAP metrics.

## IV. METHODOLOGY

This work develops and evaluates object detection systems based on structured deep learning pipeline by using SSD and YOLO architectures. The methodology is divided into five main phases that can guarantee the effective transformation of raw image data into actionable object localization.

### A. Data Collection and Annotation

I sourced labelled datasets using open repositories such as COCO and Pascal VOC. Bounding boxes and object classes were drawn using annotation tools such as Labelling. The dataset possessed high intra class variability and diverse condition such that it could generalize.

### B. Image Preprocessing

Image preprocessing consisted in resizing to 300×300 (SSD) and 416×416 (YOLO) resolution, grayscale conversion when necessary and normalization to speed up convergence. Step in response to remove the overfitting and enhanced in dynamic conditions was augmented techniques, rotation, flipping, and scaling.

### C. Model Selection

YOLOv4 and SSD were chosen for implementation because they are a demonstrated trade-off between detection speed and accuracy. Since YOLOv4 is preferred for real time efficiency and SSD is suitable for both small and large objects, it is chosen. Since it performs so well in dense object scenarios, Faster R-CNN was used as a benchmark.

### D. Training and Evaluation

Using the TensorFlow framework and GPU accelerations, transfer learning of the pre-trained weights was used to train. Mean Average Precision (mAP) and Intersection over Union (IoU) were used as evaluation metrics since they are key metrics to evaluate object localization accuracy and the overall performance.

### E. Deployment

Real-time webcam detection via webcam was done with the trained model deployed in a Python based interface using OpenCV. Inference speed and responsiveness under uncontrolled and practical condition were validated through their deployment.

## V. IMPLEMENTATION

The object detection project implementation phase had as its objective the integration of deep learning models with the ability of real time inference. This was done using modern frameworks, software environments and high-performance hardware to train, validate and deploy SSD and YOLO models to be used for object recognition.

- Frameworks and Tools Used-

The model was mainly developed on TensorFlow and Keras, and OpenCV was used for image acquisition, processing. Experimentation and benchmarking were done using PyTorch in parallel. PyCharm and Jupyter Notebook were the IDEs used for scripting and visualisation.

- Hardware and Software Environment-

To test the system, we used a workstation with an Intel i7 processor, 16 GB RAM, and GPU to accelerate training (an NVIDIA GTX GPU was optional). Webcam integration was performed for real time testing, which made the system practically applicable. The implementation in this work was carried out in a Python environment.

- Implementation Workflow-

### A. Dataset Preprocessing

Images were resized, normalized and augmented to increase training diversity. Horizontal flipping, cropping, and colour normalization were used in many techniques to improve the model's robustness to variation in the environment.

### B. Model Training

The labelled data is used to train the SSD and YOLOv4 models applying transfer learning with the pre trained weights on COCO. The hyperparameters (learning rate and batch size) used for training the models were modified to fine tune the models and increase the detection accuracy and decrease overfitting.

### C. Real-Time Detection

The models were then exported and put into a Python application to access webcam input. Real time frames were captured, passed through the detection pipeline and had bounding boxes and confidence scores dynamically annotated over them.

### D. Output Annotation

It included each detection result comprising of a labelled bounding box and class prediction with confidence levels. The result was rendered directly on the video stream, for live visualization and verification of the object localization performance.

Thus, this implementation validates the feasibility of deploying deep learning models for the deployment of real-world, real-time object detection scenarios.

## VI. RESULT AND DISCUSSION

The object detection system, made using SSD and YOLOv4, is evaluated on the basis of detection accuracy, inference speed, and robustness under different conditions. Both models showed high precision, YOLOv4 was slightly better than SSD in real time detection tasks. SSD and YOLOv4 achieved 85.6% and 91.2% average precision

respectively on a custom dataset generated from COCO and Pascal VOC. The average detection time of the system was 0.16–0.18 seconds per frame, which was enough to validate its suitability for real time applications.

Tested under such conditions as low light and partial occlusion, YOLOv4 showed better scalability and robustness compared to the other methods mentioned above. However, SSD was easier to integrate and was perfectly suited to resource constrained environments. The system was shown to be practically efficient through its webcam-based deployment, which always detected and annotated objects in live video streams with minimal latency.

Some of these included occasional misclassification in cluttered scenes and limited accuracy on small, overlapping objects, typical limitations of single shot detectors. However, data augmentation and transfer learning were used by the system to improve its generalisation.

The research provides data that confirms the feasibility and effectiveness of using deep learning object detection approach for real-time systems, while identifying areas for improvement including, among other things, pruning of the model, multi-object tracking and optimizing deploys for edge.

Future development can be future for accuracy in detection small and overlapping object, integration multi object tracking also optimization model for edge device. Further improvements in performance, scalability, as well as security will be made by investigating federated learning, adversarial defence, and domain adaptation for such real-time applications such as smart surveillance, autonomous systems, and healthcare.

## VII. CONCLUSION

The real time object detection problem has been successfully demonstrated on the practical application of deep learning-based models i.e. YOLOv4 & SSD. High detection accuracy and responsiveness in the models were ensured through the use of robust datasets, transfer learning and structuring of preprocessing. YOLOv4 was found to be superior to SSD with its high speed of inference and the ability to adapt to real world conditions, but SSD can be relied on for good, less complex datasets on less sophisticated hardware. Objective performance validation was made through the use of metrics such as mAP and IoU. The limitation of single shot models was shown, although a nice application of them, in that they were easy to implement and fast in nature but as realized the resulting false positive when there is clutter in the background as well as the problem of small object detection was identified. However, model fine tuning and augmentation strategies increased system generalization as well as scalability.

Overall, this study provides evidence validating deep learning's disruptive force on the field of computer vision and serves as a guiding point for pursuing future work at the edge in deployment, and robustness to adversarial attacks, and tracking. The system provides a solid basis for the future AI based automation in the fields of surveillance, transportation, healthcare and smart environment.

## CONFLICTS OF INTEREST

The authors declare that they have no conflicts of interest.

## REFERENCES

[1] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "YOLOv4: Optimal speed and accuracy of object detection," *arXiv preprint arXiv:2004.10934*, Apr. 2020. Available from: https://arxiv.org/abs/2004.10934

[2] R. Girshick, "Fast R-CNN," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Santiago, Chile, Dec. 2015, pp. 1440–1448. Available from: https://doi.org/10.1109/ICCV.2015.169

[3] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Columbus, OH, USA, Jun. 2014, pp. 580–587. Available from: https://doi.org/10.1109/CVPR.2014.81

[4] W. Liu *et al.*, "SSD: Single shot multibox detector," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, Amsterdam, Netherlands, Oct. 2016, pp. 21–37. Available from: https://doi.org/10.1007/978-3-319-46448-0_2

[5] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Las Vegas, NV, USA, Jun. 2016, pp. 779–788. Available from: https://doi.org/10.1109/CVPR.2016.91

[6] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," in *Adv. Neural Inf. Process. Syst. (NeurIPS)*, Montreal, Canada, Dec. 2015, pp. 91–99. Available from: https://shorturl.at/no0xn

[7] M. Tan, R. Pang, and Q. V. Le, "EfficientDet: Scalable and efficient object detection," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Seattle, WA, USA, Jun. 2020, pp. 10781–10790. Available from: https://shorturl.at/VfsFN

[8] P. Viola and M. Jones, "Rapid object detection using a boosted cascade of simple features," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Kauai, HI, USA, Dec. 2001, vol. 1, pp. 511–518. Available from: https://doi.org/10.1109/CVPR.2001.990517