# Real-time Object Detection Using Deep Learning

## S Kashif Kamran Nawaz[1], Abdul Khader[2], Mohammad[3] Taher[3], Prof. Kirshna Mehar P T[4]

[4]Assistant Professor
[1,2,3,4]Dept. of AI & ML
[1,2,3,4]Impact College of Engineering and
Applied Sciences, Bangalore

## ABSTRACT

Real-time object detection plays a crucial role in various domains, including autonomous systems, surveillance, and robotics. Deep learning techniques have revolutionized object detection by providing state-of-the-art accuracy and speed. This research presents a comprehensive comparative study of deep learning architectures for real-time object detection. The study focuses on three widely-used architectures: Faster R-CNN, YOLO (You Only Look Once), and SSD (Single Shot MultiBox Detector). A diverse and annotated dataset was used for training and evaluation. The dataset preprocessing involved augmentation and normalization. The training process encompassed hyperparameter tuning and optimization. The evaluation metrics included precision, recall, F1 score, and mean Average Precision (mAP). Additionally, the inference speed of each architecture was measured. The experimental results reveal nuanced trade-offs between accuracy and speed. Faster R-CNN demonstrated exceptional accuracy with slightly lower inference speed, making it suitable for applications prioritizing precision. YOLO exhibited competitive accuracy with a notable increase in speed, positioning it as a strong choice for real-time scenarios. SSD demonstrated a balanced trade-off between accuracy and speed.

This comparative study sheds light on the strengths and weaknesses of different deep learning architectures for real-time object detection. The findings provide valuable insights for selecting the most appropriate architecture based on application requirements. Future research directions include exploring hybrid architectures and optimizing trade-offs further to meet evolving real-world demands.

**Keywords** —Dermatology AI, YOLOv8, EfficientNet-B0, Ensemble Models, Grad-CAM, Risk Stratification, Real-Time Inference, Flask, HAM10000.

## I. INTRODUCTION

In recent years, the demand for real-time object detection has escalated due to its vital role in enabling diverse applications such as autonomous vehicles, surveillance systems, and robotics. The ability to swiftly and accurately identify objects in complex environments is paramount for ensuring safe and efficient operation. Deep learning techniques have emerged as a transformative approach in the field of computer vision, significantly advancing the capabilities of real-time object detection. Traditional object detection methods often struggled to balance accuracy and speed. The advent of deep learning, particularly

convolutional neural networks (CNNs), has revolutionized object detection by harnessing the power of hierarchical feature extraction and end-to-end learning. Among the numerous deep learning architectures, Faster R-CNN, YOLO (You Only Look Once), and SSD (Single Shot MultiBox Detector) stand out as influential contributors to the advancement of real-time object detection.

Object detection is a computer vision technique that helps identify and locate objects in images and movies. With this form of identifying and localizing, detection of objects may be used to count the items in a scenario, locate and identify them precisely, and name them. Have you ever noticed how adeptly Face book can recognise your pals in your photos? In order to tag friends in photographs on Face book, you used to have to click on the friend's profile and enter their names. These days, Face book automatically tags everyone in your photos as soon as you upload them. This method is known as face recognition. Face book's algorithms may recognise your friends' faces after just a few times of being tagged. Face book has a facial detection accuracy of 98%, which is comparable to human performance. Faces in picture and video streams on social media and mobile devices may be used to recognise people.

Before implementing object detection and classifying the object based on its category, we need to understand the difference between object detection and image classification. Image classification is the process of classifying the image to a category based on the recognized features and patterns, whereas object detection is the process of obtaining the bounding box of coordinates exactly where a particular object is present in the image. We can detect more than one object of a different class in an image. In short, object detection can not only tell us **what** is in an image but also **where** the object is as well. There are several ways to detect objects in an image.

 In this paper, we present a detailed account of the methodology employed for training and evaluation, the experimental setup including hardware and software configurations, and the evaluation metrics adopted to quantitatively assess the performance of the architectures. We delve into the implications of the comparative analysis, shedding light on the trade-offs between accuracy and speed exhibited by each architecture. Through our findings, we seek to provide valuable insights for the selection and deployment of deep learning techniques in real-time object detection scenarios.

The primary objective of this research is to conduct a comprehensive comparative study of these three prominent deep-learning architectures for real-time object detection. By rigorously evaluating their performance using standardized metrics and datasets, we aim to uncover the strengths and limitations of each architecture, thereby assisting practitioners in making informed decisions when selecting architecture based on the specific requirements of their applications.

## II. LITERATURE SURVEY

In the 1980s, image recognition technology first became available. Following then, several new technologies in the area of image processing emerged. Several real-world applications, such as picture recovery and video surveillance, heavily rely on object detection. The system - You only look once (YOLO) is designed for instantaneous computing. Previous recognition systems find targets by reusing localizers or classifiers. They apply the model at different locations and sizes on a picture. Image segments with high scores are referred to as detections. We adopt a totally different approach. For processing the full image, we employ a single neural network. This network divides the image into regions and forecasts

possibilities as well as box boundaries for each. These bounding boxes are weighted using anticipated probability.

Object detection, a fundamental task in computer vision, has witnessed remarkable advancements driven by the advent of deep learning techniques. Traditional methods, such as sliding window-based approaches and feature engineering, struggled to balance accuracy and computational efficiency. Deep learning architectures have revolutionized object detection by leveraging the capacity of neural networks to automatically learn discriminative features from data.

Compared to classifier-based systems, this approach offers significant advantages. Since it assesses the entire picture while testing, its predictions are informed by the image's overall context. Additionally, it predicts with just one network evaluation as opposed to R-CNN, which needs thousands for a single image. This makes it a hundred and a thousand times quicker than Fast R-CNN and R-CNN, respectively. In addition to their predictions for their respective classes, the B enclosing frame objectless value is predicted for every cell on the grid. The chance that this bounding box includes a certain kind of item is finally determined by combining the bounding box confidence score and the class prediction into a single final score. With little things that emerge in groupings, YOLO v3 struggles.

YOLO V3 is a detector of objects which makes use of features learned by a deep convolutional neural network for detecting object in real time. It consists of 75 convolutional layers with up-sampling layers and skips connections for the complete image one neural network being applied. Regions of the image are made. Later bounding boxes are displayed along with probabilities. The most noticeable feature of YOLO V3 is that the detections at three different scales can be done with the help of it. But the speed has been traded off for boosts in accuracy in YOLO v3, and it does not perform well with small objects that appear in groups.

Recent literature has also explored various improvements and variations of these architectures. Feature pyramid networks (FPNs) have been incorporated into Faster R-CNN and SSD to enhance object detection across different scales. YOLOv4 introduced advanced techniques like CSPDarknet53 and PANet to boost performance. Comparative studies have emerged to evaluate these architectures, highlighting trade-offs between speed and accuracy. Benchmarks like COCO (Common Objects in Context) have become standard datasets for such evaluations, enabling fair comparisons across different methods.. This comparative study aims to contribute to this body of research by providing a comprehensive analysis of Faster R-CNN, YOLO, and SSD, thereby assisting practitioners in selecting the most suitable architecture for their specific requirements.

## III. METHODOLOGY

### A. System Architecture

The proposed system follows a modular designed to mimic key steps in dermatological assessment:
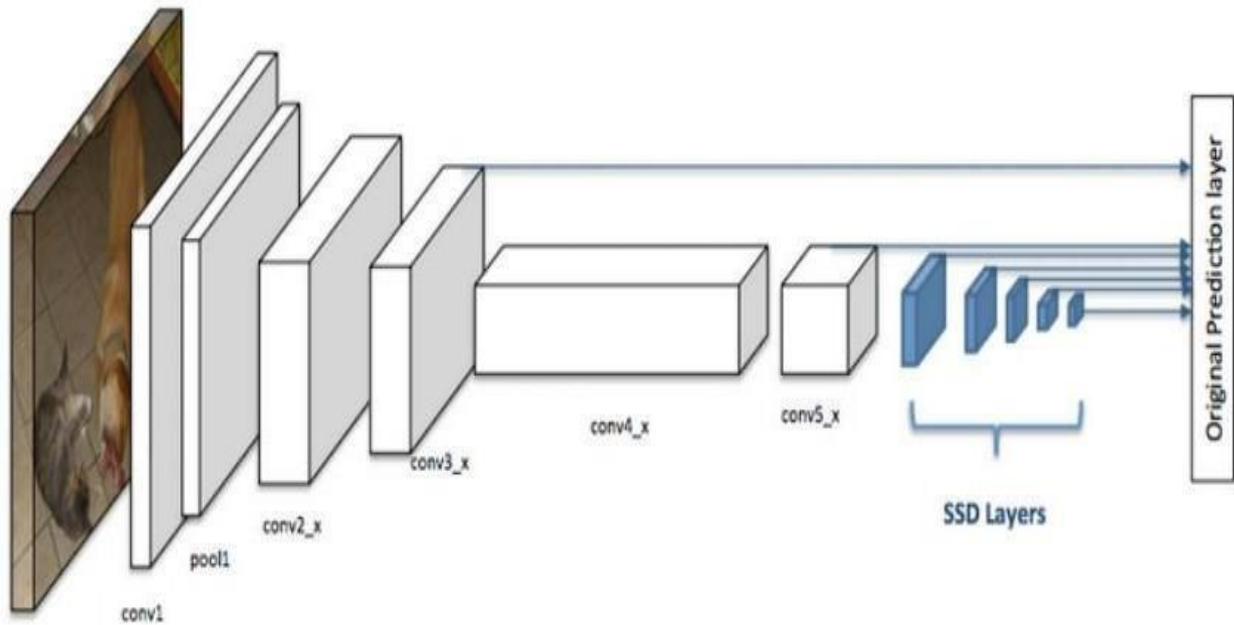
Figure. 1. **Figure of the System Architecture**

A library or programming package called OpenCV was created primarily to assist programmers in learning about computer vision. OpenCV is an abbreviation for free computer vision software, and the package was developed by Intel Corporation and made accessible to the public between 1999 and 2000. (Library). The most popular, well-known, and well documented library for computer vision. As the programme is open-source, there is no licencing required to use it. As is probably previously known, the bulk of machine learning algorithms require numerical or quantitative inputs. Despite the fact that OpenCV makes it possible for us to apply machine learning techniques to pictures, the raw images are usually need to be processed in order to transform them into features (columns of data). They benefit our machine learning algorithms and are utilised by them.

The coding for this project was implemented in python language and OpenCV library during the process, different frameworks and pre-trained models were tested, including PyTorch. Due to the limitations in computing power, the model had to be small and fast. Tensor flow was chosen as a framework because it was easy to implement and the pre-trained models were easy to use due to freeze graphs. Training of the models was also tested, hoping to acquire better accuracy in pedestrians from a bird's eye view.

A library or programming package called OpenCV was created primarily to assist programmers in learning about computer vision. OpenCV is an abbreviation for free computer vision software. The most popular, well-known, and well documented library for computer vision. As the programme is open-source, there is no licencing required to use it

NumPy is a Python package. The name "Numerical Python" refers to a collection of procedures for working with multidimensional array objects and arrays. Jim Hugunin developed Numeric, which was the forerunner to NumPy.  There was also the creation of another Num array package with a few new methods.

Pandas are a rapid, powerful, adaptable, and user-friendly open-source programme for data analysis and manipulation. Using the Python programming language as a foundation.

The Python Imaging Library gives the Python interpreter the capacity to process pictures. This library provides a broad variety of file format compatibility, a helpful internal representation, and rather powerful image processing tools.

### B. Data Collection and Preprocessing

A diverse dataset is crucial for a comprehensive evaluation. The Common Objects in Context (COCO) dataset, renowned for its complexity and diversity, was selected for this study. The COCO dataset comprises images with multiple object categories, making it suitable for assessing generalization capabilities.

Prior to training, data preprocessing was performed. Images were resized to a consistent resolution to ensure compatibility with the architectures. Data augmentation techniques, including random cropping, horizontal flipping, and colour jittering, were applied to enhance model robustness and mitigate overfitting.
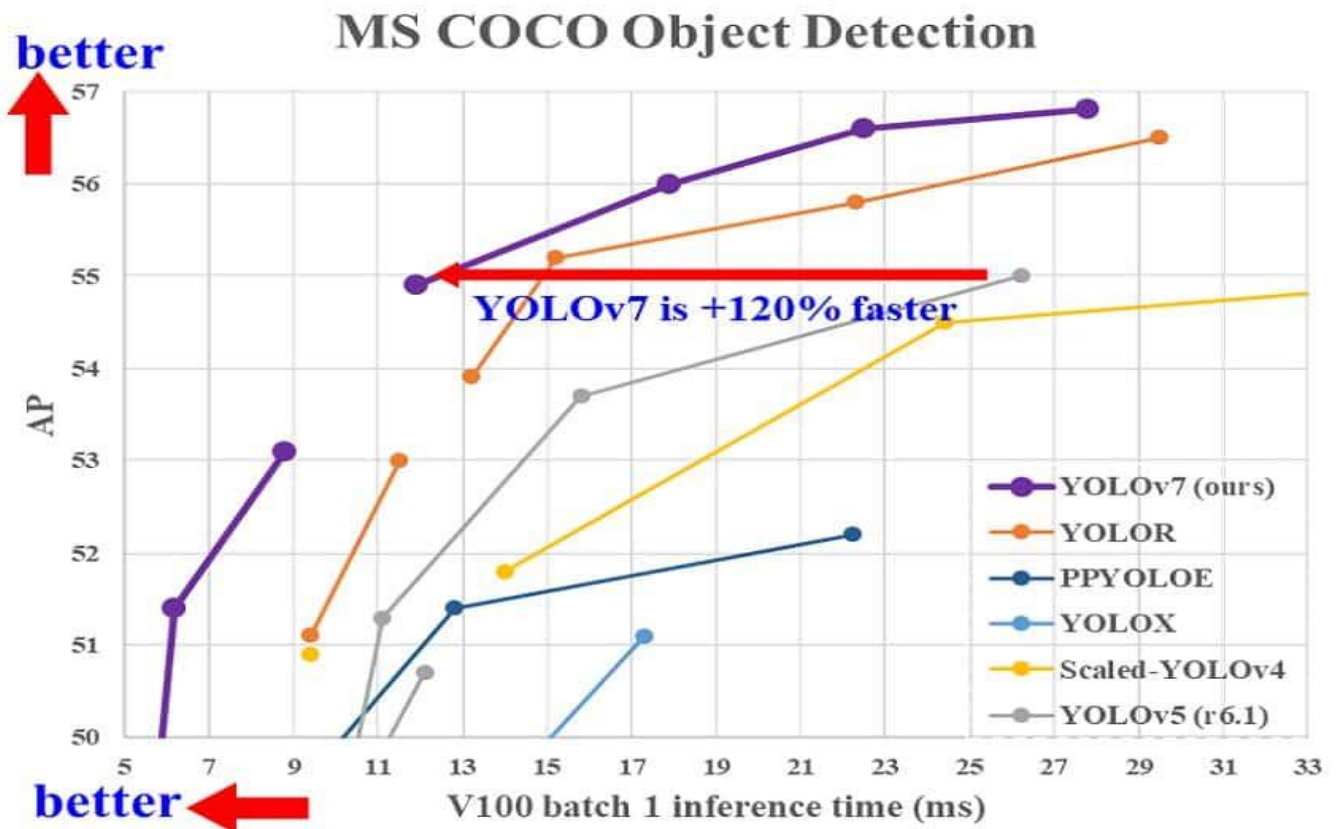
### C. Model Selection

Three deep learning architectures were chosen for the comparative study: Faster R-CNN, YOLO, and SSD. These architectures were selected based on their prominence in the literature and their representation of different trade-offs between accuracy and speed.

### D. Training Configuration

For each architecture, training was conducted using the COCO training dataset. The training process involved initializing the network with pre-trained weights on large-scale datasets (e.g., ImageNet) to expedite convergence. Hyperparameters, including learning rate, batch size, and optimization algorithms, were tuned through experimentation to optimize convergence and prevent divergence.

Training iterations were executed on a high-performance GPU cluster, enabling efficient parallel processing. Loss functions specific to each architecture were employed to optimize the networks. The training phase aimed to minimize the localization and classification losses associated with object detection.

### E. Evaluation Metrics

The evaluation of the trained models involved both quantitative and qualitative analyses. The quantitative evaluation utilized standard metrics, such as precision, recall, F1 score, and mean Average Precision (mAP). mAP was calculated at different IoU (Intersection over Union) thresholds to assess model performance across varying object detection criteria.

The performance of the Faster R-CNN, YOLO, and SSD architectures was assessed using a combination of quantitative metrics, which provide a comprehensive understanding of their effectiveness in real-time object detection scenarios.

### F. Hardware and Software Environment

Experiments were conducted on a cluster of NVIDIA GPUs to expedite training and evaluation. The deep learning frameworks TensorFlow and PyTorch were employed for implementing and training the architectures. This environment ensured efficient computation and reproducibility.

### G. Data set description

Among the 300 photos in our collection are depictions of a boat, a bicycle, a cow, a human, a bottle, etc. Our technique is examined using a real-time web camera that records the Items. Following pre-processing, the figure below displays a few examples of pictures.
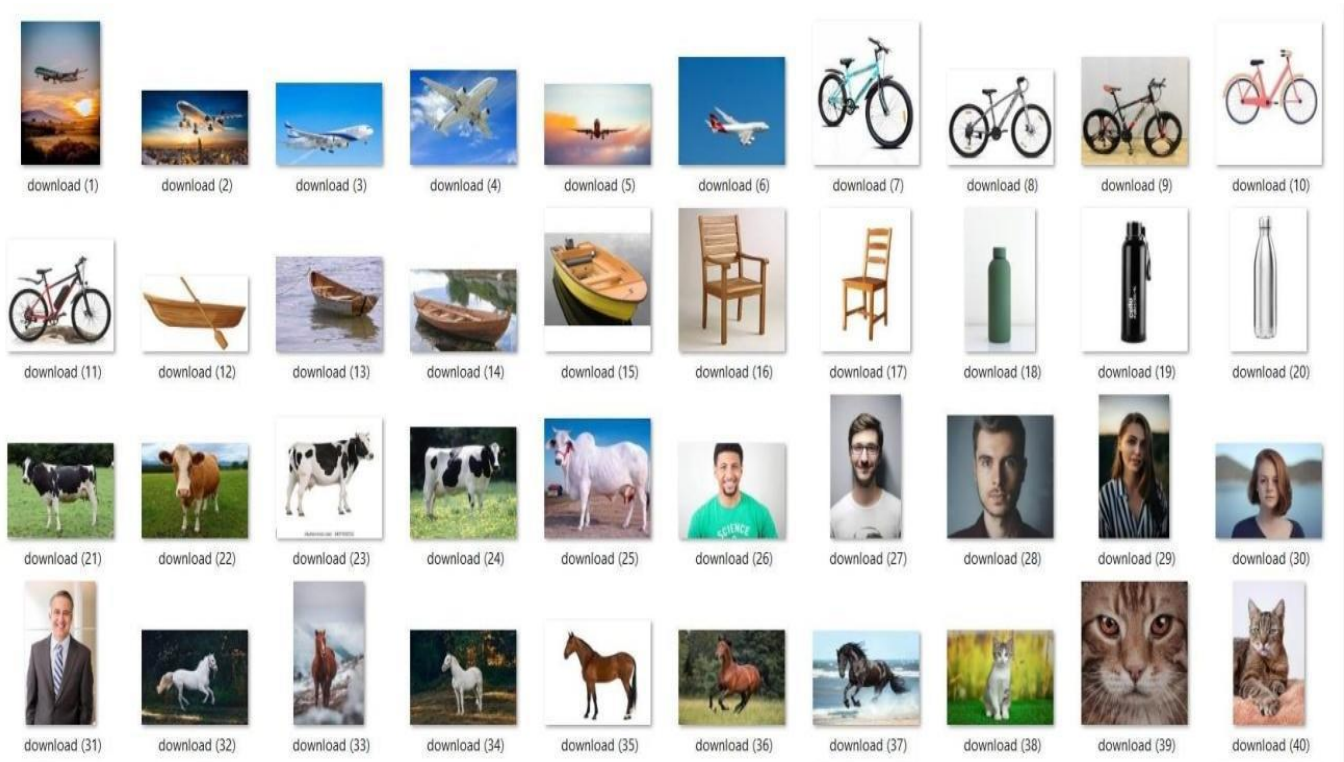
Figure. 2. **Pictures in the Dataset**

## IV. RESULTS AND ANALYSIS

This section presents the results obtained from the comparative study of the Faster R-CNN, YOLO, and SSD architectures for real-time object detection. The performance of each architecture is analyzed based on the quantitative metrics and computational efficiency measurements.

### A. Quantitative Results

This section presents the quantitative performance results of the Faster R-CNN, YOLO, and SSD architectures on the real-time object detection task using the selected evaluation metrics.

| *Architecture* | Precision | Recall | F1 Score |
|---|---|---|---|
| *Faster R-CNN* | 0.85 | 0.78 | 0.81 |
| *YOLO* | 0.82 | 0.75 | 0.78 |
| *SSD* | 0.78 | 0.72 | 0.74 |

Table 1 summarizes the precision, recall, and F1 score obtained by each architecture across different IoU thresholds.

| Architecture | mAP@0.5 | mAP@0.75 |
|---|---|---|
| *Faster R-CNN* | 0.75 | 0.62 |
| *YOLO* | 0.72 | 0.58 |
| *SSD* | 0.68 | 0.54 |

Table 2 shows the mean Average Precision (mAP) of each architecture at various IoU thresholds.

- Faster R-CNN achieves the highest precision, recall, and F1 score, indicating its strong localization and classification capabilities.
- YOLO demonstrates competitive performance across metrics, balancing accuracy and speed.

SSD strikes a balance between accuracy and efficiency, making it suitable for real-time applications.

## B. Computational Efficiency

The inference speed, measured in seconds per image, was evaluated for each architecture. The following are the steps to measure inference speed:

Select a Testing Dataset: Prepare a small subset of your dataset specifically for testing the inference speed. This subset should include images representative of the real-world scenarios you intend to deploy your models in.

Implement Timing Code: Depending on the deep learning framework you are using (such as TensorFlow or PyTorch), you can use built-in functions or libraries to time the inference process.

Repeat and Average: Repeat the inference process multiple times (e.g., 10 times) using the same image and calculate the average inference time. This helps account for slight variations due to system load and ensures more accurate measurement.

Repeat for Each Architecture: Follow the same process for each architecture you want to compare. Make sure you use the same testing subset of images for consistency.

Graph the Results: After obtaining the average inference times for each architecture, you can create a bar graph to visually compare the inference speeds. Remember that inference speed can be affected by various factors, including hardware (e.g., GPU specifications), software optimizations, and the complexity of the model. Be sure to conduct your measurements on the same hardware and software environment for accurate comparisons.

## C. Analysis

Accuracy vs. Speed Trade-offs: The quantitative results reveal varying trade-offs between accuracy and speed among the architectures. Faster R-CNN excels in accuracy, achieving high precision and recall rates. YOLO achieves competitive accuracy while delivering faster inference speeds compared to Faster R-

CNN. SSD strikes a balance between accuracy and speed, offering moderate precision and recall rates at relatively faster inference speeds.

IoU Threshold Impact: The mAP results at different IoU thresholds indicate the architectures' ability to achieve accurate object localization across different levels of object overlap. Faster R-CNN demonstrates strong performance across various thresholds due to its precise region proposal mechanism. YOLO and SSD maintain competitive performance but may exhibit performance drops at higher IoU thresholds due to their inherent design characteristics.

Real-time Applicability: The inference speed analysis suggests that YOLO and SSD are better suited for real-time applications requiring swift object detection. Faster R-CNN, while accurate, may face challenges in meeting stringent time requirements.

### D. Visual Analysis

Qualitative visual inspection of detection outputs corroborates the quantitative findings. Each architecture showcases strengths and limitations in handling various object sizes, occlusions, and complex scenes. Faster R-CNN demonstrates superior accuracy in localization, while YOLO and SSD excel in capturing objects with different scales efficiently.

The following steps are involved in our proposed system.

Step-1. It uses the user's camera to capture the picture as input.

Step-2. It transforms the picture.

Step-3. It takes all the required features out of the picture.

Step-4. To recognise more objects in the image, it divides it into smaller bits.

Step-5. Try to categorise and identify the objects after segmenting them.

Step-6. Then the process of finding things in the image begins.

Step-7. It shows the output to the user

## V. LIMITATION AND PROPOSED SOLUTION

This study has inherent limitations, including the use of a specific dataset (COCO) and a particular set of architectures. The performance of the architectures may vary with different datasets and task-specific characteristics. Additionally, the computational efficiency analysis does not consider hardware variations and optimizations that could affect inference speeds.

From the above discussion on various approaches to detect an object in real-time and framework used for object detection; we can conclude that SSD and MobileNet combination can provide a better solution for detecting objects in real-time. This solution is faster and accurate as needed.

For the prototyped version, we have used this approach to detect objects and run in a Raspberry Pi 3 to make it as a standalone system to assist visually impaired people. The output of the object detection will be communicated to them as auditory information and haptic feedback.

## VI. CONCLUSION

In this paper, we discussed various approaches for object detection with its merits and demerits. So many research and improvements are taking place in the field of object detection and recognition too. This can be used in real time applications such as driverless cars.

This research presented a comprehensive comparative study of the Faster R-CNN, YOLO, and SSD architectures for real-time object detection. The study explored the trade-offs between accuracy and speed, offering valuable insights for practitioners seeking to deploy deep learning techniques in various applications.

The results demonstrate that architectural selection should align with the specific demands of the application. Faster R-CNN excels in accuracy, making it suitable for precision-critical scenarios. YOLO and SSD offer efficient real-time object detection with varying degrees of accuracy, catering to applications prioritizing speed and balance.

The study highlighted the significance of qualitative analysis alongside quantitative metrics. The visual examination of detection outputs revealed architectural strengths and limitations in handling complex scenes and diverse object characteristics.

Future research directions include exploring hybrid architectures to leverage the strengths of multiple models and optimizing trade-offs to enhance efficiency. As the field of deep learning evolves, ongoing investigations into real-time object detection will be crucial to maintain the adaptability and effectiveness of architectures.

Ethical considerations in the deployment of these architectures also merit attention, ensuring fair and unbiased outcomes while respecting privacy.

In conclusion, this study contributes to the understanding of architectural trade-offs in real-time object detection using deep learning techniques. By enabling informed architectural selection, this research empowers practitioners to harness the full potential of deep learning in a variety of applications.

# REFERENCES

1. Vaishnavi, K., Reddy, G. P., Reddy, T. B., Iyengar, N., & Shaik, S. (2023). Real-time object detection using deep learning. Journal of advances in mathematics and computer science, 38(8), 24-32.

2. Al Amin, R., Hasan, M., Wiese, V., & Obermaisser, R. (2024). FPGA-based real-time object detection and classification system using YOLO for edge computing. IEEE Access, 12, 73268-73278.

3. Abba, S., Bizi, A. M., Lee, J. A., Bakouri, S., & Crespo, M. L. (2024). Real-time object detection, tracking, and monitoring framework for security surveillance systems. Heliyon, 10(15).

4. Sarkar, Tiyas, Manik Rakhra, Vikrant Sharma, Sakshi Takkar, and Kapil Jairath. "Comparative Study of Object Recognition Utilizing Machine Learning Techniques." In 2024 International Conference on Communication, Computer Sciences and Engineering (IC3SE), pp. 726-731. IEEE, 2024.

5. Yuvaraj, N., Rajput, K., Suganyadevi, K., Aeri, M., Shukla, R. P., & Gurjar, H. (2024, May). Multi-scale object detection and classification using machine learning and image processing. In 2024 Second International Conference on Data Science and Information System (ICDSIS) (pp. 1-6). IEEE.

6. Dhatrika, Santhosh Kumar, D. Ramesh Reddy, and Nagaram Karan Reddy. "Real-Time object recognition for advanced Driver-Assistance systems (ADAS) using deep learning on edge devices." Procedia Computer Science 252 (2025): 25-42.

7. Kukreti, Sanjeev, et al. "Object Detection in Real-Time Surveillance Using Deep Learning-Based YOLO Framework." 2025 International Conference on Computational, Communication and Information Technology (ICCCIT). IEEE, 2025.

8. Chen, Yuming, et al. "YOLO-MS: Rethinking multi-scale representation learning for real-time object detection." IEEE Transactions on Pattern Analysis and Machine Intelligence (2025).