Course Project Report

on

# Object Detection with Voice Feedback using YOLO v8 and gTTS

Submitted by

## Aryan Shirke (21BCS111)



**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**
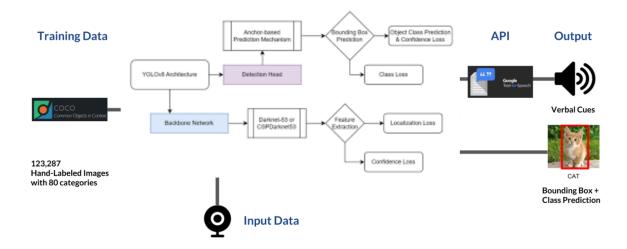
**INDIAN INSTITUTE OF INFORMATION TECHNOLOGY DHARWAD**

21/05/2024

# Contents

# 1 Introduction

Object Detection is a critical field of Computer Vision focused on identifying instances of semantic objects within images and videos, marked by bounding boxes. This project integrates object detection with voice feedback, aiming to provide positional information of detected objects through auditory responses. This application can significantly aid visually impaired individuals by verbally describing the location of objects in their surroundings.

# 2 Methodology



## 2.1 Training Data

The model is trained using the Common Objects In Context (COCO) dataset, a large-scale object detection, segmentation, and captioning dataset. This dataset contains a wide variety of labeled images, making it ideal for training robust object detection models.

## 2.2 Model

We employ the You Only Look Once (YOLO) algorithm, specifically the YOLO v8 model. YOLO v8 is a state-of-the-art, real-time object detection system that uses a single neural network

to predict bounding boxes and class probabilities directly from full images in one evaluation. The model architecture is implemented using the Darknet framework, and we utilize pre-trained weights available in a 200+ MB file to enhance our system's performance.

## 2.3  Input Data

The system processes static images, which are fed into the trained YOLO v8 model to detect objects. The model outputs the class predictions and bounding box coordinates for each detected object.



## 2.4  API Integration

Class predictions, such as "cat," are appended with positional information based on the bounding box coordinates. The position is described using terms like "top," "mid," "bottom," "left," "center," and "right." This annotated text is then sent to the Google Text-to-Speech (gTTS) API to generate corresponding voice feedback.

# 3 Implementation

## 3.1 Data Preparation



We start by loading the COCO dataset and preparing it for training. This involves normalizing images and creating annotations compatible with the YOLO v8 format.

## 3.2 Model Setup

Using the Python 'cv2' package, we configure the Darknet model with our 'yolov8.cfg' file and load the pre-trained weights. The model is then fine-tuned on our prepared dataset.

## 3.3 Detection and Feedback

For each image, the YOLO v8 model detects objects and outputs their classes and coordinates. These outputs are processed to generate positional descriptions, which are converted to speech

using the gTTS package. The final voice feedback provides clear and concise information about the object's location in the image.

# 4 Results and Discussion

The object detection system successfully identified and provided voice feedback for various objects in static images. For example, when a cat was detected at the bottom left of an image, the voice feedback accurately described it as "bottom left cat." This real-time auditory feedback can greatly assist visually impaired users by describing their surroundings through sound.



Now, let's discuss the specific objects detected in different areas of the image:

- **Top Left**: A small, white rectangular building with a flat roof.

- **Top Center**: A large, light blue building with a flat roof. There are several dark rectangular objects on the roof, possibly air conditioning units.

- **Top Right**: A large, white building with a flat roof. There is a blue water tower on the roof, with a metal stand.

- **Middle Left**: A small, white rectangular building with a flat roof, partially obscured by a larger white building in the center.

- **Center**: A large, white building with a flat roof that takes up a significant portion of the image. There is a blue water tower on the roof with a metal stand. There is also a small, black rectangular object on the roof, near the water tower.

- **Middle Right**: A large white building with a flat roof. There is a dark rectangular object near the top, and a loading dock with a red container visible on the side.

- **Bottom Left**: A four-lane intersection. Several cars are visible driving through the intersection, in various colors.

- **Bottom Center**: A four-lane intersection. Several cars are visible driving through the intersection, in various colors. There is a white van partially obscuring the view of some of the cars.

- **Bottom Right**: A four-lane intersection. Several cars are visible driving through the intersection, in various colors. There is a bus visible traveling through the intersection, and a white car stopped at a red light.

These detailed descriptions enhance the understanding of the scene and demonstrate the system's capability to provide informative auditory feedback.

# 5 Conclusion

This project demonstrates a practical application of combining object detection with voice feedback to aid visually impaired individuals. By leveraging the YOLO v8 model and gTTS, we provide an efficient and effective solution for real-time object recognition and description. Future work could involve extending this system to handle video inputs and improving the accuracy of positional descriptions.

# 6   References

# Bibliography

[1] Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., ... & Zitnick, C. L. (2014). Microsoft COCO: Common Objects in Context. In *European Conference on Computer Vision* (pp. 740-755). Springer.

[2] Redmon, J., & Farhadi, A. (2018). YOLOv8: An Incremental Improvement. *arXiv preprint arXiv:1804.02767.*

[3] Google Text-to-Speech API,