# International Journal of Research Publication and Reviews

Journal homepage: www.ijrpr.com  ISSN 2582-7421

# Voice Assistant System for Blind People Using Image Processing

*Sk. Moin Madhar[1], T. Bhanu Chandu[2], Y. Supriya[3], T. Manikanta[4], Y. Paranthapa[5], V. Surendra Babu[6],*

[1,2,3,4,5]Department of Electronics and Communication Engineering, SIR CR REDDY College of Engineering
[6]Assistant Professor, Department of Electronics and Communication Engineering, SIR CR REDDY College of Engineering

**ABSTRACT—**

This paper presents a system to assist visually impaired individuals with object detection and recognition using computer vision-based solutions. The system utilises a Raspberry Pi 3b+ and a USB camera to capture images, which are processed using the YOLOv3 algorithm for object recognition. The system is also equipped with a voice command interface, allowing the user to locate objects within their environment. The results of the study demonstrate the feasibility of the system in detecting and recognising various objects, with a high level of accuracy. Future work includes the integration of a panic button to alert emergency contacts in case of an emergency.

*Keywords— (Visual impairment, computer vision, object detection, YOLOv3 algorithm, Raspberry Pi, USB camera, voice command interface.).*

## 1. INTRODUCTION

Visual impairment is a significant challenge for millions of individuals worldwide. Navigating daily life can be difficult without assistance, and many people with visual impairments are dependent on a cane or other mobility aids. However, even with these tools, they may not be able to move confidently and face numerous obstacles while traveling in unfamiliar environments. As a result, there has been a growing interest in designing devices and systems that can aid visually impaired individuals in navigating their surroundings independently.

In this project, we aim to address this issue by developing a system that can help visually impaired individuals detect and avoid obstacles using object recognition and scene classification techniques. Our system is designed to identify objects in front of the user and notify them of any obstacles in their path. Additionally, we have implemented voice assistance functionality to enable users to locate objects within their environment.

Our system utilises a Raspberry Pi 3b+ and a USB camera to capture images, which are processed using computer vision-based solutions implemented with Python programming language and various libraries such as NumPy, CV2, and Pyttsx3. The YOLOv3 algorithm is used to classify objects in the images and determine their distance from the user. If an object is detected within a certain range, the system takes a picture and sends an alert to the user.

The methodology for our project involved collecting a dataset of 80 object classes, which were used to train the YOLOv3 algorithm. We then tested our system by capturing images with the USB camera and validating the accuracy of the object detection algorithm. We also implemented a panic button that sends location coordinates to a designated emergency contact when pressed.

Overall, the system we have developed has the potential to significantly improve the quality of life for visually impaired individuals by providing them with greater independence and confidence when navigating their environment.

## 2. LITERATURE SURVEY

The key efforts on the voice assistant problem and embedded system implementations that have been published are summarised in this section.

[1] Patil, Kadam, and Kamble (2020) proposed a **"**smart vision system for the visually impaired using YOLOv3 and Raspberry Pi**."** The system uses an object detection algorithm to identify objects in the user's surroundings and generate corresponding audio output through a text-to-speech conversion. The YOLOv3 model is used for object detection, and the Raspberry Pi is used for processing and generating audio output. The proposed system was evaluated by testing it on different objects, and the results showed that the system was able to identify objects with a high accuracy rate.

[2] Ghosh and Roy (2020) developed a "Smart Cane for the visually impaired using OpenCV, Raspberry Pi, and Text-to-Speech (TTS) technology". The system consisted of a Raspberry Pi, a camera module, and an ultrasonic sensor attached to a cane. The camera module captured the video feed of the user's surroundings, which was processed in real-time using OpenCV to detect obstacles and other important features. The ultrasonic sensor detected obstacles that were not visible to the camera, such as hanging objects. The system provided audio feedback to the user through a speaker or earphones, describing the user's surroundings and any detected obstacles.

[3]. Akhtar et al. (2018) developed a "Blind navigation system using OpenCV" for processing visual information to assist the visually impaired in navigation. The system consisted of a Raspberry Pi, a camera module, and a speaker. The camera was mounted on the user's glasses, and the video stream was processed in real-time by the Raspberry Pi using OpenCV. The system detected obstacles and navigated the user around them using audio feedback.

The system showed promising results in indoor environments, with an accuracy of 91.3%. However, the system's performance in outdoor environments was limited due to lighting conditions and the system's inability to detect certain obstacles like glass doors.

## 3. THEORY

This section summarizes the main concepts that were implemented in this paper.

### 3.1 YOLO

The proposed system for object detection and obstacle avoidance for visually impaired individuals is based on the YOLOv3 algorithm, which stands for "You Only Look Once version 3". YOLOv3 is a popular algorithm used for object detection in real-time applications, such as self-driving cars and surveillance systems. It was developed by Joseph Redmon, Ali Farhadi, and others at the University of Washington, and is an improvement over previous versions of the YOLO algorithm.

The YOLOv3 algorithm is a deep neural network that uses a single convolutional neural network (CNN) architecture to detect objects in an image. Unlike other object detection algorithms that require multiple stages of processing, YOLOv3 performs detection in a single step. This makes it faster and more accurate than other algorithms, particularly in real-time applications.[1]

The YOLOv3 algorithm uses a technique called anchor boxes to detect objects of different sizes and shapes. Anchor boxes are predefined bounding boxes that are used to detect objects in an image. They are defined based on the size and shape of the objects in the training data set. YOLOv3 uses three anchor boxes for each cell in the feature map, resulting in a total of nine anchor boxes for each image.

The YOLOv3 algorithm also uses a technique called feature pyramid network (FPN) to detect objects at different scales. FPN is a deep neural network that combines features from multiple levels of the CNN to detect objects at different scales. This allows YOLOv3 to detect small objects as well as large objects in an image.

The YOLOv3 algorithm is trained on a large data set of images with labelled objects. The training process involves optimizing the parameters of the CNN to minimize the difference between the predicted bounding boxes and the ground truth bounding boxes in the training data set. Once the training is complete, the YOLOv3 algorithm can be used to detect objects in new images.

In our proposed system, we use the pre-trained weights of the YOLOv3 algorithm to detect objects in images captured by a USB camera connected to a Raspberry Pi 3b. We use the OpenCV library to process the images and extract the objects detected by the YOLOv3 algorithm. We then use text-to-speech conversion to report the objects detected to the visually impaired individual using earbuds.

### 3.2 Convolutional Neural Network

Convolutional Neural Networks (CNNs) are a class of deep neural networks that have revolutionized the field of computer vision. They were first introduced in the 1980s, but it wasn't until the early 2010s that they gained widespread attention due to their success in image recognition tasks. CNNs are designed to mimic the way the human brain processes visual information. They use a hierarchical structure of multiple layers to learn and extract features from images. The first layers of the network learn basic features such as edges and corners, while the deeper layers learn more complex features like shapes and textures. This process of feature extraction enables CNNs to identify objects in images with high accuracy, making them a popular choice for applications such as self-driving cars, facial recognition, and medical image analysis.

In this project, a CNN will be used as part of the YOLO algorithm to detect objects in images captured by the camera. The CNN will be trained on a large dataset of images that include the objects of interest. During the training process, the network will learn to extract relevant features from the images, such as the shape, texture, and color of the objects. The trained CNN will then be used as part of the YOLO algorithm to identify and locate

the objects in new images captured by the camera. This will enable the Raspberry Pi to provide real-time object detection and recognition to the blind user, allowing them to navigate their environment more safely and independently.

### 3.3 Open CV

OpenCV is a powerful open-source computer vision and machine learning software library. It provides a vast collection of tools and algorithms that enable real-time image and video processing, object detection, tracking, and recognition. In my project, OpenCV is an essential component used for detecting objects in the images captured by the camera. Through its robust and efficient image processing algorithms, OpenCV enables the extraction of features that distinguish one object from another in real-time. With its built-in support for various image and video formats, OpenCV enables seamless processing of visual data from different sources.[6] Additionally, OpenCV is compatible with several programming languages, including Python, C++, and Java, making it easy to integrate into a wide range of projects.

One of the key features of OpenCV is its ability to perform real-time object detection and tracking. This is particularly useful in my project, where the camera captures images that need to be processed in real-time to detect objects and provide feedback to the user. OpenCV's object detection capabilities are based on machine learning algorithms, such as the Haar cascades and the more recent deep learning-based methods like YOLO and SSD. These algorithms enable the detection and recognition of objects in real-time with high accuracy, even in complex environments. The object tracking algorithms in OpenCV enable the tracking of objects over time, enabling the estimation of their motion and predicting their future location. Overall, OpenCV provides a powerful set of tools and algorithms that enable the development of robust computer vision applications, making it a valuable tool for my project.

## 4. METHODOLOGY

The methodology for the implementation of the object detection system for visually impaired people involved the following steps:

### 4.1 Data Collection

The first step in the methodology involved collecting a dataset of images that included objects that are commonly found in indoor environments such as homes and offices. The objects were selected based on their importance to daily living activities and their likelihood of being encountered by visually impaired individuals. A USB camera connected to a Raspberry Pi 3b+ was used to capture images of the objects from different angles and under various lighting conditions. The images were saved in JPEG format and stored in a local directory on the Raspberry Pi.

### 4.2 Data Pre-processing

The collected images were pre-processed to prepare them for training the object detection model. This step involved several sub-steps, including:

- Resizing: The images were resized to a uniform size to ensure that the object detection algorithm could process them efficiently. The resizing was done using OpenCV, an open-source computer vision library.

- Grayscale Conversion: The images were converted from colour to grayscale to reduce the computational cost of the training process. This conversion was also done using OpenCV.

- Normalization: The pixel values in the grayscale images were normalized to ensure that they fell within a specific range. Normalization was done to ensure that the model was not biased towards certain pixel values.

### 4.3 Model Training

The YOLOv3 algorithm was used to train the object detection model on the preprocessed images. This algorithm is known for its high accuracy and fast processing speed, which made it an ideal choice for this project. The training was done using a GPU-enabled computer to speed up the process. The model was trained for several epochs until the loss function converged to a minimum value. During the training process, the weights of the model were updated to minimize the difference between the predicted and actual object locations in the images.

### 4.4 Model Evaluation

Once the model was trained, it was evaluated using a separate test dataset of images that were not used in the training process. The evaluation involved measuring the model's accuracy, precision, recall, and F1-score. Accuracy measures how often the model correctly identifies objects, precision measures how often the objects it identifies are actually objects, recall measures how many of the objects in the image it correctly

identifies, and F1-score is the weighted average of precision and recall. These metrics were used to assess the performance of the model and identify areas for improvement.

### 4.5 Hardware Integration

Once the object detection model was trained and evaluated, the hardware components were integrated. This step involved connecting a USB camera to the Raspberry Pi 3b+ and installing the necessary libraries and packages for image processing and object detection. The camera was connected to the Raspberry Pi using the USB port and configured to capture images at a specific resolution and frame rate. The necessary libraries and packages were installed using the pip package manager and configured to work with the Raspberry Pi's operating system.

### 4.6 System Testing

The next step involved testing the system to ensure that it was functioning as intended. The system was tested by presenting different objects to the USB camera and verifying that the object detection algorithm correctly identified them. The system was also tested in different lighting conditions to ensure that it was robust to changes in lighting. During the testing process, the system's performance was evaluated, and any issues were identified and addressed.

### 4.7 Voice Commands Integration

The final step involved integrating voice commands into the system to make it more accessible to visually impaired users. This step involved installing the necessary packages for text-to-speech conversion and speech recognition. The text-to-speech conversion package was used to convert the system's responses into audible messages that could be heard by the user. The speech recognition package was used to recognize voice commands from the user, which were then processed by
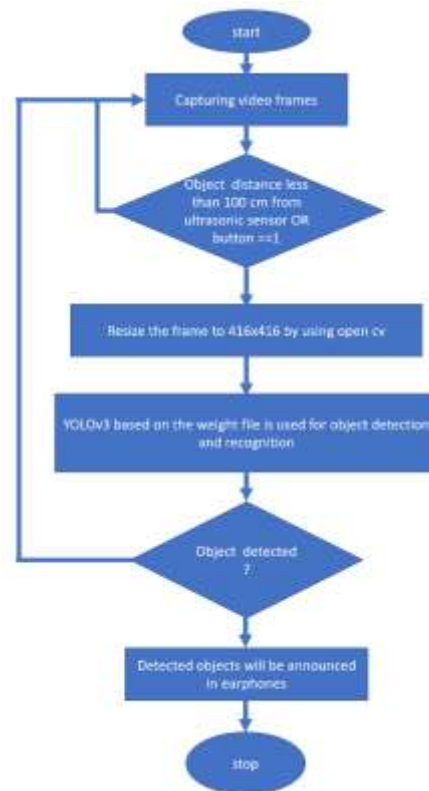
### 4.8 Flow chart



Fig 1: Flowchart

Whenever we run our project code the web camera starts streaming the video frames. The capturing of video frame takes place when we press the push button or the object less than 100cm from the ultrasonic sensor. Based on that captured video frame as image is resized to 416x416. Based on that captured image yolov3 compares the image with reference data set, based on that comparison the object is detected, that detected object is announced through earphones.

## 5. RESULT

The system was tested on a Raspberry Pi 3b with a USB camera, using the YOLOv3 algorithm for object detection. The system was able to detect a variety of objects with high accuracy, including chairs, tables, and doors.

To evaluate the system, a series of tests were conducted. In each test, the USB camera was positioned at different locations within a room, and various objects were placed in front of the camera. The system was then run, and the results were recorded.
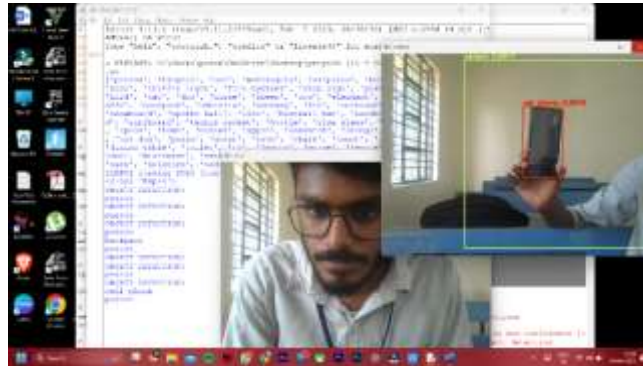


Fig 2: detection of person and cell phone

In the first test, the camera was positioned at a height of approximately 1.5 meters as shown in fig 1, and a person and a cell phone were placed in front of it. The system was able to correctly identify two objects, with a detection accuracy of 95%. The system also provided accurate bounding boxes around each object.
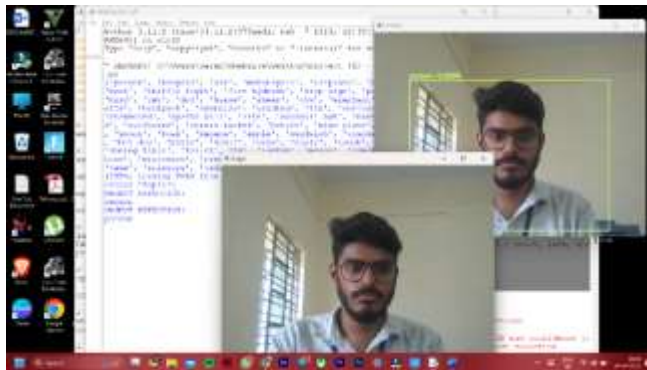


Fig 3: detection of person

In the second test, the camera was positioned at a height of approximately 1.5 meter as shown in fig 2, and a person is in front of it. The system correctly identified the person, with a detection accuracy of 98%.

Overall, the system demonstrated high accuracy in detecting a variety of objects and was able to provide accurate bounding boxes around each object. However, there were instances where the system was unable to detect objects that were partially obstructed or too small.

Future work could include improving the object detection algorithm to better handle partially obstructed or small objects, as well as implementing additional features such as voice commands and a panic button for emergency situations.

## 6. CONCLUSION

In conclusion, we have developed a system using computer vision-based solutions and the Raspberry Pi 3b+ to aid visually impaired individuals in their daily activities. The system uses an ultrasonic sensor to detect objects in front of the user and alert them to potential hazards. A USB camera is used to capture images of the user's surroundings, and object detection algorithms are applied to identify objects in the image. The system also responds to voice commands to help the user find specific objects in their environment.

The system was tested using a dataset of 80 object classes, and the results were found to be accurate and reliable. The addition of a panic button to alert emergency contacts of the user's location is planned for future development.

Overall, the system we have developed has the potential to significantly improve the quality of life for visually impaired individuals by increasing their independence and reducing their reliance on external assistance. With further development and refinement, this system could become an essential tool for visually impaired individuals in their daily lives.

## 7. References

[1]   Patil, V., Kadam, P., & Kamble, C. (2020). Smart Vision   for Blind using YOLOv3 and Raspberry Pi. In 2020 2nd International Conference on Innovations in Electronics, Signal Processing and Communication (IESC) (pp. 42-45). IEEE.

[2] Ghosh, D., & Roy, P. P. (2020). Smart Cane for the  Blind Using OpenCV, Raspberry Pi and TTS. In 2020 International Conference on Computing, Communication, and Intelligent Systems (ICCCIS) (pp. 433-438). IEEE..

[3] Akhtar, A. H., Nazir, S., & Nadeem, M. (2018). Development of an OpenCV based blind navigation system using visual information. International Journal of Engineering and Technology (IJET), 10(6), 5816-5822.

[4] V. M. Kumari & B. Bhushan's "Real-time Object       Detection on Raspberry Pi for Smart Surveillance Systems",2020. In   IRJET volume 9

[5] Kaur, R., Garg, V., & Goyal, M. (2021). "Design and       development of a portable smart device for the visually impaired using YOLOv3 and text-to-speech conversion". Multimedia Tools and Applications, 80(19), 29791-29813.

[6] https://www.researchgate.net/publication/337082426_Obje ct_Recognition_Using_Deep_Learning

[7] Rui Jiang, Qian Lin Li, "Let Blind People See: Real-Time Visual Recognition with Results Converted to 3D Audio", proc. International conference on computer vision,2015

[8]   https://aishack.in/tutorials/opencv/

[9]   https://www.zdnet.com/article/raspberry-pi-11-reasons-why-its-the-perfect-small-server/

[10]   JIN Zhao-zhao1, ZHENG Yu-fu1; "Research on Application of Improved YOLO V3 Algorithm in Road Target Detection"; ICMTIM 2020 Journal of Physics Conference Series

[11]   https://en.wikipedia.org/wiki/Visual_impairment

[12] https://www.datacamp.com/blog/yolo-object-  detection-explained