

Assistive Object Recognition System for Visually Impaired

Shifa Shaikh

Electronics and Tele-communication
Vivekanand Education Society Institute of Technology
Mumbai, India

Vrushali Karale

Electronics and Tele-communication
Vivekanand Education Society Institute of Technology
Mumbai, India

Gaurav Tawde

Electronics and Tele-communication
Vivekanand Education Society Institute of Technology
Mumbai, India

Abstract— The issue of visual impairment or blindness is faced worldwide. According to statistics of the World Health Organization (WHO), globally, at least 2.2 billion people have a vision impairment or blindness, of whom at least 1 billion are blind. In terms of regional differences, the prevalence of vision impairment in low- and middle-income regions is four times higher than in high-income regions.[6] Blind people generally have to rely on white canes, guide dogs, screen-reading software, magnifiers, and glasses to assist them for mobility, however, To help the blind people the visual world has to be transformed into the audio world with the potential to inform them about objects as well as their spatial locations. Therefore, we propose to aid the visually impaired by introducing a system that is most feasible, compact, and cost-effective. So, we implied a system that makes use of Raspberry Pi in which you only look once (YOLO v3) machine learning algorithm trained on the coco database is applied. The experimental result shows YOLO v3 achieves state-of-the-art results of 85% to 95% on overall performance, 100% (person, chair, clock, and cell-phone) recognition accuracy. This system not only provides mobility to the visually impaired with that it provides the term that ahead is an XYZ object rather than a sense of obstacle.

Keywords— Visual Impairment, Raspberry Pi, YOLO v3 Algorithm, Computer Vision, Object Recognition, voice output.

I. INTRODUCTION

“ONLY BECAUSE ONE LACKS THE USE OF THEIR EYES DOES NOT MEAN THAT ONE LACKS VISION.”

Eyesight is one of the essential human senses, and it plays a significant role in human perception about the surrounding environment. For visually impaired people to be able to provide, experience their vision, imagination mobility is necessary. The International Classification of Diseases 11 (2018) classifies vision impairment into two groups, distance and near presenting vision impairment.[6] Globally, the leading causes of vision impairment are uncorrected refractive errors, cataract, age-related macular degeneration, glaucoma, diabetic retinopathy, corneal opacity, trachoma, and eye injuries. It limits visually impaired ability to navigate, perform everyday tasks, and affect their quality of life and ability to interact with the surrounding world upon unaided. With the advancement in technologies, diverse solutions have been introduced such, as the Eye- ring project, the text recognition

system, the hand gesture, and face recognition system, etc. However, these solutions have disadvantages such as heavyweight, expensive, less robustness, low acceptance, etc. [2] hence, advanced techniques must evolve to help them. So, we propose a system built on the breakthrough of image processing and machine learning.

The proposed system captures real-time images, then images are pre-processed, their background and foreground are separated and then the DNN module with a pre-trained YOLO model is applied resulting in feature extraction. The extracted features are matched with known object features to identify the objects. Once the object is successfully recognized, the object name is stated as voice output with the help of text-to-speech conversion.

The key contributions of the paper include:

- Robust and efficient object detection and recognition for visually impaired people to independently access familiar and unfamiliar environments and avoid dangers.
- Offline text-to-speech conversion and speech output.

II. RELATED WORK

1) Real-Time Objects Recognition Approach for Assisting Blind People:

In this paper, two cameras placed on blind person's glasses, GPS free service, and ultrasonic sensors are employed to provide information about the surrounding environment. Object detection is used to find objects in the real world such as faces, bicycles, chairs, doors, or tables that are common in the scenes of a blind. Here, GPS service is used to create groups of objects based on their locations, and the sensor detects an obstacle at a medium to long distance. The descriptor of the Speeded-Up Robust Features (SURF) method is optimized to perform the recognition. The use of two cameras on glasses can be sophisticated. [2]

2) Wearable Object Detection System for the Blind:

In this paper, the RFID device is designed as a support for the blind for the disclosure of objects; especially, it is developed for searching the medicines in a cabinet at home. This device can provide information about the distance of a defined object, how near or far it is and simplifies the search. For identifying the medicines, the device can provide the user with an acoustic signal to find the desired product as soon as possible. The

The diagram illustrates the system architecture for a video game. It features a central Raspberry Pi board. To the left, a cartoon character holding a camera is connected to the Pi via a USB cable, with the label "Real-time Images from Camera". Above the Pi, a black power bank is connected via a USB cable, with the label "Power supply from Power bank". Below the Pi, a USB cable connects to a pair of headphones, with the label "Audio output from Headphones/Speaker".

Diagram illustrating the VGG-16 architecture, showing the sequence of convolutional and pooling layers. The input is a 224x224x3 image.

The layers are defined as follows:

- Conv. Layer: $7 \times 7 \times 64 + 2$
- Conv. Layer: $3 \times 3 \times 128$
- Maxpool Layer: $2 \times 2 + 2$
- Conv. Layer: $3 \times 3 \times 256$
- Conv. Layer: $3 \times 3 \times 256$
- Conv. Layer: $3 \times 3 \times 512$
- Maxpool Layer: $2 \times 2 + 2$
- Conv. Layer: $3 \times 3 \times 1024$
- Conv. Layer: $3 \times 3 \times 1024$
- Conv. Layer: $3 \times 3 \times 1024 + 2$
- Conv. Layer: $3 \times 3 \times 1024$
- Conv. Layer: $3 \times 3 \times 1024$
- Conn. Layer: 4096
- Conn. Layer: 1000

737

$$\begin{aligned} & \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\ & + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[(\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2 \right] \\ & + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} (C_i - \hat{C}_i)^2 \\ & + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} (C_i - \hat{C}_i)^2 \\ & + \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \end{aligned}$$

Fig 4: Loss function expression

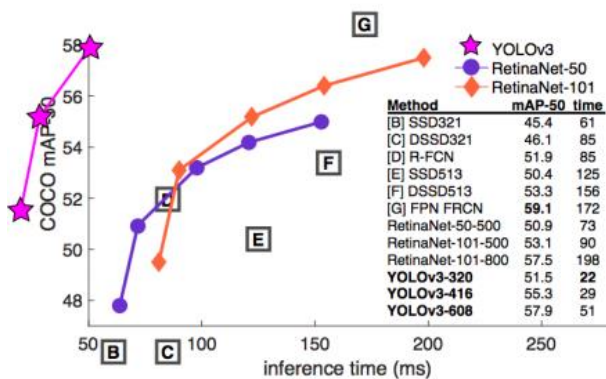


Fig 5: Comparison of YOLO with other algorithms.

Opencv (Open source computer vision): is a library of programming functions mainly aimed at real-time computer vision. The library has more than 2500 optimized algorithms. [12] These algorithms can be used to detect and recognize faces, identify objects, classify human actions in videos, track camera movements, etc.

Dnn module (Deep Neural Network): dnn is the module in OpenCV, which is responsible for all deep learning related concepts.

Type	Filters	Size	Output
1x	Convolutional	32 3 x 3	256 x 256
	Convolutional	64 3 x 3 / 2	128 x 128
	Convolutional	32 1 x 1	
	Convolutional	64 3 x 3	
2x	Residual		128 x 128
	Convolutional	128 3 x 3 / 2	64 x 64
	Convolutional	64 1 x 1	
	Residual		64 x 64
8x	Convolutional	256 3 x 3 / 2	32 x 32
	Convolutional	128 1 x 1	
	Convolutional	256 3 x 3	
	Residual		32 x 32
8x	Convolutional	512 3 x 3 / 2	16 x 16
	Convolutional	256 1 x 1	
	Convolutional	512 3 x 3	
	Residual		16 x 16
4x	Convolutional	1024 3 x 3 / 2	8 x 8
	Convolutional	512 1 x 1	
	Convolutional	1024 3 x 3	
	Residual		8 x 8
Avgpool		Global	
Connected		1000	
Softmax			

Fig 6: DARKNET-53 Backbone Architecture

In our system, images are pre-processed, background and foreground are separated, and then the DNN module is applied, resulting in feature extraction. Darknet-53 is used as the feature extractor, as it achieves the classification accuracy 2x

faster. And the extracted features are matched with known object features to recognize them. A DNN-based algorithm is more robust and accurate on a wide range of faces. DNN module only allows forward propagation on the pre-trained model.[13] In DNN module of OpenCV, it requires your input transform to a blob, or tensor in other neural network frameworks.

Pytsx3 lib: For audio output, we use the “pytsx3” library. The main advantage of using this library is that it does not require any internet, and it converts text-to-speech very fast and effectively.

V. WORKING & IMPLEMENTATION

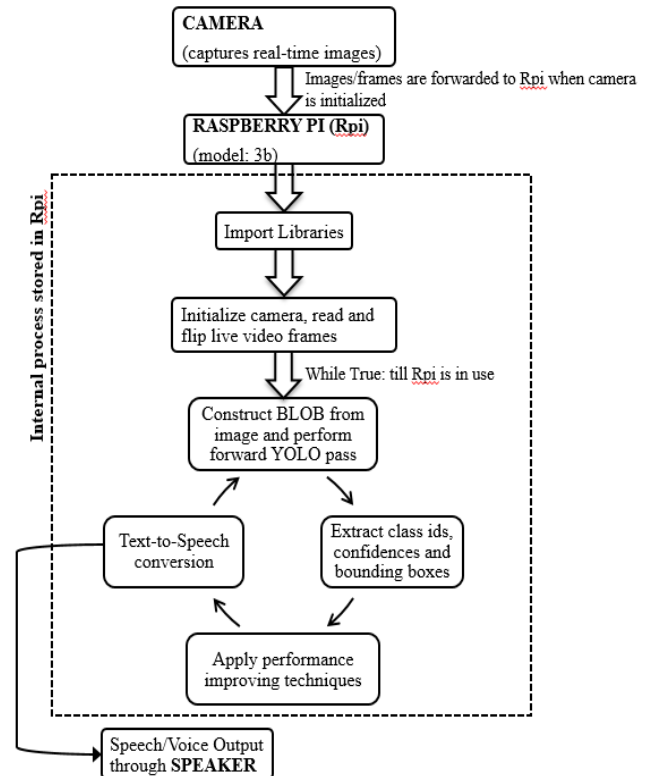


Fig 7: Flowchart of system

The Flowchart is a communication outline that shows how objects work with each other and in what order. The above Flowchart of our framework clarifies the stream that is first, the user starts and wears the system. Once the Raspberry Pi (Rpi) is on, it will implement its internal process/code. The code keeps on executing till the Rpi is on. Initially, Rpi will import all the libraries that are: OpenCV, Pytsx3, Time, and NumPy and will read the text file containing class names, YOLO weights, and configuration files. After that, the code will initialize the camera connected to it. The camera will capture real-time frames at 1fps (frame per second), then the code will read the input image/frame and get its width and height to an adequate level. Then an object detection algorithm in our case YOLO is applied to this altered frame. Before forward passing this altered image to YOLO weights and YOLO configuration files, a 'BLOB from image' is constructed. To obtain (correct) predictions from deep neural networks such as YOLO, you first need to pre-process your data. In the context of deep learning, feature extraction, and

image classification, we have used the OpenCV function blobFromImage. This function performs the following:

1. Mean subtraction - is used to help combat illumination changes in the input images in our dataset.
2. Scaling by some factor - is used to scale the input image space into a particular range
3. Optionally channel swapping. [8]

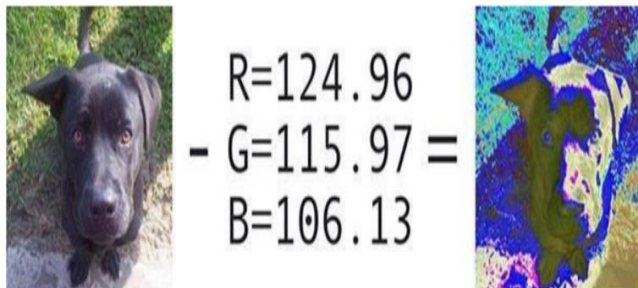


Fig 8: A visual representation of mean subtraction where the RGB mean (centre) has been calculated from a dataset of images and subtracted from the original image (left) resulting in the output image (right).

Then the code performs a forward pass of the YOLO object detector, giving us our bounding boxes, class ids, and associated class probabilities.

Another advantage of YOLO other than being fast is that it provides three methods to improve its performance:

- **Intersection over Union (IoU)** decides which predicted box is giving a good outcome. It calculates the IoU of the actual bounding box and the predicted bounding box.
- **Non-max suppression** suppresses weak, overlapping bounding boxes.
- **Anchor Boxes** detects multiple objects in a single grid.[7]

Further, the frames are divided into a 3x3 grid, which helps in finding the position of objects. Our system aims to produce an audio output for the visually impaired. The Detected object labels are converted into speech using the pyttsx3 library.

Lastly, Upon successful recognition of an object and as per grids, the system will provide speech output stating the name of the object along with its grid name, for e.x. 'Mid left car', 'Mid right car'. Hence helping the visually impaired people in recognizing the objects in the field of view.

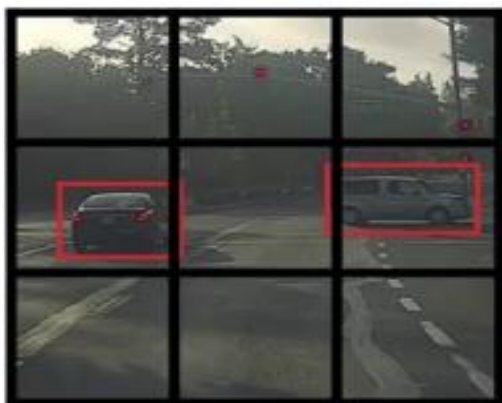


Fig 9: Division of Image into grids

VI. RESULTS AND DISCUSSION

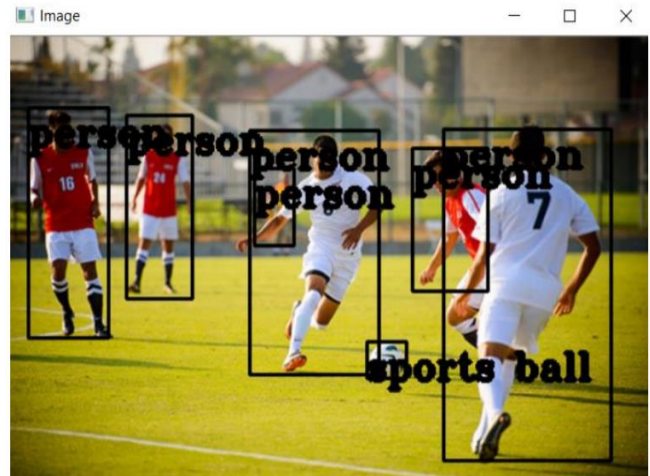


Fig 10: Image with bounding boxes and class label

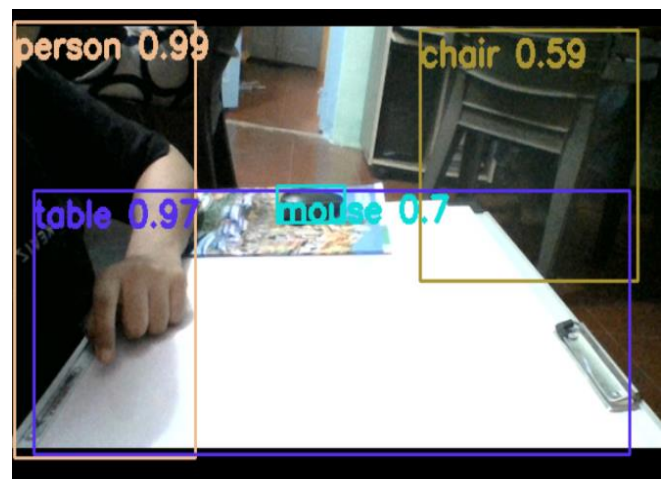


Fig 11: Real-time Object Detection with Multiple Bounding Boxes 1

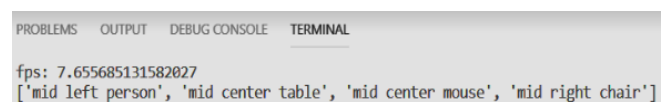


Fig 12: Real time YOLO Object detection with text output 1

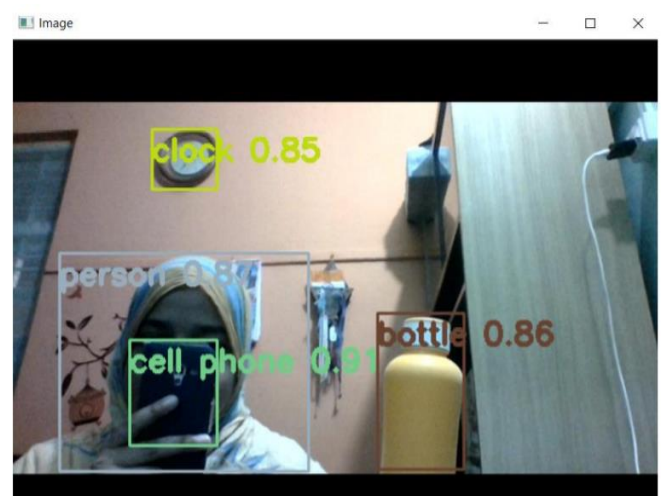


Fig 13: Real-time Object Detection with Multiple Bounding Boxes 2

```
fps: 7.733104841379943  
['bottom left cell phone', 'mid left person', 'bottom center bottle', 'top left clock']
```

Fig 14: Real time YOLO Object detection with text output 2

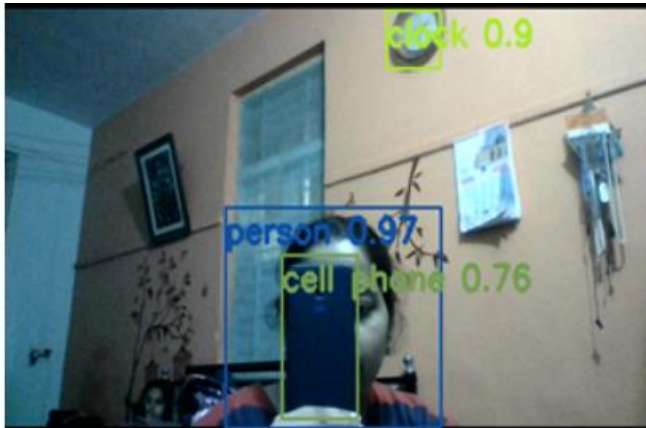


Fig 15: Real-time Object Detection with Multiple Bounding Boxes 3

```
fps: 8.380539030786958  
['bottom center person', 'top center clock', 'bottom center cell phone']
```

Fig 16: Real time YOLO Object detection with text output 3

Fig 10 illustrates object detection and recognition of the already acquired image. The system has successfully recognized every object present in the image based on the trained coco dataset. Similarly, Fig 11, Fig 13, and Fig 15 illustrate real-time object detection and recognition, along with their confidences of class recognition. Also, Fig 12, Fig 14, and Fig 16 illustrate the text form, which is converted into speech. We have successfully achieved a speed of 7 fps to 9 fps in this CPU based system. The speed of detection and recognition can be increased with the use of a GPU based system.

VII. FUTURE SCOPE

The future perspective of this project is to increase the object recognition rate which can be achieved by using the TensorFlow library and to provide an exact distance measurement between the people and object. However, for developing an application that involves many objects that are fast-moving, you should instead consider faster hardware. Further, we can implement face recognition and text recognition in the same system. Thus, making the system compatible overall.

VIII. CONCLUSION

In recent years, some solutions have been devised to help blind or visually impaired in recognizing objects in their environment but they are not efficient. Our purpose is to provide a robust and comfortable system for the blind to recognize their surrounding objects. Our advanced system uses a USB camera to seize real-time images in front of the users. The machine learning and feature extraction technique used here is YOLO. The YOLO framework trades with object detection by choosing the entire image in a single instance, and splits the image into grids, then predicts the bounding box coordinates and class probabilities for these boxes. The biggest advantage of sing YOLO is its excellent speed – it's incredibly

fast and YOLO also understands generalized object representation. This system will make visually impaired virtually visible also it innovatively uses the text-to-speech technology which provides audio descriptions of their surroundings and helps them to travel with self-confidence. The proposed system is mobile, robust, and efficient. Also, it creates a virtual environment and this system provides a sense of assurance as it voices the name of the object recognized.

REFERENCES

Book,

- [1] Peter Harrington, "Machine Learning in Action, pp."by Manning Publications-1st edition.

Papers,

- [2] Zraqou, Jamal & Alkhadour, Wissam & Siam, Mohammad. (2017). "Real-Time Objects Recognition Approach for Assisting Blind People." Multimedia Systems Department, Electrical Engineering Department, Isra University, Amman-Jordan Accepted 30 Jan 2017, Available online 31 Jan 2017, Vol.7, No.1
- [3] A. Dionisi, E. Sardini and M. Serpelloni, "Wearable object detection system for the blind," 2012 IEEE International Instrumentation and Measurement Technology Conference Proceedings, Graz, 2012, pp. 1255-1258, doi: 10.1109/I2MTC.2012.6229180.
- [4] Daniyal, Daniyal & Ahmed, Faheem & Ahmed, Habib & Shaikh, Engr & Shamshad, Aamir. (2014). "Smart Obstacle Detector for Blind Person." Journal of Biomedical Engineering and Medical Imaging. 1. 31-40. 10.14738/jbemi.13.245.
- [5] Christian Szegedy Alexander Toshev Dumitru Erhan, "Deep Neural Networks for Object Detection."
- [6] N.Saranya, M.Nandinipriya, U.Priya,"Real Time Object Detection for Blind People",Bannari Amman Institute of Technology, Sathyamangalam, Erode.(India).
- [7] Rui (Forest) Jiang,Qian Lin,Shuhui Qu,"Let Blind People See: Real-Time Visual Recognition with Results Converted to 3D Audio",Stanford,2018.

Website,

- [8] <https://www.who.int/news-room/detail/08-10-2019-who-launches-first-world-report-on-vision>
- [9] <https://www.analyticsvidhya.com/blog/2018/12/practical-guide-object-detection-yolo-framework-python/>
- [10] <http://datahacker.rs/face-detection-opencv-images/>
- [11] <https://towardsdatascience.com/yolo-v3-object-detection-53fb7d3bfe6b>
- [12] https://medium.com/@jonathan_hui/real-time-object-detection-with-yolo-yolov2-28b1b93e2088
- [13] <https://opencv.org/about/>
- [14] https://github.com/TheNsBhasin/DNN_Object_Detection