

# Object Recognition App for Visually Impaired

Sumitra A. Jakhete  
Prof. IT dept

Pune Institute of Computer Technology  
Pune  
sajakhete@pict.edu

Pranit Bagmar  
Student, IT dept

Pune Institute of Computer Technology  
Pune  
prbagmar78@gmail.com

Avanti Dorle  
Student, IT dept

Pune Institute of Computer Technology  
Pune  
avantidorle98@gmail.com

Atharva Rajurkar  
Student, IT dept

Pune Institute of Computer Technology  
Pune  
atharvarajurkar97@gmail.com

Piyush Pimplikar  
Student, IT dept

Pune Institute of Computer Technology  
Pune  
piyushpimplikar049@gmail.com

**Abstract**—Vision is one of the most important senses that help people interact with the real world. There are nearly 200 million blind people all over the world, and being visually impaired hinders a lot of day to day activities. Thus it is very necessary for blind people to understand their surroundings, and to know what objects they interact with. This project proposes an android application to help blind people see through handheld device like mobile phone. It integrates various techniques to build a rich android application that will not only recognize objects around visually impaired people in real time but also give an audio output to assist them as quickly as possible. SSD (Single Shot Detector) Algorithm is used for the object recognition as well as detection. Also this algorithm gives nearly accurate results for real time object detection and is proven to be faster than other relative algorithms. The application further uses android tensorflow APIs and android TextToSpeech API to give audio output.

**Keywords**—Blind, Object, Detection, SSD, Tensorflow, TextToSpeech

## I. INTRODUCTION

It is easier for vision enabled people to carry out their everyday activities since they can clearly see all the objects in their surroundings, any obstacles they come across, other people and hence is easy to interact with these objects. Whereas, visually impaired people have to struggle a lot to deal with real world due to their everyday chores and jobs. There are more than a million blind people all over the world, and being visually impaired hinders a lot of day to day activities. Thus it is very essential for blind people to know their surroundings, and what objects they interact with to prevent accidents and make their life simpler. Visually impaired people always need someone to guide them throughout their day such as to cross roads, catch a bus, and many such activities. The main motive of developing this application is to assist blind people.

This application aims to help the visually impaired people to know their surrounding objects that could be just basic everyday objects or can create obstacle in their activity. The application is build to recognise or detect some household objects like chair, table, bed, refrigerator, laptops etc and some on outdoor objects like cars, motorbikes, potted plant, people etc.

The application will use mobile phone camera to scan the surrounding in real time and take the frames from the ongoing video. The frames will be sent to the next module where the SSD algorithm will create bounding boxes around the objects in the frame and classify them into given categories.

At last the application will produce an audio output of the object detected which has the maximum confidence score among all other present in the frame. The frames are selected at a particular time interval to avoid the hindrance in the audio output.

## II. LITERATURE SURVEY

The thesis deals with the use of a mobile camera as an assistive tool in a special phone for the visually impaired people, based on the Android system and developed in the Department of Cybernetics located in the Czech Technical University, Prague. In the introduction, it discusses the applicability of general tasks of computer vision in the area of assistive technologies for the blind and visually impaired. It deals with both the specific use of the camera (for identifying objects, banknotes, reading text labels), as well as accessing the basic functionality of the camera itself. A novel banknote recognition algorithm based on the BRISK descriptor and the Gradient Boosted Trees Classifier is implemented as a part of this work. Implementation of an accessible mobile user interface to Google's existing real-time text recognition library, as well as the implementation of a simple camera and image gallery application is also discussed[1].

According to [2], we read a strategy for identifying objects as well as detecting them in pictures using only one DNN (Deep Neural Network). This methodology, named SSD, represents the bounding box outputs in discrete form and modifies into many default boxes based on different aspect ratios as well as scales for each map of a feature. During the time of prediction, the neural network creates score values if there exists each category of object in every single default box, and also delivers changes in accordance to the box that will match the shape inference of an object. There exists 74% of mAPI for an input of 300 by 300, on the dataset VOC2007 at 59 Frames per second on an Nvidia

Titan X and for an input of 512 by 512, 76% mAP, is achieved by Single Shot Detector outperforming Compared to other models, Single Shot Detector has the best accuracy even when the input size of image is smaller.

### III. ALGORITHM

#### SSD (Single Shot Detector) :

There exists 2 sections for an object detection system using SSD algorithm:

- One is to extract the feature maps
- Second is to apply filters of convolution to detect objects.

SSD uses VGG16 (a CNN Architecture Used in SSD. It consists of 16 Layers) to extract the feature maps. Then Conv4\_3 layer is used to detect objects. For each cell (also called location), it makes 4 predictions of the objects.



Fig. 1. 4 Predictions at each cell

There exists a boundary box for each class and also 21 scores for the same for each prediction. Also there is one extra class for no object. Any class having a bounded object with highest score gets chosen. A total of 38 by 38 by 4 predictions are created by Conv4\_3. Where regardless of the feature maps depth, there exists 4 predictions for each cell where most predictions have no object. Class "0" is reserved by SSD to indicate that it has no objects.

The location and the class scores are computed by using small filters of convolution. Once the feature maps extraction is done, a 3 by 3 convolution filter is applied by SSD for every cell in order to predict the result. Each of these filters gives an output in the form of 25 channels. These channels include 21 scores one for every class as well as one for the boundary box.

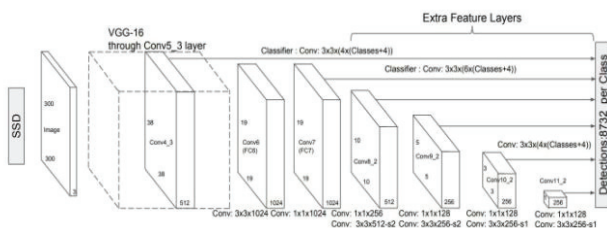


Fig. 2. SSD Architecture

Multi-scale feature maps is used for independent object detection. The feature maps resolution gets decreased since spatial dimension is reduced gradually by the CNN. To detect objects with large scale, SSD uses lower resolution layers.

#### CNN (Convolutional Neural Network):

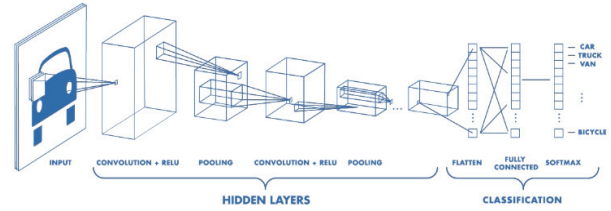


Fig. 3. CNN Architecture

Convolutional Neural Network (CNN) is a neural network with multiple layers. It is a class of Deep Neural Network, that is mostly used for the purpose of analysis of visual images. CNN contains one layer for input, called as input layer, next it has one or more middle layers which are called as convolutional layers, next the network consists of one or more fully connected layers and lastly the network is completed by adding an output layer in the respective order. We are using CNN as a building block for creating our models. It works as follows, a convoluted process is applied to the incoming data by the convolutional layers after which the outcome that is generated is passed to the next layer. This convolution layer has reaction similar to the human's neuron to vision process. Every neuron in the convolution uses information that it receives and processes the information that it is responsible for. Although there exists feedforward neural networks that can classify the data based on the features, it is not feasible to use such systems when it comes to processing pictures. Convolutional systems might also include pooling layers with the scope of global or local. Pooling layer is one of the building blocks that helps to reduce progressively the dimensionality of representation, due to which the complexity of computation is reduced to a great extent. It does so by grouping some neurons from one layer. Local pooling mostly groups neurons in small number like 2x2. Global pooling applies the process to every neuron. Layers that have complete associativity consists of every single neuron from one layer connected to every single neuron from another layer. This is sometimes also similar to the multi-layered neural network on fundamental level. The flattened grid or matrix undergoes changes through completely associated layer for image classification.

### IV. IMPLEMENTATION

#### Dataset :

Our Android application uses Neural Networks for object recognition. This requires an image dataset of the objects to train the classifier. In this project we have used COCO (Common Objects in Context) 2014 Database with 80 different object classes which have 83K training images, 41K Testing images. The dataset used is the labelled dataset which is useful to train the model. Some of the objects among 80 classes are as follows:

- Person: person
- Animal: cat, cow, dog, horse, sheep etc
- Vehicle: bicycle, boat, bus, car, motorbike, train, truck etc
- Indoor: bottle, chair, dining table, potted plant, sofa, bed, TV/monitor, laptop etc
- Other: bench, trafficlight, fire hydrant, stop sign etc.[10]

### Data Preparation:

The COCO dataset was downloaded from [cocodataset.org](http://cocodataset.org)[10].

### Data Labeling :

The images are labeled by using Labellmg software. For some images the annotations file is downloaded with the dataset itself. Annotation file contains parameters object\_class, unique object\_id, x\_coordinate for centre, y\_coordinate for centre, width and height for each image.

### Train-Test Split:

After collecting and annotating the dataset, we randomly shuffle the data to select 80% of the data on which we train the model. The remaining 20% of the data, unseen by the model, is used for the testing of the model.

### Model training:

The main idea behind making object detection or object classification model is Transfer Learning which means using an efficient pre-trained model.

Here we have using three models : Object Detection API provided by Tensorflow (uses SSD mobilenet v1), MULTIBOX and YOLO. By default Object Detection API by Tensor is used since it was found to be most efficient.

### Real Time Video Processing:

The frames are captured at the rate of - frames per second with preview size of 640 x 680.

The stable output is generated for the real-time input.

### Object Detection:

Bounding boxes are generated which predicts the certainty called as confidence score. This score lets us know that the bounding box consists of some object. For every bounding box, the cell predicts a class of that object which gives a distribution of probability among all the available classes in the given model. The confidence score along with the probability just calculated, gives us the final score which lets the user know how likely it is that the bounding box contains some specific object.

Minimum detection confidence to track a decision:

For Tensorflow Object Detection API : 0.6f

For MULTIBOX : 0.1f

For YOLO : 0.25f

Bounding boxes whose score is more than the threshold is given as an output along with the respective class name.

Finally, the generated text output is converted into audio by using TextToSpeech API.

### Android Application

The model is integrated into Android Application. The app uses a rear camera of the smartphone for real time processing.

## V. RESULTS AND EVALUATION

### Detail discussion of experiments carried out

Since there has been a lot of research and experiments in Computer Vision and particularly its domain Object Detection and Recognition, there exists many different algorithms and models that strive to get nearly accurate results. There are following models available:

- Single Shot Multibox Detector (SSD)
- RCNN
- You Only Look Once (YOLO)

Out of which we decided to work on SSD due to its benefits and high processing on real time data.

Firstly, we gathered the COCO 2014 data for various objects including 80 classes. At first we decided to use trained YOLO model. We considered the YOLO implementation to get a detailed insight of how object detection gets carried out in various models. We also tried SSD implementation of which we achieved better accuracy and performance than YOLO. Hence SSD was chosen as the primary model of project. Next we were successful in developing a basic prototype in python language and which could detect and recognize objects using a webcam, and provide voice output as well. Further this module was integrated in Android studio, and apk was generated.



Fig. 4. Output: Chair and person detection



Fig. 5. Output potted plant and dining table



## VI. CONCLUSION

Visually impaired people today can read using Braille script but it is still tough for them to recognise and interact with household objects and also on roads. If they can recognise the object more easily or can know if car or bike is coming in their way it would be easy for them to handle them or act according to.

In this project we developed an android application to aid the visually impaired people which helps them to recognise the objects they come across and sends them an audio of the label based on the confidence score of the predicted object in the frame.

## REFERENCES

- [1] Bc. Jan Hadaček, "Application of a Camera in a Mobile Phone for Visually Impaired People." Masters thesis, Czech Technical University in Prague, May 2017.
- [2] "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks." Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun, IEEE transactions, Dec 2016
- [3] "You Only Look once : Unified, Real-Time Object Detection ." J Redmon, S Divvala, R Girshick, A Farhadi, IEEE transactions, May 2016
- [4] "SSD: Single Shot MultiBox Detector Wei Liu." , Dragomir Anguelov , Dumitru Erhan , Christian Szegedy, Scott Reed , Cheng-Yang Fu , Alexander C. Berg, IEEE transactions, Jan 2016
- [5] "Support Vector Machines for 3D Object Recognition" , Massimiliano Pontil and Alessandro Verri , IEEE transactions on pattern analysis and machine intelligence, vol:20, no. 6
- [6] Vicky Mohane, Chetan Gade " Object Recognition for Blind people Using Portable Camera" WCFTR World conference 2016
- [7] Meghajit Mazumdar, Dr. Sarasvathi V, Akshay Kumar "Object Recognition in Videos by Sequential Frame Extraction using Convolutional Neural Networks and Fully Connected Neural Networks" International Conference on Energy, Communication, Data Analytics and Soft Computing 2017
- [8] Yide Ma, Dong Hwan Kim, and Sung-Kee Park "Region-Based Object Recognition by Color Segmentation Using a Simplified PCNN" IEEE transactions on neural network and learning system, Vol, 26 No. 8 Aug 2015
- [9] Mingjie Liang, Huaqing Min, Ronghua Luo, and Jinhui Zhu "Simultaneous Recognition and Modeling for Learning 3-D Object Models From Everyday Scenes" IEEE transaction on cybernetics, Vol 45 No. 10 Oct. 2015
- [10] <http://cocodataset.org/#home>
- [11] <https://www.kaggle.com/jessicali9530/coil100>
- [12] <http://www.vision.ee.ethz.ch/en/datasets/>
- [13] <https://blog.statsbot.co/real-time-object-detection-yolo-cd348527b9b7>