# Android Travel Mate Application using Real Time OCR

**Qayyum Sayed[1], Bhagyesh Padate[2], Ashika Mistry[3], Jayshree Vanmali[4]**

B.E Scholar, Information Technology Department, Universal College of Engineering, Thane, India [1, 2, 3]

Assistant Professor, Information Technology Department, Universal College of Engineering, Thane, India [4]

**Abstract:** Proposed system that can extract a text from the image, the extracted text is translated into a specific language. This application is useful for Tourists and travelers to understand the native country language. They simply use their mobile camera, simply click the image of the signboards, menus etc. The OCR engine used in the system extracts the text from the image. Tourist and Travelers simply use their Android phone camera for clicking the image.

**Keywords:** Android, OCR, Tesseract Engine.

## I. INTRODUCTION

The people now a days use mobile phones, which are the smartphones. So the application has a great demand in the mobile market. Optical character recognition (OCR)[4] is a tool that scans image to find a text in it. It recognizes text form natural scenes, hand-written forms etc.[2]

The internet and smartphones are connecting people and let them exchange information. Tourist and Traveller can translate the text into 12 different languages.

### a. Android
Android is software for mobile devices which means a set of application programs that form a complete system. This software platform provides a foundation for applications just like a real working platform. Android is an operating system based on Linux with a Java programming interface. It provides tools, e.g. a compiler, debugger.

### b. OCR
Optical Character Recognition(OCR)[1] It is one such system that allows us to scan printed, typewritten or handwritten text convert scanned images in to a computers able format, either in the form of a plain text or a word text. This improves the accuracy of recognizing the character during document processing compared to various existing available character recognition methods.

### c. Tesseract Engine
The Binarization[7] of Captured Image takes place, after that the text layout is analysed, Blobs are detected and finally words and lines are detected. The words are sent to a number of passes. In these passes each word is chopped into characters and characters are checked for the need of joining the broken characters or the breaking of associated characters. Finally chopped characters are recognized with the help of inbuilt fuzzy features matched to language specific training data of Unicode characters. After each pass the words are matched back and forth with the Language specific Dictionary words.

## II. EXISTING SYSTEM

Existing OCR have only capability to convert and recognize only the documents of English or a specific language only. In case of handwritten characters recognition, models should be used to constrain the character choices to overcome the wide variability of hand printing and cursive script. A pattern recognition algorithm is used to extract shape features and assign the observed character into the appropriate class. It overcomes the existing problems with OCR technology i.e. limited memory and limited processing power challenge.

## III. PROPOSED SYSTEM

Proposed system enables Travelers and Tourists to easily capture the native country language Books pages, signboards, banners and hotel menus etc. The OCR [1] converts the text in the captured image into Unicode text format. It also provides translation facility so that Tourists can translate this Unicode text into their own country language. Also android platform has been increasingly being common in accordance with its features like low-cost, customizable, lightweight operating system and more.

## IV. METHODOLOGY

### a. Camera Capture Module
In this module the user is allowed to resize the camera capture box by touching the box corners on the screen so as to capture the only concerned text image from signboard, banner and book pages. Once the capture button is pressed the beep sound plays and the captured image is sent to Tesseract OCR engine module.[3]

### b. Dictionary words Matching Module
In this module each group of sequential characters is searched for a dictionary based word match, which helps in identifying the word more accurately rather than just giving an earning less word as result. Finally the
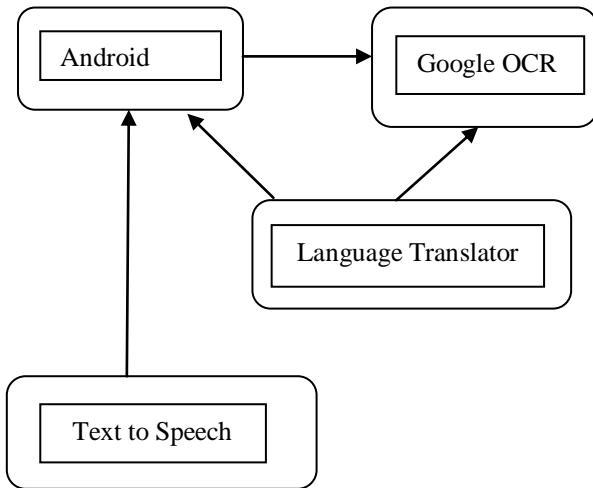
recognized text is transferred to Unicode text Post processing Module.

### c. Unicode Text Post processing Module

In this module, the recognized characters are displayed as Unicode characters and the user is allowed to translate the recognized text into his desired language available in the drop down list from settings.



**Fig 1. Internal processing**

## V.  CONCLUSION

The proposed system tells about OCR system for Signboard character recognition. The systems have the ability to serve excellent results. This application provides fast and robust high Quality performance because of having improved Auto focus, continuous dynamic preview, improved noise tolerance feature and no remote computing overhead.

## ACKNOWLEDGMENT

## REFERENCES

[1]   Smith, R. "An Overview of the Tesseract OCR" in proc.  ICDAR 2007, Curitiba, Paraná, Brazil.
[2]   Bansal, V. and Sinha, R.M.K. "A Complete OCR for Printed Hindi Text in Devanagari Script", Sixth International Conference on Document Analysis and Recognition, IEEE Publication, Seatle USA, 2001, Page(s):800-804.
[3]   Pal, U., Chaudhuri, B. B. "Indian Script Character recognition: A survey", Pattern Recognition, vol. 37, pp. 1887-1899, 2004.
[4]   Saba, T., Sulong, G. and Rehman, A. "A Survey on Methods and Strategies on Touched Characters Segmentation", International Journal of Research and Reviews in Computer Science (IJRRCS) Vol. 1, No. 2, June 2010.
[5]   Simple Android Photo Capture. Available at: http://labs. makemachine.net/2010/03/simple-android-photo-capture/
[6]   Microsoft Translator Java API. Available at: http://code.google.com/p/microsoft- translator-java-api.
[7]   Tesseract OCR training data downloads. Available at http://code.google.com/p/tesseract-ocr/download/list