



IJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 10 **Issue:** V **Month of publication:** May 2022

DOI: <https://doi.org/10.22214/ijraset.2022.42476>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Optical Character Recognition for Image & Handwriting to Text Conversion

Vasundhara Pathak¹, Shriyansh Sharma², Tanishka Goel³

^{1, 2, 3}Meerut Institute of Engineering and Technology, Meerut, India

Abstract: *This paper combines the functionality of Optical Character Recognition and speech synthesizer. The idea is to develop stoner friendly operation which performs image to text conversion.*

Objective

The advantage of proposed system that overcomes the disadvantage of the prevailing system is that it supports multiple functionalities like editing and searching. It also adds benefit by providing heterogeneous characters recognition.

I. INTRODUCTION

OCR is the acronym for Optical Character recognition. This technology allows to automatically recognizing characters through an optical mechanism. In case of human beings, our eyes are optical mechanism. Optical Character Recognition (OCR) market size is predicted to be USD 13.38 billion by 2025 with a year on year growth of 13.7 %. This growth is driven by rapid digitization of business processes using OCR to scale back their labor costs and to save lots of precious man hours.

Although OCR has been considered a solved problem there's one key component of it, Handwriting Recognition or Handwritten Text Recognition (HTR) which remains considered a challenging problem statement. The high variance in handwriting styles across people and poor quality of the handwritten text compared to printed text pose significant hurdles in converting it to computer readable text. Nevertheless it is a crucial problem to unravel for multiple industries like healthcare, insurance and banking

II. APPLICATIONS FOR OCR

A. Matlab

What's Matlab?

It's a programming language for specialized computing. It integrates calculation & visualizing, programming in an easy-to-use terrain where problems and results are expressed in familiar form memorandum.

1) Scientific and engineering plates

2) The development of the operation, including graphical stoner interface structure

MATLAB is an interactive system whose introductory data element is an array that doesn't bear dimensioning. By using this we can break colorful specialized & mathematical computing problems, especially those with matrix and vector phrasings, in a veritably lower time it would take to write a program in a scalar non-interactive languages like in similar as. C or Fortran.

It's firstly written to give easy access to matrix software developed by the LINPACK and EISPACK systems. Moment, MATLAB uses software developed by the LAPACK and ARPACK systems, which together represent the state-of-the-art

B. Tesseract

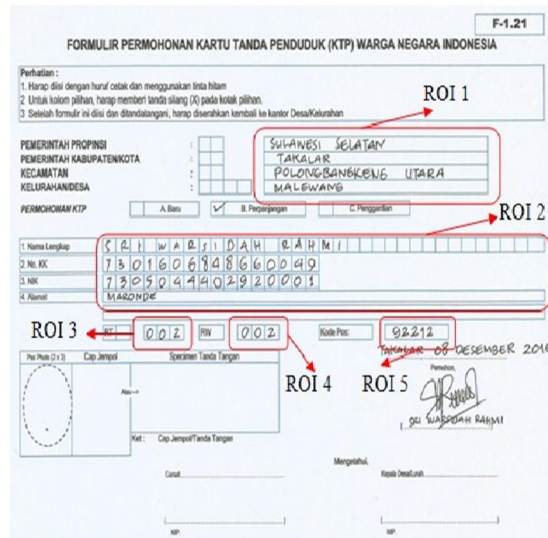
What's Tesseract ?

It's OCR machine which is open source. Tesseract has been developed in 1984 to 1994. In 1995, it was transferred to UNLV for Annual Test of OCR Accuracy after the common design between HP Labs Bristol and HP's Scanner Division in Colorado.

2005, It is released as an open source by HP and can be downloaded from-. code.google.com/p/tesseract-ocr (2).

Tesseract works with singly developed Page Layout Analysis Technology. Hence Tesseract accepts input image as a double image. Tesseract can handle both, the traditional Black on White textbook and also inverse-White on Black textbook. Silhouettes of element are stored on connected Element Analysis. Nesting of outlines is performed That Combines the outlines together like Blob. Similar Blobs are organized into textbook lines. Text lines are anatomized for fixed pitch and commensurable textbook. Also the lines were splited in words by analysis according to the character distance. Fixed pitch is diced in character cells and commensurable textbook is broken into words by definite spaces and fuzzy spaces.

Tesseract performs exertion to fete words. This recognition exertion is substantially consists of two passes. The first pass tries to fete the words. Also satisfactory word is passed to Adaptive Classifier as training data, which recognizes the textbook more directly. During alternate pass, the words which weren't honored well in first pass arehoned again through run over the runner.



III. TECHNIQUES

The original approaches of working handwriting recognition involved Machine Literacy styles like Hidden Markov Models (HMM) SVM etc. Once the initial text is pre-processed, feature extraction is performed to spot key information like loops, inflection points, ratio etc. of an individual character. These generated features are now fed to a classifier say HMM to urge the results.

The performance of machine learning models is pretty limited thanks to manual feature extraction phase and their limited capacity of learning. Feature extraction step varies for each individual language and hence isn't scalable. With the arrival of deep learning came tremendous improvements in accuracy of handwriting recognition. Let's discuss few of the prominent research within the area of deep learning for handwriting recognition

IV. MULTI-DIMENSIONAL RECURRENT NEURAL NETWORKS

RNN/LSTM as we all know can affect sequential data to spot temporal patterns and generate results. But they're limited to handling 1D data and hence won't be directly applicable to image data. To solve this problem, the authors in this paper proposed a multidimensional RNN/LSTM structure as can be seen in the figure below

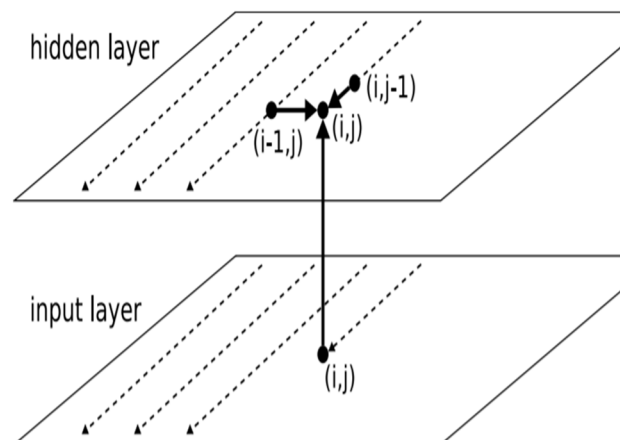


Figure show Two dimensional MDRNN. The thick lines show connections to the current point (i, j). The connections within the retired sub caste aeroplane are intermittent. The dashed lines show the scanning strips along which former points were visited, starting at the top left corner.

V. PROBLEMS IN HR

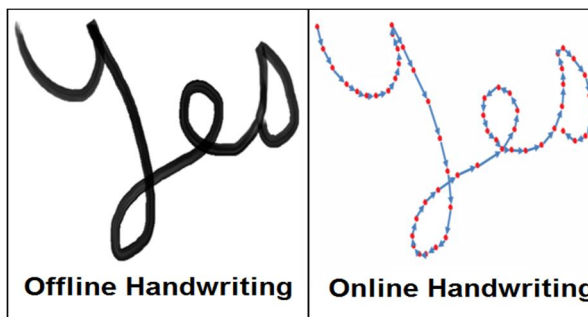
- 1) Huge variability and nebulosity of strokes from person to person
- 2) Handwriting style of an individual person also varies time to time and is inconsistent
- 3) Poor quality of the source document/ image due to declination over time
- 4) Text in published documents sit in a straight line whereas humans need not write a line of textbook in a straight line on white paper
- 5) Cursive handwriting makes separation and recognition of characters challenging
- 6) Text in handwriting can have variable gyration to the right which is in discrepancy to published textbook where all the textbook sits up straight
- 7) Collecting a good labelled dataset to learn isn't cheap compared to synthetic data Online Styles-Online styles involve a digital pen/ stylus and have access to the



VI. METHODS

Handwriting Recognition methods can be broadly classified into the below two types

- 1) *Online Methods:* Online methods involve a digital pen/stylus and have access to the stroke information, pen location while text is being written as the seen in the right figure above. Since they have a tendency to possess tons of data with regards to the flow of text being written they will be classified at a reasonably high accuracy and therefore the demarcation between different characters in the text becomes much more clear
- 2) *Offline Methods:* Offline methods involve recognizing text once it's written down and hence won't have information to the strokes/directions involved during writing with a possible addition of some background noise from the Source i.e paper



VII. TEXTUAL SIMILARITY

When talking about text similarity, different people have a slightly different notion on what text similarity means. In essence, the goal is to compute how 'similar' two pieces of text are in (1) meaning or (2) surface closeness. The first is referred to as semantic similarity and the latter is referred to as lexical similarity. The methods for lexical similarity are often used to achieve semantic similarity (to a certain extent), achieving true semantic similarity is often much more involved. In this article, I mainly focus on lexical similarity as it has the most use from a practical stand-point and then I briefly introduce semantic similarity

The conversion of handwritten content to text format after this process compare to the data set we have to find the similarity between those two

Content

VIII. ADVANTAGE

The conversion of handwriting into text format and textual similarity make human effort easy in many ways

Checking student exam's copy in colleges and school of bulk of students is hard and time consuming but with the help of handwriting conversion and textual similarity. First step is to convert the hand written content of exam copy into computer text than compare to the original context /original answers with the help of textual similarity than according to the similar content the software give marks on that copy automatically

IX. DATASETS

- 1) **IAM:** IAM dataset contains about 100 Thousand images of Words from the english language with Words written by 656 different authors. The Trained, Tested and conformed set of words written by mutually exclusive authors Link-www.fki.inf.unibe.ch/databases/iam-handwriting-database
- 2) **CVL:** The CVL dataset consists of seven handwritten documents written by about 310 actors, performing in about 83 Thousand Word crops, divided into train and test sets Link-<https://cvl.tuwien.ac.at/exploration/cvl-databases/an-off-line-database-for-pen-reclamation-pen-identification-and-word-finding/>

REFERENCES

- [1] "Reading Machine Speaks Out Loud", February 1949, Popular Science.
- [2] ^ Holley, Rose (Apr 2009). "How Good Can It Get? Analysing and Improving OCR Accuracy in Large Scale Historic Newspaper Digitisation Programs". D-Lib Magazine. <http://www.dlib.org/dlib/march09/holley/03holley.html>. Retrieved 5 Jan 2011.
- [3] ^ Suen, C.Y., et al (1987-05-29), Future Challenges in Handwriting and Computer Applications, 3rd International Symposium on Handwriting and Computer Applications, Montreal, May 29, 1987, users.erols.com/rwservices/pens/biblio88.html#Suen88
- [4] ^ Tappert, Charles C., et al (1990-08), The State of the Art in On-line Handwriting Recognition, IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol 12 No 8, August 1990, pp 787-ff, <http://users.erols.com/rwservices/pens/biblio90.html#Tappert90c>, retrieved 2008-10-03
- [5] ^ LeNet-5, Convolutional Neural Networks
- [6] ^ Milian, Mark (December 20, 2010). "New iPhone app translates foreign-language signs". CNN: Tech. www.cnn.com/2010/TECH/mobile/12/20/word.lens.iphone.app/index.html. Retrieved December 20,



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)