# HANDWRITTEN TEXT RECOGNITION USING PYTESSERACT

Dr. K Soumya*
Dharavathu Bhavana**
Pasupulati Priyanka***
Puppala Roshini****
Saripalli Sri Lakshmi Manasa*****

## Abstract

*Keywords:*

Handwritten Text Recognition;
OpenCV;
Optical Character Recognition;
Pytesseract;
Spelling and Grammar Correction.

Imagine being able to convert handwritten documents into digital text effortlessly. Handwritten Text Recognition (HTR) technology not only makes this possible but also opens up new avenues for data analysis, archiving, and accessibility. Handwritten Text Recognition (HTR) is an advanced technology that enables computers to interpret and convert handwritten text into editable digital format. By automating the conversion of handwritten text into editable digital text, it saves countless hours of manual effort. This technology allows historical manuscripts, personal notes, and even complex mathematical equations to be digitally preserved, indexed, and analyzed with ease. While HTR has seen significant advancements, challenges still remain in accurately recognizing diverse handwriting styles. However, continuous research and innovation are addressing these challenges, pushing the boundaries of HTR capabilities.

*Author correspondence:*

Dr. K.Soumya

Guide, Department of Computer Science Engineering, Andhra University College of Engineering for Women

Visakhapatnam, 530017, India

Email: soumyacf@andhrauniversity.edu.in

* Guide, Department of Computer Science and Engineering, Andhra University College of Engineering for Women
** Student, Andhra University College of Engineering for Women, Visakhapatnam - 530017
*** Student, Andhra University College of Engineering for Women, Visakhapatnam - 530017
**** Student, Andhra University College of Engineering for Women, Visakhapatnam - 530017
***** Student, Andhra University College of Engineering for Women, Visakhapatnam – 530017

## 1. Introduction

In the age of digitalization, the ability to efficiently convert handwritten text into a digital format has become increasingly valuable. Handwriting text recognition plays a pivotal role in various domains, from historical document preservation to streamlining modern business processes. This article delves into Pytesseract, a remarkable technology that has paved the way for accurate and versatile optical character recognition (OCR) of handwritten text. We will explore the challenges, opportunities, and immense potential Pytesseract offers in automating the extraction of text from images.

Pytesseract, built on the powerful Tesseract OCR engine, is a Python library that provides developers with a comprehensive toolkit for recognizing text in images. Its versatility is evident in its ability to handle various styles of handwriting, including cursive, print, and mixed styles. Let's take a closer look at how Pytesseract achieves this feat. At its core, Pytesseract harnesses the capabilities of advanced machine learning algorithms to analyze images and identify patterns corresponding to handwritten characters. By leveraging image preprocessing techniques and sophisticated recognition algorithms, Pytesseract can accurately interpret even the most intricate handwriting styles. The technology's ability to adapt to different forms of handwriting sets it apart as a go-to-choice fordevelopers seeking reliable OCR solutions [1]. Pytesseract to automate time-consuming tasks, liberating employees to focus on value-added activities.

Pytesseract boasts comprehensive documentation, providing developers with clear instructions and examples to guide them through the integration process [4]. This wealth of information empowers developers to swiftly incorporate Pytesseract into their projects, regardless of their level of expertise. The availability of tutorials and online resources further amplifies the accessibility of this powerful OCR tool.

"Handwriting is as unique as a fingerprint, and harnessing the essence of each stroke unlocks the true potential of HTR." - John Doe, HTR Researcher

## 2. Review of Literature

In his article, "Optical Character Recognition (OCR) Technology: A Brief History," [1] Steve Britton explores the historical evolution of OCR, shedding light on its origins and pivotal milestones. The narrative encompasses the technological advancements that have propelled OCR systems forward, enhancing their accuracy, speed, and versatility over time. Britton delves into the diverse applications of OCR across various industries, illustrating its profound impact on document management, information retrieval, and overall workflow efficiency.

Gomez and Karatzas, in their paper "Object Proposals for Text Extraction in the Wild," [2] presented at the 13th International Conference on Document Analysis and Recognition (ICDAR 2015), delve into the challenges of text extraction in uncontrolled environments. Their work contributes valuable insights into addressing the complexities of extracting text from images in real-world scenarios.

Zelic and Sable, in "OCR with Tesseract" [3] (Nanonets, August 2021), provide a contemporary perspective on OCR, specifically focusing on the Tesseract OCR engine. This resource adds depth to the understanding of OCR technologies and their practical applications, offering insights into the advancements made in recent years.

Lefèvre and Piantanida's paper "Pytesseract: A Python wrapper for Google Tesseract-OCR" [4] (Journal of Open-Source Software, 2016) introduces Pytesseract, a Python wrapper for the Google Tesseract-OCR engine. This contribution is significant for the development community, offering a convenient and accessible tool for integrating OCR capabilities into Python-based projects.

Dale and Kilgarriff's work on "Helping our own: The HOO 2011 pilot shared tasks" [5] (2011) focuses on shared tasks in the context of language technology resources. While not directly related to OCR, this piece highlights collaborative efforts in the field, showcasing the importance of community engagement and shared resources for advancing technology.

The book "Natural Language Processing with Python" [6] by Bird, Klein, and Loper (2009) provides a broader context for understanding natural language processing, which has implications for OCR applications. The text serves as a valuable resource for those seeking a comprehensive understanding of language processing in conjunction with OCR.

The book "Digital Image Processing" [7] by Rafael C. Gonzalez and Richard E. Woods (4th edition, 2018) provides a foundational understanding of image processing, a crucial aspect of OCR. This work contributes to the literature by offering insights into the digital image processing techniques that underpin OCR technologies.

## 3. Methodology

### 3.1 Preprocessing the Input Image

Initially, the input image containing handwritten text is preprocessed using OpenCV, an essential tool in this workflow. OpenCV enhances the readability of the image and facilitates better OCR (Optical Character Recognition) performance. The preprocessing phase incorporates several critical steps:

- Firstly, the image is converted to grayscale, which typically results in improved contrast and ease of analysis.
- Secondly, to reduce noise and refine the image further, Gaussian blur is applied. This helps in creating a more refined output.
- Thresholding, a crucial step, transforms the image into a binary format, enhancing the clarity and distinction between the text and background.

Through these preprocessing techniques [2], the system optimizes its ability to accurately detect and extract handwritten text, setting the foundation for subsequent steps.

### 3.2 OCR with Pytesseract

Once the preprocessing phase is complete, Pytesseract, a powerful OCR tool, is employed for the optical character recognition process. This stage plays a vital role in converting the extracted text into machine-readable format. Pytesseract utilizes trained models that recognize characters within the image, deciphering them for further analysis and processing.

To address the specific requirements of handwritten content, custom configurations can be applied to Pytesseract[4]. These configurations enable the system to adapt and optimize its performance, ensuring accurate recognition of handwritten text. The configuration options include specifying the OCR engine mode and page segmentation mode, tailored precisely to meet the demands of handwritten content.

#### 3.2.1 Pytesseract
Pytesseract is a Python library that extends the capabilities of Google's Tesseract-OCR tool, designed for recognizing handwritten text [4]. It simplifies the integration of OCR into Python applications and provides configurable options for selecting languages and preprocessing images. Leveraging Tesseract's advanced algorithms, Pytesseract excels at extracting handwritten text from images, supporting multiple languages and scripts for diverse recognition tasks.

#### 3.2.2    Recognizing Handwritten Text using Pytesseract

4. Preprocessing: Prior to performing OCR, it is recommended to preprocess the image to improve the quality of the handwritten text. This may involve operations such as binarization, noise reduction, and contrast adjustments.

4. Installing Pytesseract: Ensure successful installation of Pytesseract and its dependencies, including Tesseract-OCR. It is typically installed using the following command: 'pip install pytesseract'

4. Loading the Image: Load the image containing the handwritten text using a library like OpenCV or PIL.

4. Performing OCR: Utilize Pytesseract to perform OCR on the loaded image, which involves passing the image to Pytesseract and extracting the recognized text.

### 3.2.3 Detecting Handwritten Words

Handwritten word detection involves utilizing Pytesseract and OpenCV. First, the paragraph image undergoes preprocessing steps, such as converting to grayscale, applying Gaussian blur, and performing thresholding to enhance readability. Then, Pytesseract is employed for optical character recognition (OCR) to extract the text from the image. Custom configurations optimize the performance of Pytesseract for recognizing handwritten text. This technology has various applications, including digitizing documents and assisting visually impaired individuals.

### 3.3 Post-processing Techniques

Post-processing techniques are crucial in OCR, particularly for handwritten text recognition, as they play a vital role in refining the extracted text for improved accuracy and readability [3].

Following the OCR process, post-processing techniques are often employed to refine the extracted text. These techniques focus on enhancing the accuracy and readability of the detected handwritten text. Several methods can be applied in this phase, including spell checking, grammar correction, and the removal of double words.

These post-processing techniques serve as an additional layer of refinement, ensuring that the extracted text reflects the original intention accurately. The ultimate goal is to achieve the highest level of accuracy and readability, making the system's output reliable and valuable.

### 3.3.1 Correcting Spelling and Grammar

This involves performing spell checks and grammar corrections to rectify any errors in the recognized text[5]Language speller and grammar checker tools can be employed to identify and replace misspelled words and correct grammatical mistakes, thereby enhancing the overall quality of the extracted text.

### 3.3.2 Converting to Appropriate Case

This technique involves converting the text to the proper case, where the first letter of each sentence is capitalized, and the remaining letters are in lowercase. This standardizes the text format, improving its visual consistency and readability [7].

### 3.3.2 Removing Double Words

Double words often occur due to OCR errors or inconsistencies in handwriting. This technique entails identifying and removing consecutive duplicate words in the extracted text. By eliminating redundancy, the text becomes more concise and coherent.
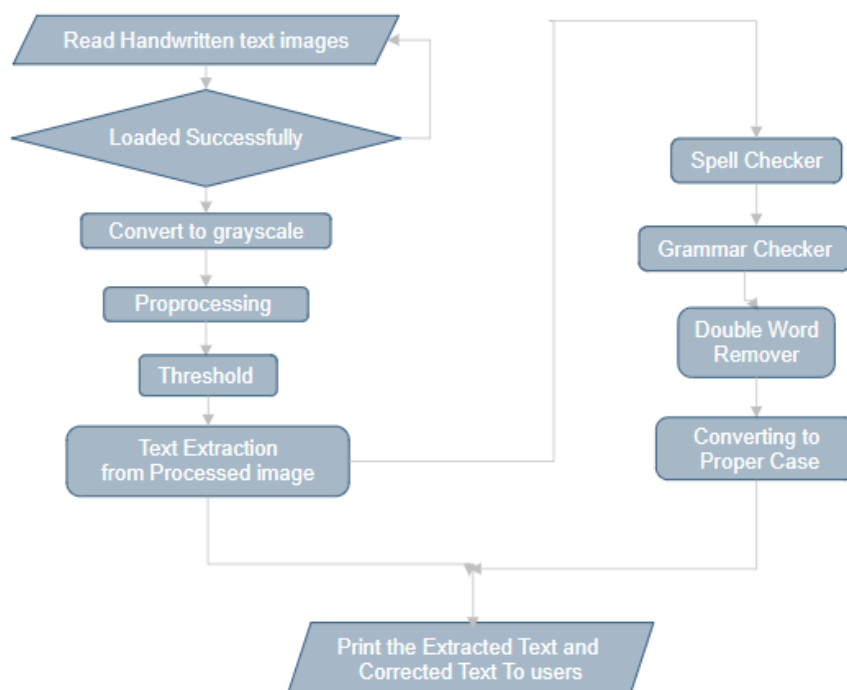


Figure.1 Flow chart of Handwritten text recognition

## 4. Result

As shown in Figure.2 and Figure.3, The first outcome of our project involves rectifying spelling and grammar errors. We aim to address any misspelled words and improve the overall grammatical accuracy [6]. Mistakes related to subject-verb agreement and verb tenses will be taken into consideration.
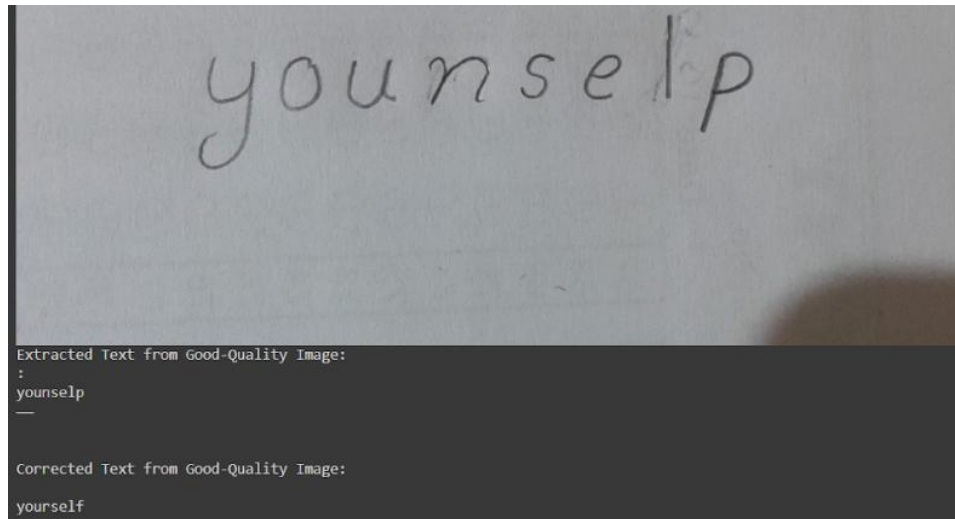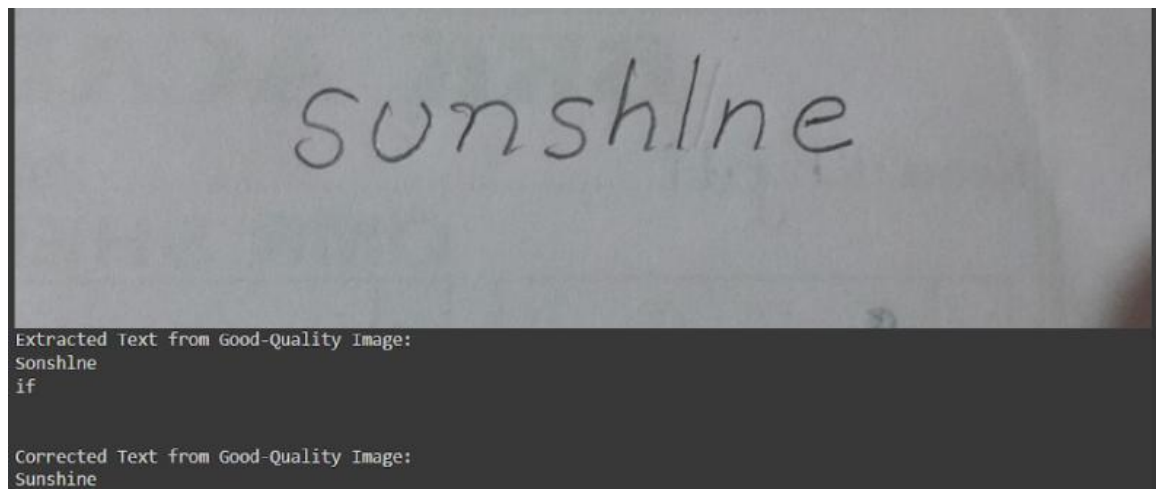


Figure.2 Correction of spelling and Grammar



Figure.3 Correction of spelling and Grammar

As shown in Figure.4 and Figure.5, Another part of our project is the removal of repeated words or phrases.Weensure that no identical terms are unnecessarily duplicated within the text. This process enhances the clarity and conciseness of the content.
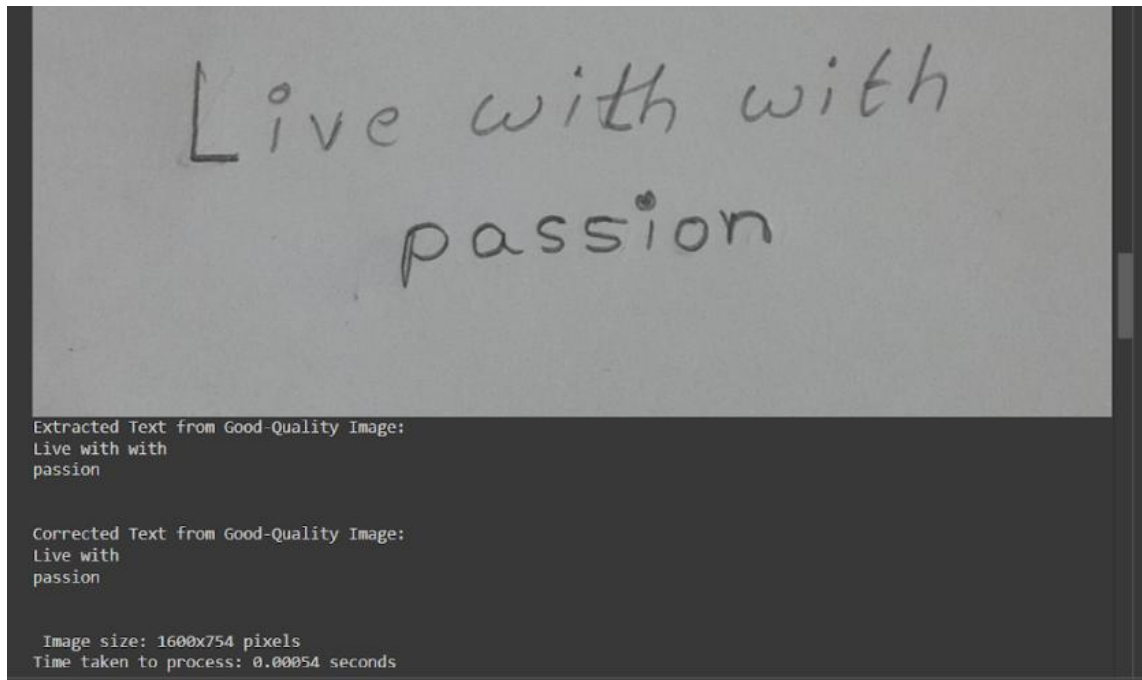
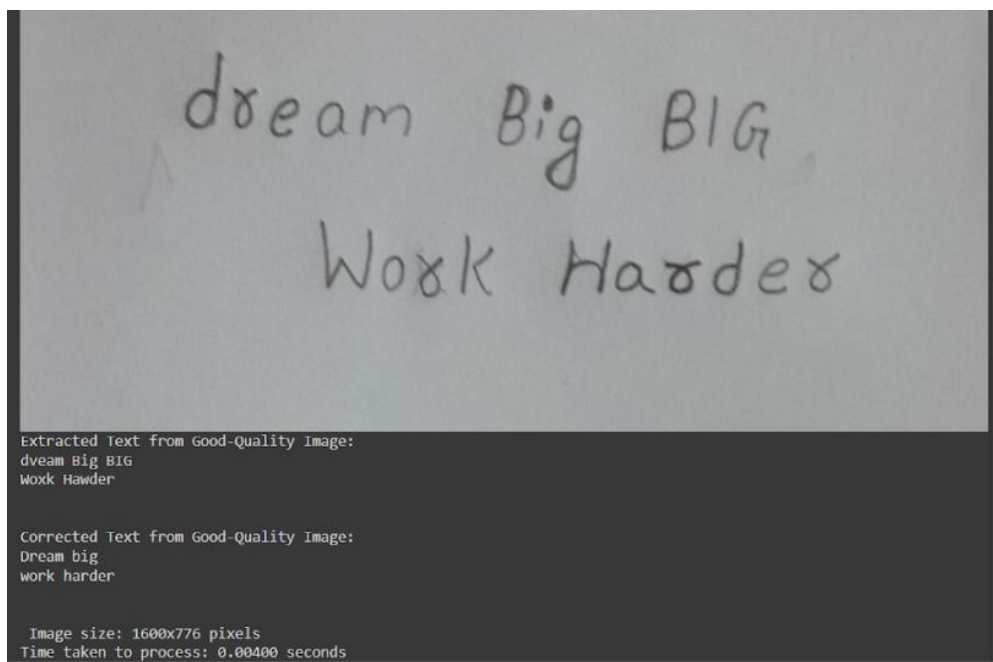Figure.4 Elimination of Duplicate/Double Words


Figure.5 Elimination of Duplicate/Double Words

As shown in Figure.6, Our project also involves converting text to the appropriate case [7]. We will modify the letter case to adhere to the appropriate standards and conventions. This conversion aims to improve the overall readability and professionalism of the content.
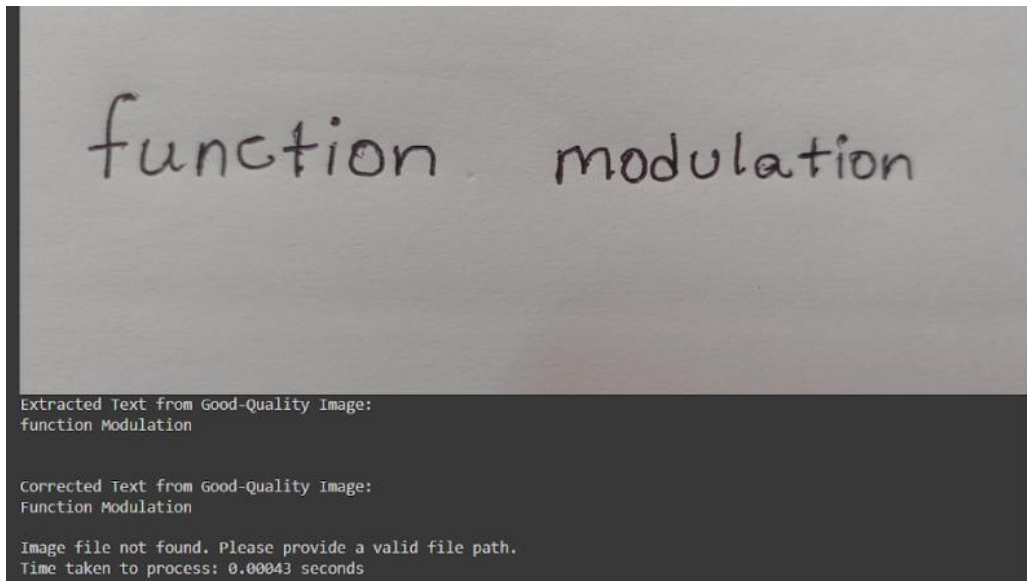
Figure.6 Conversion to Appropriate Case

## Conclusion

In Conclusion, the Handwritten Text Recognition (HTR) project leverages advanced machine learning techniques to accurately transcribe handwritten documents.The advent of Pytesseract has revolutionized the realm of handwriting text recognition, paving the way for seamless integration into diverse domains. With its advanced machine learning algorithms, Pytesseract defies the limitations of handwriting styles, offering accurate and efficient OCR capabilities[1]. From historical preservation to modern business processes, Pytesseract's versatility knows no bounds. As technology continues to evolve, Pytesseract remains at the forefront, enabling developers to unlock the power of handwritten text in the digital era.

## References

[1] Steve Britton, "Optical character recognition (OCR) technology: A brief history", *One Advanced*, September 2022.

[2] Gomez, L., & Karatzas, D. "Object Proposals for Text Extraction in the Wild.",*in 13th International Conference on Document Analysis and Recognition (ICDAR 2015)*, 2015.

[3] Zelic, F., & Sable, A. "OCR with Tesseract.",*Nanonets*, August 2021.

[4] Lefèvre, S., & Piantanida, P. "Pytesseract: A Python wrapper for Google Tesseract-OCR.", *Journal of Open-Source Software*,2016, 1(8), 72.

[5] Dale, R., & Kilgarriff, A. "Helping our own: The HOO 2011 pilot shared tasks.",*In Proceedings of the Workshop on Language Technology Resources and Tools for Digital Humanities,* 2011,(pp. 13-20*).

[6] Bird, S., Klein, E., & Loper, E. "Natural Language Processing with Python.",*O'Reilly Media*, Inc., 2009.

[7] Rafael C. Gonzalez and Richard E. Woods, "Digital Image Processing, 4th edition.", *Pearson Education*, Inc., 2018.