# Smart Handwriting-to-Digital Note Converter with Auto-Formatting for Students

Parnika Thakur[1], Chetan Thakur[2], Prajwal Telang[3], Sanket Tendulkar[4], Snehal Thawase[5], Satyam Thete[6]

[1,2,3,4,5,6]*Department of Computer Engineering, Vishwakarma Institute of Technology, Pune, Maharashtra, India*

*Abstract*—**Hand-written notes remain a common practice in classes and meetings for quickly capturing ideas. However, they are often messy and difficult to revise and reuse. For students in particular, transforming multiple pages of handwriting into clean digital text is challenging, as manually re-typing is time-consuming and handwriting OCR often is error-prone or places the original layout.**

**To achieve this, we developed an intelligent system that would transform handwritten notes into clean, editable, and well-formatted text. This idea was inspired by students who usually struggle to keep track of their notes between multiple notebooks and digital files. Our system used OCR-Tesseract, some image cleanup steps, and formatting using regular expressions. It can identify headings, bullet points, and line breaks, and set everything into a neater format.**

**The interface is simple: you upload a scanned image or photo of handwritten notes, and the system processes it. The result is a well-structured text file that can be exported as .txt or .doc, based on your need. It's useful both for students in lectures and for professionals jotting down ideas. It helps turn messy handwritten pages into organized digital notes that are easier to search, reuse, and share, with little effort from the user.**

*Index Terms*—**Handwriting, OCR, Note Digitization, Smart Formatting, Digital Notes, Students, Productivity, Revision.**

## I. INTRODUCTION

IN the digital era, learners and individuals in the professional world are increasingly relying on technology to sort information. However, there is one area that is still lagging which is the use of hand written notes. Although there are digital devices, some people still prefer to note down using their hands because they are comfortable, fast or used to it. Nonetheless, such notes are usually challenging to edit, group or recycle more so when they are not well written or in a hurry at a lecture hall or meeting.

The process of typing pages of handwritten data into the computer manually is time-consuming and subject to mistakes particularly in the quest to ensure that the information is well arranged and formatted. Although there are OCR (Optical Character Recognition) tools, the majority of them are designed to work with printed text or are not intelligent and produce unstructured and difficult to use results.

It takes a more intelligent solution than just being able to identify content written by hand with a reasonable degree of accuracy but also automatically convert it into clean and well-structured and edible digital documents. To fill this gap, a student and professional-specific system was constructed by reviewing this gap. Our tool will make use of Tesseract OCR, which extracts text and is further assisted with image processing methods to achieve better recognition. Users are able to upload pictures of handwritten pages a structured, editable output having the option of being exported as either .txt or.docs file. It is an efficient solution that enhances the digitization of the life of traditional handwritten note-taking and integrates it with the life of digital workflow. It makes it easier to digitize notes, enabling the users to pay more attention to learning and collaboration and less to tiresome formatting.

With the increasing demand of information exchange and systematic record keeping in the academics and professional settings, the inefficiencies of hand-written notes are a grave threat to productivity. Unstructured handwritten material is hard to share in group projects, during study sessions or in the workplace meeting. Its long-term usability is diminished and the good resources remain inaccessible in hardcopy pages. Although digital

tables and stylus-based applications are meant to resolve this, they need expensive hardware and continue to have challenges with the correct conversion of handwriting to text especially when using with different with accurate handwriting styles, irregular line spacing or unprofessional formatting. In versus, our solution is meant to be a hardware-free, lightweight solution that can take a scanned image or a photo on a smartphone.

Formatting intelligence is one of the main assets of our system. In comparison to the traditional OCR tools that produce plain text as their output, our tool identifies and maintains structural understanding. It is a good thing to have since students who are about to take exams need to have well-arranged notes and this can make a difference when it comes to revision. Image preprocessing in the system is done using OpenCV in order to improve the quality of the input. Text recognition is based on the Tesseract OCR and the engine is followed by a regex-based tool that reforms the raw text obtained into a read and straightforward format. Export formats will include .txt and .docs so that they can work across platforms.

The work helps narrow the analog-digital gap in note management and suggests a scalable basis on which the note management can further extend, including real-time recognition of handwriting, multi-language, and additional integration with popular apps like Google Docs or Notion. As the educational tools and workflows are becoming more digitized, the problem of this project is timely and relevant, and the solution is easy, accessible, and effective.

## II. LITERATURE REVIEW

Although there has been tremendous advancement in optical character recognition (OCR) and the transformation of documents into a digital form, the transformation of handwritten notes into editable and digital files is not a well-researched issue. Handwritten text recognition (HTR) is a field of research that has been examined by many researchers using deep learning, self-attention models, and hybrids of AI. However, very few consider such smart formatting or student features as automatic conversion of notes and sorting them.

The review of the existing techniques of handwritten recognition presented by Alhmad, Shehab, et al. [1] indicated that the models based on symmetry could

help OCR to be more accurate and structured. Their analysis allowed to find the boundaries of OCR with messy or ill-formatted handwritten notes. It also identified variations among recognition styles offline, which is significant in case of student notes.

Huh, Sen, Obaidullah, et al. [2] researched online handwritten recognition systems and the development of deep-learning designs over the last several decades. They discovered that convolutional techniques are powerful, and they are the best when trained over big and diverse data sets. Nevertheless, they observed such issues as noisy input, uneasy handwriting, and the inability of the systems to comprehend the meaning of certain text as well. To correct these, we will preprocess with regular expression rules to sanity check and add sense to the output.

A self-attention model of handwritten recognition was proposed by Nam Tuan Ly, Trung Tan Ngo and Masaki [3]. They demonstrated that transformer architectures outperform the conventional CNN+RNN configurations on long text using a better context. This inspired us to add layout-aware structuring together with OCR to keep meaning intact.

Benelarbi and Sajaj [4] developed a hybrid model of AI, which combined both the neural nets and rule-based corrections, which achieved more than 98.5% accuracy in handwritten recognition. However, their model does not include the formatting of the documents and our system will provide it through clever structuring.

Agrawal and Jagtap [5] highlighted the advantages of ensemble deep learning because it is faster and more reliable. They have found that our idea of using Tesseract OCR together with preprocessing and regular expression formatting to achieve less messy results is valid.

Suen, Legault and other [6] paved the way to future handwriting systems, including context aware cutting and multi-language support. Although the paper is not new, its main concepts are still useful to solidify the system and make it widely applicable in general.

Language models were applied to offline handwritten recognition by Zamora-Martinez and Frinken [7] and demonstrated that semantic context could be effective in enhancing accuracy. Our formatting idea, based on the spotting of keywords to group the extracted text, is supported by their work.

In brief, despite the significant improvement in recognition accuracy, there are not many solutions that

can format and extract structures. Our project goes a notch higher to include handwriting recognition with smart formatting that enables one to easily convert handwritten notes into practical electronic files.

## III. METHODOLOGY/EXPERIMENTAL

### A. Materials/Components/Flowchart/Block Diagram/Theory

The system converts handwritten notes into edited and clean digital files. It operates in three primary processes, i.e. taking the picture, attempting to read the text, and formatting it. It was written in Python and Streamlit and relies on free software like Tesseract OCR, OpenCV, Pillow and regex.

The users begin by uploading a scan or a photo of their notes. With Streamlit, we are able to build a basic web page which can run in any computer. It can take normal picture files (such as .jpg or .png) and demonstrates the results within a short time. This allows the student and worker to convert handwriting to digital text without using special devices.

Written notes are usually dark or are poorly lit, particularly on phone pictures. The image will be processed first with the assistance of OpenCV and Pillow. We put it in grayscale to enable the OCR to read it. This is followed by bluring to remove noise and yet maintain sharp text.

And then we convert the image to pure black and white to bring out the letters. Dilation and erosion are used to unite broken pieces of letters. At last, we change the picture perspective. These are measures that assist Tesseract to understand the text more.

Tesseract OCR takes care of the primary task of reading the letters. Once we are done editing the picture, we give it to Tesseract. we adjusted its options such as page mode and engine mode so that it reads whole paragraphs, rather than individual words. This preserves the original significance of the notes.
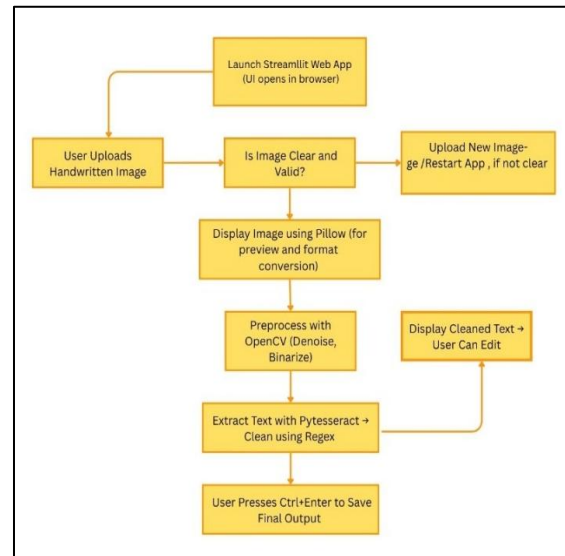
The formatting engine is another important component and it renders the text in an appearance as the original. It is based on Python regex to identify patterns, such as all-caps headings or underlined ones, bullets, lists, and paragraphs. It also eliminates OCR errors and junk symbols, transforming sloppy output into clear readable text.

Notes that are written may include either doodles or symbols not connected to the text. It is not considered by the system to make the output clean. When the handwriting is not clear or the letters are touching, certain errors may occur but the preprocessing assists in minimizing the errors. It is our goal to bring the digital notes to something that is as close to the original as possible, even though they may not be perfect.

In case, the initial attempt is not satisfactory the users can re-upload the photos and the text can be corrected. They have an option of having another snapshot or to enhance lighting and then uploading again. This makes the system handy and convenient.

After reading and formatting, the text is displayed in the screen. The digital notes can be viewed by the users and downloaded in the form of a .txt or.doc file. This allows them to make or store the notes conveniently. The feature is very handy with students, workers, and any other person who takes notes.



## IV. RESULTS AND DISCUSSIONS

Benefits:

High Usability Recognition of Real-Time Handwriting:

The system has the ability of converting the handwritten notes to properly arranged digital documents. It reads the text correctly and is very functional. It also fashions the text in a tidy manner. Success is immediately displayed as the result after uploading an image in the Streamlit interface, thus providing you with instant feedback. This is simple to operate, particularly in cases where the students and employees desire to get their notes into a digital format

without the use of massive software or special hardware.

Efficient Preprocessing Pipelines:
The preprocessing phase of the image is created using OpenCV and Pillow, which are necessary to make handwritten notes more legible. The elimination of noise, thresholding, and skew correction enhance the quality of the picture that facilitates the OCR to read the text more. It is even easier to read notes written in poor light or using plain handwriting and Tesseract can identify the characters correctly.

Tesseract Fine-tuning OCR:
Tesseract is an open-source OCR program, and when some settings are adjusted, it works fine. It becomes easier to read the text by setting the page segmentation mode (--psm) and OCR engine mode (--oem). This can be used to sustain paragraph structures, sentence structures and general accuracy particularly in long or restrictive writing. There is also a very obvious preprocessing pipeline that prevents Tesseract to confuse words or introduce too many errors.

Intelligent Auto-Format Laurence Regex:
One of the key characteristics of the system is the formatting component that can redefine the structure of the text after the recognition with the help of the use of the regular expressions (regex). It searches through such patterns as headings, bullet points, and numbered lists even when handwriting is somewhat uncharacteristic. This is useful in converting the sloppy OCR output into a more superior document. Due to this, the resulting digital text appears quite similar to the way the original handwritten notes were organized, including the right sequence and the focus.

Flexible Export Choices as well as Disciplined output:
The end product can be previewed on live basis and can be downloaded in form of.txt or.docs files, a feature that makes the system very handy. Users have the opportunity to directly edit the notes in word processors or share it across the platforms. This flexibility allows real world applications including digital archive, sharing of notes and working with academic processes.

Challenges:
Handwriting variability and the OCR limitations:

The quality of the OCR is also highly determined by the quality and clearness of the handwriting despite the good preprocessing. Tesseract does not always have the letters read out properly in case the handwritten text is quite cursive or uneven. It may lead to errors such as incorrect characters or omission of words. This demonstrates that the system is effective in most of clean handwritings but extremely stylized or poorly scanned entries might still be inaccurate.

Rule Based Formatting Limits:
Although the regexes are useful in making the text look like a document, the range of handwriting and note formats makes them less efficient. As an illustration, headings which lack clear indicators (such as underlining or all caps) might not be identified. Formatting errors can also be experienced due to overlapping bullets and spacing. It might be examined in the future that a more dynamic machine-learning-based layout analysis is possible to enhance flexibility.

Minimal Language/ Multi-line Support:
The system is currently primarily compatible with English handwritten notes and presuppose the left-to-right orientation. It does not decipher other languages and finds it difficult to do math equations or diagrams. Supporting regional languages, right-to-left writing or mixed content would enable more individuals to apply the system to more scenarios.

The interface has limitations in streamlit:
Streamlit does not have a sophisticated layout control and performance tuning despite being quick and interactive. As an illustration, the upload of large images and high-resolution scans can be delayed. Performance and user experience could be enhanced with a more scalable deployment with a Flask or FastAPI application using frontend frameworks.

V.CONCLUSION

The paper presents an intelligent system of handwritten text recognition. It makes the transition between rough handwritten and structured computer-based papers. The system converts scanned handwritten notes into a digital format that can be edited and formatted using such free tools as Tesseract OCR, OpenCV, and Streamlit.

This prevents a manual and lengthy procedure. The system employs preprocessing methods, intelligent rule based formatting and interactive web interface such that users can view the results in realtime without intense technical effort.

The key advantage of the system is that it is well-designed and easily accessible. It is no complex business product. No additional hardware and subscription are required. It is particularly convenient among students and teachers who might have poor light or have uneven handwritings due to the fact that the formatting engine does not alter the original structure.

Nonetheless, they have certain limits. Very cursive/irregular handwriting can be a problem to OCR. The formatting based on rules is effective only in case the input is rather consistent. The system is also paper-only and basic formatting. It does not allow various languages and complicated contents like math equations, diagrams, and handwritten tables.

Machine learning will be used in the future to comprehend layout and meaning. Being able to support more languages and accommodate varied writing styles will make it stronger and assist more people.

It would also be more helpful to add real-time scanning on the phone, cloud processing, and automatic saving to make it even more practical among students and workers. The project demonstrates the way in which such combination of OCR and image processing, as well as rule-based logic, can produce a very helpful and user-friendly solution.

Document storage is also facilitated by the system that converts messy unstructured hand written notes into clean and searchable, editable files. With the increased demand on seamless analog to digital transitions, particularly in blended learning and remote working, this technology is a potentially promising advancement to smart handwriting digitization.

## ACKNOWLEDGMENT

## REFERENCES

[1] S. Alhamad, A. Alabdulatif, N. Alzahrani, and A. Alshehri, "Handwritten Recognition Techniques: A Comprehensive Review," Symmetry, vol. 16, no. 6, pp. 1012, Jun. 2024. [Online]. Available: https://www.mdpi.com/2076-3417/16/6/1012.

[2] A. Gosh, M. Sen, M. Obaidullah, and P. Bera, "Advances in Online Handwritten Recognition in the last Decades," Computer Science Review, vol. 46, Nov. 2022. [Online]. Available: https://doi.org/10.1016/j.cosrev.2022.100492.

[3] N. T. Ly, T. T. Ngo, and M. Nakagawa, "A Self-Attention Based Model for Handwritten Text Recognition," Lecture Notes in Computer Science, Springer, May 2022. [Online]. Available: https://doi.org/10.1007/978-3-031-06430-2_16.

[4] M. Benelarbi and B. Sajaj, "Enhancement of Handwritten Text Recognition Using AI-Based Hybrid Approach." MethodX, vol. 11, pp. 102345, 2024. [Online]. Available: https://doi.org/10.1016/j.mex.2024.102345.

[5] V. Agarwal, J. Jagtap, A. Sharma, and P. Shah," Exploration of Advancements in Handwritten Document Recognition Techniques," Intelligent Systems with Applications, vol. 22, Jun. 2024. [Online]. Available: https://doi.org/10.1016/j.iswa.2024.200201.

[6] C. Y. Suen, R. Legault, T. M. M. Li, and M. Mai, "Building a New Generation of Handwriting Recognition Systems," Pattern Recognition Letters, vol. 14, no. 4, pp. 303-315, Apr.1993. [Online]. Available: https://doi.org/10.1016/0167-8655(93)90005-D.

[7] F. Zamora-Martinez, V. Frinken, A. H. Toselli, and E. Vidal," Neural Network Language Models for Off-line Handwriting Recognition," Pattern Recognition, vol. 47, no. 1, pp.164-176, Apr. 2014. [Online]. Available: https://doi.org/10.1016/j.patcog.2013.07.003.