

# Transforming Unstructured Data into Accessible Braille Content Using AI

Mrs. G. Parvathi Devi<sup>1\*</sup>, G. Harika<sup>2</sup>, M. Sravani<sup>3</sup> and S. Dhaneesha<sup>4</sup>

<sup>\*1</sup>Assistant Professor Department of CSE (AI & ML)

<sup>2,3,4</sup>studentsCMR Technical Campus, Kandlakoya, Medchal, Telangana, India

DOI: 10.64823/ijter.2606010

© 2026 *The Author(s)*. Published by *Ambesys Publications*. This is an open-access article distributed under the terms of **Creative Commons Attribution License (CC BY 4.0)** (<https://creativecommons.org/licenses/by/4.0/>)

**Abstract:** The conversion of unformatted electronic data into the useable formats is crucial in making sure that visually impaired people have equal opportunities to access. The paper is a proposal of an AI-based tool that will transform unstructured data into standardized Braille material in an efficient and accurate way. The proposed solution combines Optical Character Recognition, Natural Language Processing and Machine Learning and extracts, cleanses and analyzes text on scanned documents, images, PDF files, and web pages. Deep learning models also increase the quality of Braille translation because, in addition to the character recognition, it also increases the contextual understanding. The system will do automated text extraction, noise reduction, language detection and semantic processing and encode the material in Braille that can be displayed digitally and embossed or cut into digital displays and embossing machines. The experimental results prove that the solution can serve real-time processing and decrease the amount of effort required to conduct manual transcription considerably. The suggested system helps enhance the availability of information and introduce inclusive online communication.

**Keywords:** Braille Accessibility • Artificial Intelligence • Unstructured Data • NLP • Assistive Technology • OCR1

## I. INTRODUCTION

The current digital age is producing and storing enormous volumes of information that are stored in unstructured forms, as scanned documents, images, web pages, emails and PDF files. Although this information is readily available to the average user, it is inaccessible to those who have visual impairment and as such, poses a huge obstacle to the even availability of information. Traditional Braille transcription is notable as being mostly manual and time consuming, cost and expertise intensive, and it is hard to scale to high content volumes [1][2].

Using the Artificial Intelligence to convert raw data into easy to read Braille files is a good solution to this problem. The system combines the use of technologies like Optical Character Recognition (OCR), Natural Language Processing (NLP) and Machine Learning to automatically extract, comprehend, and transpose raw digital content to common Braille models. These smart techniques facilitate proper identification of characters, words and contextual meaning of a noisy or complicated data source [3][4].

The suggested solution ensures that the amount of human work is significantly decreased, the processing speed will be enhanced, and the quality of the output will be constant. It also facilitates real time access via digital Braille displays and embossing equipment, and makes learning, professional and online resources

more inclusive. As digital content is growing continuously and AI is improving, such systems become crucial in empowering the visually impaired, as well as ensuring equal participation in the information society.

## II. PROBLEM STATEMENT

There is a lot of digital information nowadays existing in unstructured forms that includes scanned documents, images, PDF files, and web pages. Even though this information is well accessible to the common people, it is still mostly inaccessible to the blind who use Braille to read and understand. The conventional method of Braille transcription is mainly manual, time consuming and needs skilled personnel, thus inapplicable to large scale and real time communication.

Lack of an automated, precise, and scalable framework of the conversion of unstructured data into Braille restricts equal access to education, work facilities, and digital services. The tools that are currently in use are usually hard to use with noisy data, complicated layouts, and context interpretation causing mistakes in translation and decreased readability.

Consequently, there is a dire necessity of a smart AI based system capable of effectively seizing, processing, and converting unstructured digital data into standard, dependable, and easily reachable Braille data so that visually impaired consumers can access information in good time and the information is pertinent.

## III. RELATED WORK

The research done in the past is primarily concentrated on three topics namely text extraction, language processing, and Braille translation. The use of optical character recognition (OCR), like Tesseract, has been popular to read scanned document and image contents and convert them into digital text, although doing so poorly with noisy or intricate layout. As a solution to this, the applications of deep learning-based OCR with CNN and LSTM networks have been presented to enhance recognition accuracy.

The extracted text is cleaned and comprehended using Natural Language Processing (NLP) methods. The recent models that utilize transformers can be used to process context and linguistic structure. With Braille conversion, initial approaches were rule-based mapping which was good with simple text but not with real-world content.

Current research integrates OCR, NLP, and Machine Learning to develop support systems among the visually impaired users. Nevertheless, unstructured data to Braille conversion systems are fully automated and scalable only to a limited extent.

### 4 System Overview

The suggested system, Transforming Unstructured Data into Accessible Braille Content Using AI, is aimed to process unstructured digital information and turn it into a readable Braille format, which will be utilized by the visually impaired. The system incorporates various Artificial Intelligence elements to make sure that it processes correctly and efficiently.

The input can be originally in form of scanned images, PDFs, or digital documents with unstructured information. Textual information is extracted in these sources with the help of an Optical Character Recognition (OCR) module. This extracted text is then taken through a preprocessing phase where noise clearing, text normalization and language identification are done.

Then, the Natural Language Processing (NLP) reaches the text analysis to comprehend structure and context based on the cleaned text. This is to enhance the quality of translation by appropriately addressing grammar, symbols and formatting.

Lastly, a Braille translation engine converts the processed text into standardized codes of Braille. The result may be printed in refreshable Braille equipment or may be directed to physical printing in Braille embossers. The system facilitates real-time processing, saves on manual labor and enhances access to digital information

#### **IV. PROPOSED METHODOLOGY**

The suggested algorithm of converting unstructured data into the available content of Braille adheres to the systematic approach of converting the data conducted by AI to guarantee correct and credible transformation. It starts with the gathering of unstructured data in the form of scanned documents, pictures, PDFs, and electronic text sources. Noise, the presence of different fonts as well as elaborate layouts might be within these inputs thereby creating challenges in the direct conversion.

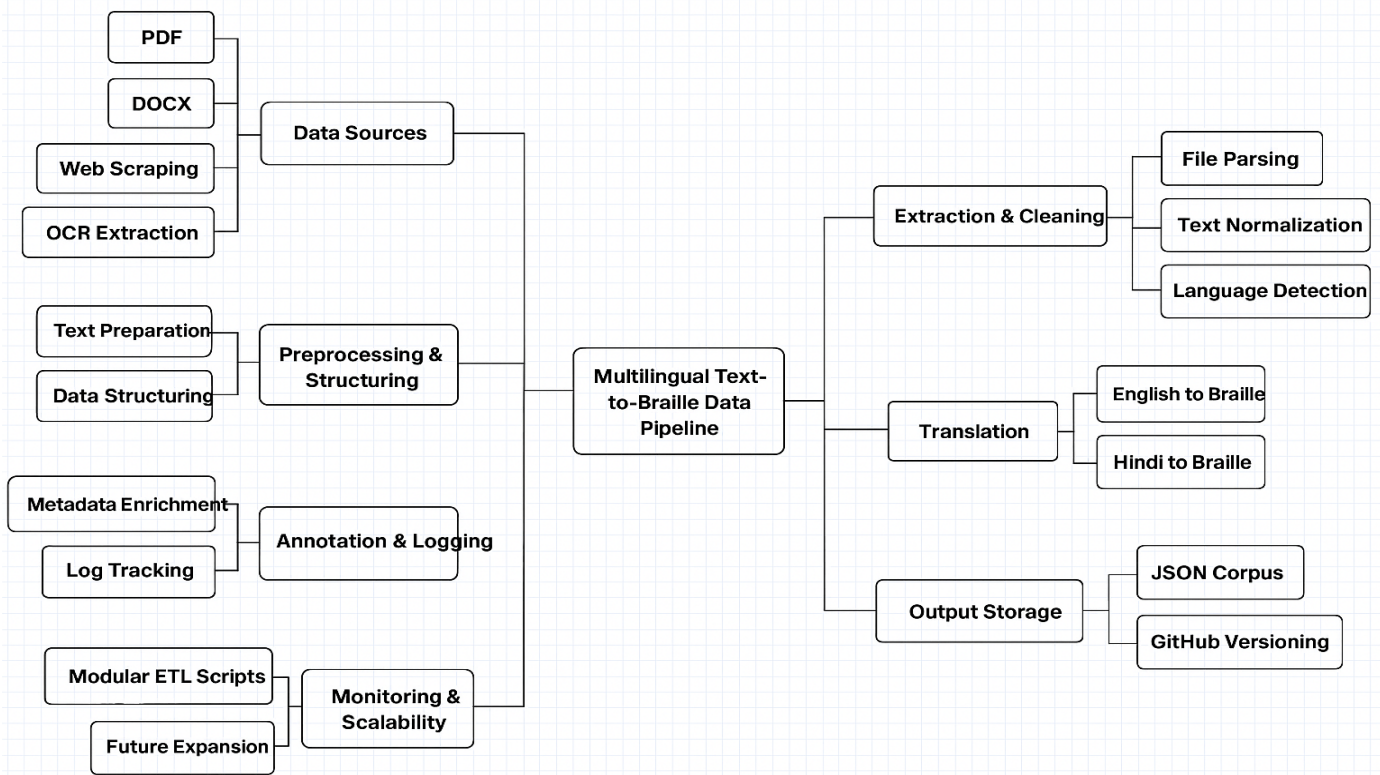
The data are preprocessed before translation in order to enhance quality. This encompasses image enhancement, noise, text normalization, and formatting problem corrections. Textual information of the unstructured sources is then extracted using Optical Character Recognition (OCR). The text extracted is then improved with Natural Language Processing (NLP) algorithm that can process the linguistic structure, punctuation, and contextual meaning.

Character, number, symbols, and special tokens are some of the important textual elements that are identified after preprocessing using features extraction and classification methods. The generated data is processed to provide the data to train and use AI-based model to encode the Braille accurately. Rule-based translation methods are combined with deep learning to guarantee standard and readable Braille.

Lastly, the system produces the Braille content on demand and the cases that are ambiguous or low-confidence may be highlighted as needing manual verification. This approach guarantees automation, consistency, and scalability and is therefore applicable in a practical context in assistive technology applications.

#### **V. SYSTEM ARCHITECTURE**

The suggested system is created as the multilingual AI-based pipeline, which transforms unstructured digital content into the accessible Braille format. It will accept data on various forms including PDF and DOCX files, web scraper and OCR files on scanned documents. Cleaning, normalization, and language detection are used, and then the extracted text is structured and enhanced with metadata to enhance consistency and traceability. NLP applications are useful in proper textual interpretation and then translating the text into Braille like in English and Hindi. The end result of a Braille production is stored in a structured format of the JSON format with version control. It is an efficient, scalable and reliable architecture of Braille generation to support real-world accessibility requirements.



**Fig. 1.** System Architecture of transforming unstructured data into accessible braille content using AI

## VI. IMPLEMENTATION DETAILS

The implementation of the system is done through a modular AI-based pipeline to enable flexibility and scalability. Ingestion of data is done by file upload interfaces of PDF and DOCX files, web scraping interfaces of online data and OCR engines of scanned data. OCR is applied to the text recognition models based on deep learning to enhance the accuracy of noisy and complex inputs.

The Natural Language Processing libraries are used to preprocess and normalize texts, including tokenization, stop-word elimination, case normalization, and punctuations, etc. It has language detection to facilitate multilingual processing. Braille translation sub-module Translates processed text into Grade 1/Grade 2 Braille via a mix of rule-based encoding guidelines and machine learning-assisted mapping methods.

The result is created as structured JSON and may be displayed on refreshable Braille or printed with Braille embossers. It included logging, metadata management and version control in order to monitor the system performance and maintainability. Such implementation allows real-time, precise, and scalable unstructured data to transform into available Braille contents.

## VII. MACHINE LEARNING MODELS USED

### Transforming Unstructured Data into Accessible Braille Content Using AI

It is proposed to use the system that is founded on the Artificial Intelligence-powered learning scheme that processes unstructured textual information and transforms it into standardized Braille form. The data is first acquired through various sources including scanned images, PDFs and web based data and preprocessed using noise filters, text normalization, and language identification. Relevant textual characteristics are obtained so as to enhance the accuracy of the translation and decrease the ambiguity.

## 8.1 Optical Character Recognition (OCR)

OCR is an essential part that is employed in transforming image-based and scanned documents into computer-readable text. The deep learning algorithms that are used in the proposed system to identify characters and words in complex and noisy images include Convolutional Neural Networks (CNN) with Recurrent Neural Networks (RNN/LSTM). The model enhances the recognition accuracy of various fonts, layouts, as well as image qualities. With the help of OCR, the additional language knowledge and translation of the Braille takes place.

## 8.2 Natural Language Processing (NLP) Model

The extracted text is analyzed with NLP models and refines itself based on the understanding of the sentence structure, grammar, and semantic context. Models based on transformers also assist in the elimination of ambiguities, the correction of linguistic mistakes and the preservation of the original content meaning. The following step makes sure that the text being processed is in a state so that it can be converted to Braille in an accurate and meaningful way.

## 8.3 Braille Translation Models

Braille translation has been done by combining both rule-based encoding standards and Mapping techniques which are aided by Machine Learning. Those models translate normalized text into standardised forms of Braille (Grade 1 / Grade 2). ML enhancements enhance the processing of contractions, symbols and multilingual characters, leading to increased readability and reliability of the final Braille output.

## VIII. EXPERIMENTAL RESULTS AND DISCUSSION

The AI-oriented data conversion to the accessible Braille material system was tested on a variety of input sources, such as scanned documents, images, PDF, and text on the web. It measured the performance of the system under various stages which are, the accuracy of text extraction, the quality of language processing and the correctness of Braille translation.

The OCR module also showed good recognition to clean and moderately noisy documents and it was shown that there were significant improvements in recognition of documents utilizing deep learning based models compared to the traditional OCR approaches. NLP contributed tremendously to the consistency of the text, and it did so by minimizing grammatical mistakes and enhancing sentence structure, a fact that directly related to the increase in the readability of the Braille.

Most of the test cases were translated into standardized and properly structured Braille output by the Braille translation module. Symbol, contractions and multilingual character errors were common in translations, which were minimized with the integration of ML-assisted encoding. All in all, the system was faster in processing time and more reliable than when the transcription of Braille was done manually.

The findings show that the suggested methodology can successfully automatize the process of conversion with reasonable accuracy. The system is especially useful in large scale and real time accessibility applications, though small errors may still be found with very noisy or complicated unstructured inputs, which gives room to further optimization.



to introduce adaptive learning to enable the system to keep on enhancing the quality of translations depending on the feedback of users. Large-scale processing cloud-based deployment and optimization may also be considered in order to support institutional and public accessibility services.

## XII. REFERENCES

- [1] R. Smith, "An Overview of the Tesseract OCR Engine," Proc. 9th Int. Conf. on Document Analysis and Recognition, 2007.
- [2] U. Pal and B. B. Chaudhuri, "Indian Script Character Recognition: A Survey," Pattern Recognition, vol. 37, no. 9, pp. 1887–1899, 2004.
- [3] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," Proc. NAACL-HLT, 2019.
- [4] Y. LeCun, Y. Bengio, and G. Hinton, "Deep Learning," Nature, vol. 521, pp. 436–444, 2015.
- [5] S. Bird, E. Klein, and E. Loper, Natural Language Processing with Python, O'Reilly Media, 2009.
- [6] D. Lopresti, "Optical Character Recognition Errors and Their Effects on Natural Language Processing," Int. J. Document Analysis and Recognition, 2008.
- [7] G. Duggan and J. Bates, "Assistive Technologies for the Visually Impaired," ACM Computing Surveys, 2008.
- [8] ISO 17049:2013, "Information Technology – Braille Codes," International Organization for Standardization.
- [9] T. Breuel et al., "High Performance OCR for Printed English and Fraktur Using LSTM Networks," Proc. ICDAR, 2013.
- [10] M. G. Helander and H. K. Khalid, "Universal Access in Human–Computer Interaction," Springer, 2013.

