

Machine Transliteration of Handwritten MODI Script to Devanagari using Deep Neural Networks

SOLLEY JOSEPH¹, JOSSY GEORGE²

¹Carmel College of Arts Science and Commerce for Women, Goa, India

²Christ University, Delhi NCR Campus, India

Corresponding author: Solley Joseph (e-mail: solley.joseph@res.christuniversity.in).

ABSTRACT The transliteration process involves transcribing words from the source language into the target language that uses a different script. Language and scriptural hurdles can be overcome via transliteration systems. There is a demand for automated transliteration systems due to the existence of several languages and the growing number of multilingual speakers. This study focuses on the Machine Transliteration of handwritten MODI script to Devanagari. MODI script was the official script for Marathi till 1950. Although Devanagari has, since then, taken over as the Marathi language's official script, the MODI script has historical significance as large volumes of its manuscripts are preserved in libraries across different parts of India. However, MODI into Devanagari transliteration is a difficult task because MODI script documents are complex in nature and there is no standard dataset available for the experiment. Machine Transliteration can be approached either as a Natural Language Processing task or as a pattern recognition task. In this research work, the transliteration task is carried out using the pattern recognition technique. The transliteration of MODI script to Devanagari is implemented using Convolutional Recurrent Neural Network (CRNN) based Calamari OCR, which is open-source software. An accuracy of 88.14% is achieved in character level matching of each word in the MODI to Devanagari transliteration process. When considering the entire word matching, the accuracy achieved is 61%. Machine Transliteration of MODI script documents results in the retrieval of large repositories of knowledge from ancient MODI manuscripts.

KEYWORDS Machine Transliteration; MODI script; Calamari OCR; CRNN; Deep Neural Networks.

I. INTRODUCTION

THE process of transcribing words from one language to another language (which uses a different script) is referred to as transliteration. It entails expressing words in one language with the phonetic or spelling counterparts of another language that are fairly accurate. Language and scriptural limitations can be overcome with the help of transliteration systems [1, 2]. The application areas of transliteration include natural language processing tasks and machine translation in particular. The prevalence of many languages, as well as the growing number of multilingual speakers, calls for the development of automated transliteration systems. This study focuses on the Machine Transliteration of handwritten MODI script to Devanagari.

MODI was intended for fast and easy writing and was adapted for writing Marathi during the 12th century. It was an official script for all the administrative work in the state of Maharashtra and was in use until 1950 [3]. The advent of printing technology in India brought down the usage of the

script significantly as it was fairly difficult to typeset the script and Devanagari was promoted as the Marathi language's official script in the 1950s.

A. MOTIVATION

Currently, the Devanagari script is used for writing Marathi, but the MODI script has significant historical importance as large volumes of MODI documents are collected and kept in reserve at various locations, both in and out of the country. They are stored in various libraries and temples in the form of official letters, land records, and other administrative documents. Significant amounts of MODI documentation are reportedly kept in India and other Asian and European nations. There are sizable collections of these records preserved at Bharat Itihas Sanshodhan Mandal in Pune, Saraswati Mahal in Tanjavur, Rajwade Sanshodhan Mandal in Dhule (Maharashtra), etc. The State Archive Department's Pune section has rare manuscripts from Chhatrapati Shivaji's reign and the Peshwa empire [4]. A collection of MODI records,

including land and income records for the state's various talukas, is kept by the Directorate of Archives and Archaeology in Goa state as well. These are unbound manuscripts that are kept in book form for preservation. (The repository contains 13,000 pages spread across 113 books in a single taluka in addition to other undocumented loose manuscripts). Due to inadequate facilities, several of these documents are about to deteriorate. The preservation of this vast knowledge base will be aided by the digitization of these texts. Historical sites also have inscriptions in MODI script in addition to manuscripts. There are some MODI document collections in many other nations, including Denmark. Furthermore, there are a vast number of individual land records that are in private possession. Some of these materials are on the brink of deterioration due to inadequate facilities [5]. Digitizing these documents will help preserve this vast repository of information. The use of the script covered a wide range of activities, including education and journalism. Due to the scarcity of MODI professionals, most of the MODI manuscripts stored in various places are still untouched and there is a need to extract the information stored in them by transliterating them to the Devanagari script. As of now, manual transliteration is the only means available to do so. The price associated with the rendering process is highly unaffordable and the fact that there are only a handful of professionals capable of undertaking the task itself makes it even more cumbersome. These factors reinforce the idea of having an automated system for Machine Transliteration (MODI to Devanagari) that would address all the existing issues and would be of great benefit to society.

A lot of work is being done to bring the MODI script back using publishing prepared scripts for printed content. These initiatives started to show promise with the publication of the script in print during the 1800s. The state of Maharashtra currently publishes a large number of books, journals, and newspapers on a regular basis. In addition to this, introductory books in MODI script were also instituted in the educational curriculum for the purpose of teaching school children. MODI enthusiasts from different parts of Maharashtra have taken various initiatives to promote MODI scripts on social media. On July 16, 2020, it was reported that a number of professionals from various fields gathered in Maharashtra to study the MODI manuscript, where they emphasized the historical significance of the manuscript.

In addition to efforts to bring the script back to life, a lot of MODI papers that are kept in different libraries are being categorized and managed. A project to encode MODI script in Unicode was started and finished by the University of California, Berkeley's Script Encoding Initiative. The MODI user community is fortunate that the MODI script may now be represented in plain text due to Unicode's inclusion of the script. This establishes a uniform framework that can be used to produce additional resources for the script [3]. Centre for Advances Computing (CDAC) Mumbai has started certain initiatives targeted at improving the script in response to the increased interest of MODI enthusiasts in learning the script. One of their initiatives is the MODI Script learning software – (which is designed for teaching MODI script to new learners of the script). There are a lot of MODI documents in Thanjavur, Tamil Nadu, and many of them were authored by intellectuals of Maharashtra. Tamil University has begun digitizing and cataloging these MODI records through the Government of India-funded effort by transforming them into a Portable

Document Format [5]. The Maharashtra government set aside 80 lakhs in 2013 to digitize MODI manuscripts. According to the reports, a significant effort is being made to preserve MODI scripts. Archives and Archaeology Department of Goa has been working on transliterating the historical literature available in the MODI script to get a clear picture of the pre-independence coastal state. According to Archives and Archaeology Department sources, approximately 2,500 documents were transcribed and more than 2500 are to be outsourced [5].

Machine Transliteration of MODI Script into Devanagari Script is necessary because of the great significance of MODI script repositories stored at various places. The historical data contained in these numerous MODI manuscripts must be revealed, hence creating an effective machine transliteration system is the need of the hour.

B. PROBLEM FORMULATION

Text recognition of MODI script and its transliteration to Marathi script is the major objective of this research work. Machine Transliteration can be approached either as a Natural Language Processing (NLP) task or as a pattern recognition task. NLP based approach demands a properly annotated text corpus. In this study, the transliteration task is carried out using the pattern recognition technique, mainly due to the absence of MODI script text corpus. Deep Learning based methods are reported as the best-performing methods in text recognition and transliteration tasks. The transliteration of MODI lines of text was implemented using Calamari OCR (Optical Character Recognition) which is Deep Learning based open-source software.

MODI to Devanagari transliteration is a very difficult task primarily due to the absence of word demarcation symbols and also because of the complexity of the script [6]. A major problem in using Deep Learning based methods is that a large volume of data is needed for training the network. In the transliteration task, the dataset includes MODI script words/lines as well as its equivalent transliterated Devanagari words/lines (ground truth). Data preparation is a lengthy and cumbersome task in such cases. Only a limited number of MODI experts are available for ground truth verification and that makes the data preparation task expensive as well. As a part of this study a database of MODI script text lines and its equivalent Marathi script 'ground truth' was generated.

The remaining part of the paper is organized as follows: the state of the art is presented in Section 2. The description of the MODI script is given in Section 3 and Section 4 details the machine transliteration system. In Section 5 the implemented model of MODI to Devanagari Machine Transliteration system is explained. Section 6 concludes the work with recommendations for future work.

II. STATE OF THE ART

Automatic transliteration is a component of a Machine Transliteration system and it has applications in areas such as multilingual text processing and cross-lingual information retrieval and extraction [7]. Machine Transliteration has mostly been approached as a Natural Language Processing problem and many researchers [8-15] have tried to solve this problem using NLP based methods and built Machine Transliteration systems for various languages other than MODI script. A brief review of various Machine Transliteration techniques was given by Kaur [1] in which the author described the existing

approaches of Machine Transliteration systems and concluded that most of the existing Machine Transliteration systems are based on statistical and hybrid transliteration methodologies. A comparison of four Machine Transliteration models, such as grapheme-based, phoneme-based, correspondence-based and hybrid models, was made by Jong-Hoon et al. [2]. Based on their comparison of the four models within the same framework and using the same data, they concluded that the hybrid and correspondence-based models were the most effective.

Some researchers chose a different approach instead of using conventional NLP techniques [16-19]. They consider text recognition and transliteration as a Machine Learning/Pattern Recognition task and try to solve the problem using pattern recognition algorithms. Recurrent Neural Network and Bidirectional Long Short Term Memory (BLSTM) are the two models that are extensively used in text recognition and transliteration, and they give promising results. Deep Learning based Machine Transliteration was performed by Soumyadeep Kundu et al. [16], using Recurrent Neural Network and convolutional sequence-to-sequence based Neural Machine Translation. They tried their method on thirteen pairs of languages. Nachum and Turner [17] implemented a Recurrent Neural Network based model for the automatic transliteration of Judeo Arabic script to Arabic script. Text recognition was performed on printed text for seven Indian languages by Chavan et al. [18] using Bidirectional Long Short Term Memory and a Multidimensional Long Short Term Memory based model. Deselaers et al. [19] performed Arabic-English transliteration using Deep Belief Networks. They stated that these models can give promising results when combined with state-of-the-art systems. Sankaran et al. [20] designed a transcription based method for Devanagari text recognition. A Bidirectional Long Short Term based method was used by them and the model was trained using feature sequence and their corresponding Unicode representation. This way, the system could learn the mapping to convert the word into the Unicode sequence. Kehri et al. [21] designed a Devanagari word recognition system using Recurrent Neural Network and BLSTM models. The Bidirectional Long Short Term Memory based transcription method was also implemented by Krishnan et al. [22] for printed word recognition of seven Indian languages such as Hindi, Malayalam, Tamil, Kannada, Telugu, Gurumukhi and Bengali. They compared it to the currently existing Optical Character Recognition system for these seven languages and reported that the above mentioned BLSTM based method was superior compared to the other methods.

Some of the researchers explored the possibility of using existing open-source software to build effective text recognition and transliteration systems. Calamari is an example of such open-source OCR and it is based on Deep Neural Network (DNN) architecture [23]. Wick et al. [24] conducted a comparative study of the OCR engines Calamari and OCRopus. Their research concluded that Convolutional Neural Networks and Long Short Term Memory (LSTM) Networks of Calamari OCR gave promising results in the case of printed documents. Reulet et al. [25] made a similar comparison of Calamari OCR (an open source) with a commercially available OCR and reported that Calamari OCR outperforms the other commercial OCR engines considerably. The experiment was performed for the character recognition of Fraktur scripts.

The extraction of information from ancient MODI manuscripts is a difficult task. A survey of the related works

shows that very limited published work reported on MODI script character recognition. Due to the complexity of the task, MODI to Devanagari (Marathi) transliteration is still an untouched area of study. MODI manuscripts stored at different locations are about to degrade and performing document analysis on them is a complex task [6].

III. DESCRIPTION OF MODI SCRIPT

MODI Script was developed in the twelfth century in Devagiri. It is based on Brahmic scripts and originated from the family of Nagari scripts. In addition to Marathi, languages such as Hindi, Konkani, Kannada, Urdu and Tamil used the MODI script for writing [3, 6]. Most of the ancient literature sources claim that Hemadri Pandit, the last king of the Yadav Family, introduced the MODI script. He drew inspiration from Sinhala cursive writing and contributed to the development of similar writing for Marathi, which involved less lifting of hand during the writing and could be written very quickly [3]. MODI was written as characters hanging on a horizontal line, which were drawn across the page. In the MODI script, there are 46 unique characters; 36 consonants and 10 vowels as shown in Figure 1 [4, 26].



Figure 1. Basic Character Set of MODI Script

It is believed that the word ‘MODI’ originates from the Marathi verb *modane* (Marathi: मोडणे), and it means “to bend or break” [3]. ‘Lekhan’ or Boru (a pen made of ‘Bamboo’) was used to write the script. Despite the fact that the MODI script is based on Devanagari, there are some significant differences between these two scripts. The differences are obvious in the lettering, rendering behaviors, and orthography of both characters. The behaviors of these letters differ from the Devanagari script in some situations, such as consonant-vowel combinations and consonant-conjuncts, which are standard features of MODI orthography. The word/sentence termination symbol was not used in MODI script and therefore word segmentation is a very difficult task in the case of MODI script documents [6, 26]. There are a number of MODI script styles available based on the time period. The various MODI script styles and time periods are presented in Table 1.

Table 1. MODI Scripts Styles

Sr. No.	Styles of MODI Script	Time Period
1	AdyaKalin (proto-MODI)	12th century
2	YadavKalin	13th century
3	Bahamanikalin	14th to 16th century
4	Sivakalin	17th century
5	Chitnisi	18th century
6	PeshveKalin	till 1818
7	Anglakalin	1818 to 1952

It is observed that a style called School Book script was also in use during the 19th and 20th centuries. It was mainly used in primary school books. MODI script has some unique features that make it different from other Brahmi based Indian scripts [2]. These are discussed in the following subsections:

i) usage of long Valenti: MODI does not distinguish between normal and long forms of the letters ‘i’ and ‘u’. The character may be represented as both ‘i’ and ‘ī’ and the character may be used for writing both ‘u and ‘ū’;

ii) Consonant-Vowel Combinations: another unique characteristic of the MODI script is when the consonants are part of a cluster: in consonant-vowel combinations. Some consonants do not change their shape when combined with certain vowels. On the contrary, some consonants appear in different shapes when they are part of a cluster (contextual form). Seven characters take contextual form in the cluster as shown in Table 2 [3];

Table 2. List of Characters that Take Contextual Forms in Combination with Certain Vowels

	regular	contextual	occurs with
THA	ਠ	ਠ	ੳ -AA, ੁ -U, ੂ -UU, ੋ -O, ੌ -AU
DA	ਢ	ਢ	ੳ -AA, ੋ -O, ੌ -AU
DHA	ਢ	ਢ	ੳ -AA, ੂ -UU, ੋ -O, ੌ -AU
PA	ਘ	ਘ	ੳ -AA, ੁ -U, ੂ -UU, ੋ -O, ੌ -AU
MA	ਘ	ਘ	ੁ -U, ੂ -UU
YA	ਘ	ਘ	ੁ -U, ੂ -UU
RA	ਠ	ਠ	ੳ -AA, ੁ -U, ੂ -UU, ੋ -O, ੌ -AU

iii) Calligraphic Uniformity: the similarity in the shape of certain basic characters of MODI script makes both character recognition and transliteration a difficult task. A group of similar-looking characters is depicted in Figure 2 [3].

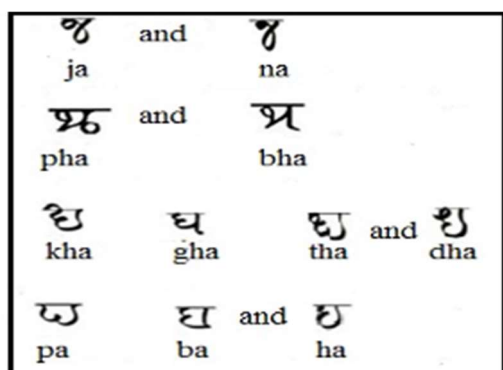


Figure 2. Shape Similarities of some MODI Characters

IV. MACHINE TRANSLITERATION SYSTEMS

The task of machine transliteration is highly difficult and time-consuming due to the numerous types of challenges that exist. The work of transliteration is an intricate task due to the fact that pronunciation differs between languages and dialects of the same language [19]. Transliteration is used as part of many multilingual applications, corpus alignment, multilingual text

processing, cross-lingual information retrieval and extraction [16]. It is a prominent research area in the field of Machine Translation, which helps to obtain accurate output [16].

Before the introduction of Deep Learning, Machine Transliteration was done mainly using different statistical methods, in a traditional way. With the advent of Deep Learning methods, researchers in the field of Machine Transliteration have implemented these techniques in Machine Transliteration tasks of different languages and have fetched better results. Some of the researchers try to solve the Machine Transliteration task using Natural Language Processing (NLP), while others choose Pattern recognition techniques. Either way, it is observed that Deep Learning techniques are being used in Machine Transliteration.

V. TRANSLITERATION USING DEEP LEARNING BASED CALAMARI OCR

Deep Neural Network-based techniques can be effectively used for Machine Transliteration tasks. The proposed model for the transliteration of the handwritten MODI manuscript to Devanagari text is shown in Figure 3.

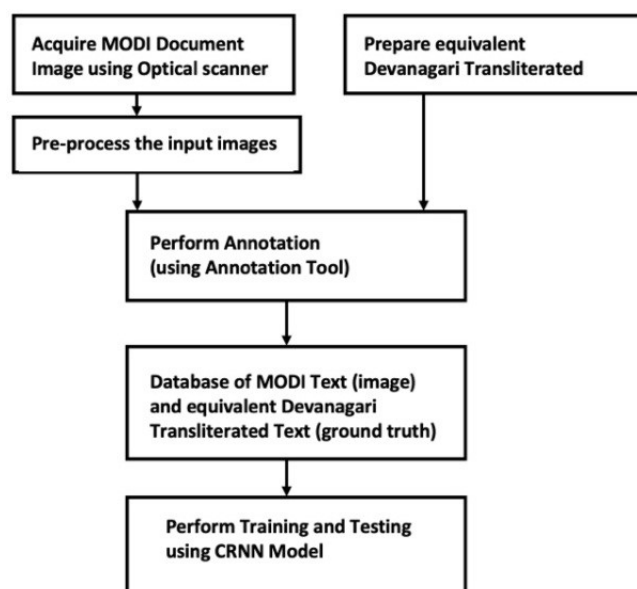


Figure 3. Process Flow of the Proposed Transliteration Model

MODI handwritten text and its Devanagari transliteration were used as the dataset for this model.

A. DATASET FOR MODI TO DEVANAGARI TRANSLITERATION

There are no standard datasets available for MODI script transliteration and therefore the required dataset was generated by the researchers. The MODI manuscripts (handwritten) and the corresponding transliterated Devanagari manuscripts were used as the dataset. The manuscripts that are used for MODI to Devanagari transliteration were gathered from various libraries such as Bharat Itihas Sanshodhak Mandal, Pune, Maratha History Museum (Deccan College, Pune) and Directorate of Archives and Archaeology, Goa. Some of the MODI manuscripts preserved in Bharat Itihas Sanshodhak Mandal were already transliterated to Devanagari script and we were

able to use the transliterated documents as our ground truth data. However, the majority of MODI manuscripts were yet to be transliterated and therefore those had to be transliterated manually by MODI experts. In addition to the above-mentioned dataset of MODI manuscripts collected from libraries, a specially prepared dataset was also generated with the help of MODI script experts from Sattari and Mumbai. The MODI script experts were given MODI manuscripts and they were asked to write the transliteration in Devanagari script documents. The MODI manuscripts were then scanned at a resolution of 300 dpi and stored as images. The Devanagari transliteration of each MODI manuscript served as the ground truth.

The scanned handwritten documents were segmented into lines and stored as images (.jpg). The scanning resolution is 300 dpi. Equivalent Devanagari transliterations of each of the MODI texts were entered using an annotation tool (which was specially developed by researchers) as the ground truth (.gt.txt) of the respective images (.jpg). The ground truth database was prepared with the help of MODI script experts during the data collection phase of the research work. The annotation tool was capable of displaying each of the MODI script images which were stored as .jpg files and the user could enter the corresponding Devanagari script (.gt.txt). Both the images (.jpg) and the corresponding Devanagari script (.gt.txt) would then be stored using the same filename (e.g., filename.jpg and filename.txt). Images of 4334 MODI script lines and their transliteration in Devanagari (in Unicode) were generated for the transliteration task.

B. PROPOSED MODEL

The proposed model uses Convolutional Recurrent Neural Network (CRNN) based Calamari OCR for the transliteration of the MODI script to Devanagari. Calamari is open-source OCR recognition software that is based on the most advanced deep neural network, implemented by Tensorflow and Python3. An extensive variety of CRNN designs can be specified by the user with a system such as Calamari. It is used for historical and contemporary fonts. It is designed for training and applying OCR models on text lines including several latest techniques to optimize the computation time and the performance of the models. Since it is based on TensorFlow, it facilitates the usage of DNN along with Convolutional Neural Networks and LSTM structures that are proven to acquire the best outputs in the case of OCR.

The customizable network architecture of Calamari OCR is comprised of a Convolutional Neural Network Long Short-Term Memory and the Connectionist temporal classification (CTC) layer. The default network of this Convolutional Neural Network model consists of a stack of two Convolutional Neural Network and Pooling-Layers, respectively and a following LSTM layer.

Fine-tuning is performed on the customizable Calamari OCR to suit the purpose of MODI to Devanagari Transliteration. The basic architecture of the network is shown below in Figure 4.

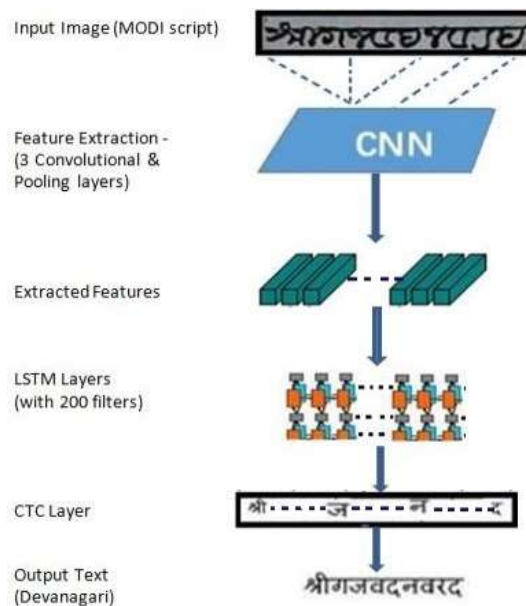


Figure 4. The Proposed Transliteration Model's Basic Architecture

The fine-tuned network in the proposed model consists of three CNN & pooling layers followed by an LSTM layer. The specification details are shown below:

- CNN layer 1 consists of 256 filters (2*2) with a dropout rate of 0.5 and a pooling layer (2*2)
- CNN layer 2 consists of 512 filters (2*2) with a dropout rate of 0.5 and a pooling layer (2*2)
- CNN layer 3 consists of 1024 filters (2*2) with a dropout rate of 0.5 and a pooling layer (2*2)
- These 3 layers are followed by an LSTM with 200 filters and a dropout rate of 0.5.

C. RESULTS AND DISCUSSION

The dataset consists of scanned images of the lines of MODI script text and its corresponding Devanagari text, as the ground truth data. The ground truth of each of the individual images is entered using the annotation tool. The annotation tool designed for the data entry is shown in Figure 5. The image which is displayed at the bottom part of the figure is the MODI script image (.jpg) and the equivalent Devanagari text (ground truth – .gt.txt) is displayed on the upper part of the figure (Figure 5).



Figure 5. The Annotation Tool's User Interface

Images of 4334 MODI script lines and their transliteration in Devanagari (in Unicode) as the ground truth were used as the dataset for the experiment. The dataset was divided into 845 for testing and 3379 for training (20:80 ratios). For training the network, the dataset of 3379 images (.jpg) and the ground truth (.gt.txt) files were used as the input to the model. After the completion of training, the prediction is performed on the test dataset (evaluation dataset) of 845 images in the test dataset. The

ground truth data and its predicted text are then compared and the total number of wrongly predicted characters is calculated (character errors). The evaluation criteria are the Character Error Rate (CER) and accuracy. The Character Error Rate is calculated using the following equation:

$$CER = \left(\frac{\text{Character_errors}}{\text{Total_characters}} \right) * 100. \quad (1)$$

The accuracy at the character level prediction is then computed based on the CER. An accuracy of 88.14% was achieved at the character level matching after fine-tuning the network. Similarly, the computation of word-level matching was also performed. If all the characters are predicted correctly in a specific word, then that word is qualified for the word level match. At the word level matching, an accuracy of 61% was achieved. A prediction sample of character-by-character matching of the predicted text against the ground truth text is depicted in Figure 6.

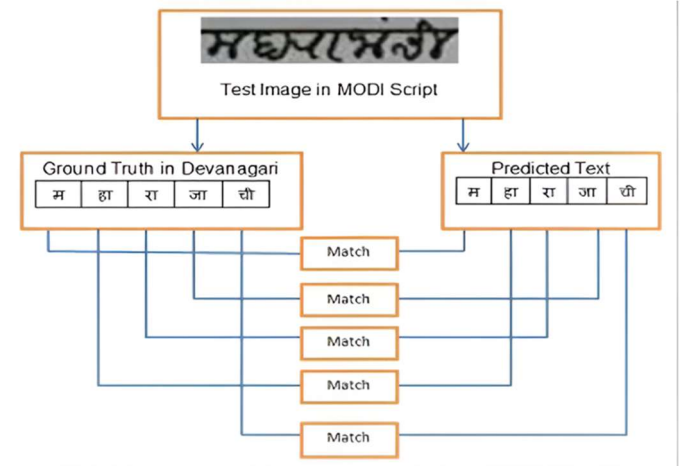


Figure 6. Sample of a Correctly Transliterated Line using the CRNN Method

A sample list of test images and their ground truth along with the predicted text is given in Table 3.

Table 3. Samples of Test Images with their Ground Truth and Predicted Text

Sr. No	Test Image (MODI)	Ground truth (Devanagari)	Predicted Text	Total no. of Chr.	Wrongly Pred. Chr.
1		हस्तलिखितात	हस्तलिखितात	10	0
2		इतिहास	इतिहास	6	0
3		भारतीय	भारतय	6	1
4		आदिशक्ती	आदिशकती	8	0
5		येणेप्रमाणे	येेप्रमाणे	10	1

The input test image is in MODI script and the ground truth is the transliteration of that text in Devanagari script. The predicted word is the output from the model.

D. TOOLS USED

The model is implemented using the Python environment with TensorFlow and Keras. Python is an interpreted language used for implementing Machine Learning methods by using Scikit-Learn and other libraries. TensorFlow is an open-source software library used for Machine Learning. It can be applied for a variety of tasks, but its main focus is on deep neural network training and inference. Keras is a Python-based open-source neural network library that runs on top of Theano or TensorFlow. Keras creates network models very quickly. The process of creating models, specifying layers, and configuring multiple input-output models is handled by the Keras High-Level API (Application Programming Interface). Calamari OCR which is open-source software is used for building the model. The customizable Calamari OCR was finetuned to suit the needs of the proposed model.

VI. CONCLUSION AND FUTURE WORK

In India, as well as in other countries in Asia and Europe, an extensive collection of MODI records has been archived. The historical information recorded in various MODI manuscripts needs to be unveiled and there is a huge demand for automating

the process of transliteration. Machine Transliteration of MODI script documents will result in the retrieval of large repositories of knowledge from ancient MODI manuscripts. In this study, a deep learning based method is used for MODI text recognition and transliteration. The transliteration of MODI lines of text is implemented using Calamari OCR, which is open-source software. An accuracy of 88.14% was achieved in character level matching of each word in the MODI to Devanagari transliteration process. When considering the entire word matching the accuracy achieved was 61%.

As a part of the study, a database of handwritten MODI text lines and transliterated Devanagari text lines were generated. Such a database is useful in training deep neural networks in the Machine Transliteration process. The lack of a standard dataset is one of the major challenges in MODI script character/text recognition based research. We aim to make this database available to other researchers by publishing it in the public domain. Machine Transliteration of MODI script documents can result in the retrieval of large repositories of knowledge from ancient MODI manuscripts. This research work will contribute considerably to the field of MODI script document analysis.

As a future work, the model used for MODI to Devanagari transliteration can be modified so that it can transliterate a complete sentence. Thus, it can be an initial effort to transliterate the entire text from ancient MODI manuscripts. The aim is to experiment with other Deep Learning based models to improve the accuracy of MODI/ Devanagari transliteration.

The lack of a MODI script data repository and an adequate number of research studies on the transliteration of MODI script to Devanagari are two of the study's limitations. A significant amount of dataset is required for the deployment of Deep Learning models in order to enhance the performance and generalization capacity of the model. The significant drawback was the lack of a consistent dataset for transliteration and character identification in MODI script.

A database of handwritten MODI text lines as well as transliterated Devanagari (Marathi) text lines were produced as another outcome of this study. A database like this is helpful for training deep neural networks for machine transliteration tasks. One of the main issues with MODI script character/text recognition based research is the absence of a publicly available dataset. By releasing this database into the public domain, we hope to enable access to it for more academics.

References

- [1] K. Kaur and P. Singh, "Review of machine transliteration techniques," *Int J Comput Appl*, vol. 107, no. 20, pp. 13–16, 2014, <https://doi.org/10.5120/18866-0061>.
- [2] J. H. Oh, K. S. Choi, and H. Isahara, "A comparison of different machine transliteration models," *Journal of Artificial Intelligence Research*, vol. 27, pp. 119–151, 2006, <https://doi.org/10.1613/jair.1999>.
- [3] A. Pandey, "Final Proposal to Encode the MODI Script in ISO / IEC 10646," pp. 26–27, 2011, [Online]. Available at: <https://www.unicode.org/L2/L2011/11212r2-n4034-MODI.pdf>
- [4] S. Joseph and J. George, "Handwritten character recognition of MODI script using convolutional neural network based feature extraction method and support vector machine classifier," *Proceedings of the 2020 IEEE 5th International Conference on Signal and Image Processing (ICSIP), Nanjing, China, 2020*, 2020, pp. 32–36, <https://doi.org/10.1109/ICSIP49896.2020.9339435>.
- [5] S. Joseph and J. George, "Feature extraction and classification techniques of MODI script character recognition," *Pertanika J Sci Technol*, vol. 27, no. 4, pp. 1649–1669, 2019.
- [6] S. Joseph, J. George, "Convolutional Autoencoder Based Feature Extraction and KNN Classifier for Handwritten MODI Script Character Recognition," in S. Shukla, A. Unal, J. V. Kureethara, D. K. Mishra, and D. S. Han (Eds) *Lecture Notes in Networks and Systems book series (LNNs)*, vol. 290, Springer, 2021, pp. 142–149, https://doi.org/10.1007/978-981-16-4486-3_15.
- [7] N. T. A. N. Le, F. Sadat, L. Menard, and D. Dinh, "Low-Resource Machine Transliteration Using Recurrent Neural Networks," *ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP)*, vol. 18, no. 2, pp. 1–14, 2019, <https://doi.org/10.1145/3265752>.
- [8] M. S. H. Ameur, F. Meziane, A. Guessoum "Arabic machine transliteration using an attention-based encoder-decoder model," *Procedia Comput Sci*, vol. 117, pp. 287–297, 2017, <https://doi.org/10.1016/j.procs.2017.10.120>.
- [9] G. S. Josan and G. S. Lehal, "A Punjabi to Hindi machine translation system," *Proceedings of the COLING'2008 22nd International Conference on Computational Linguistics, Proceedings of the Conference*, vol. 1, no. 2, pp. 157–160, 2008, doi: 10.30019/IJCLCLP.201006.0001.
- [10] E. K. Vellingiriraj, M. Balamurugan, and P. Balasubramanie, "Text analysis and information retrieval of historical Tamil ancient documents using machine translation in image zoning," *International Journal of Languages, Literature and Linguistics*, vol. 2, no. 4, pp. 164–168, 2016, <https://doi.org/10.18178/IJLL.2016.2.4.88>.
- [11] A. Das, A. Ekbal, T. Mandal, and S. Bandyopadhyay, "English to Hindi machine transliteration system at NEWS 2009," *Proceedings of the 2009 Named Entities Workshop: Shared Task on Transliteration (NEWS 2009)*, 2009, pp. 80–83, <https://doi.org/10.3115/1699705.1699726>.
- [12] P. H. Rathod, M. L. Dhore, and R. M. Dhore, "Hindi and Marathi to English machine transliteration using SVM," *International Journal on Natural Language Computing*, vol. 2, no. 4, pp. 55–71, 2013, <https://doi.org/10.5121/ijnlc.2013.2404>.
- [13] M. Alkhatib and K. Shaalan, "Boosting Arabic named entity recognition transliteration with deep learning," *Proceedings of the Thirty-Third International FLAIRS Conference (FLAIRS-33)*, vol. 6, pp. 484–487, 2020.
- [14] P. Sanjanaashree and M. Anand Kumar, "Joint layer based deep learning framework for bilingual machine transliteration," *Proceedings of the 2014 International Conference on Advances in Computing, Communications and Informatics, ICACCI'2014*, pp. 1737–1743, 2014, <https://doi.org/10.1109/ICACCI.2014.6968553>.
- [15] Y. Shao and J. Nivre, "Applying neural networks to English-Chinese named entity transliteration," no. 2011, pp. 73–77, 2016, <https://doi.org/10.18653/v1/W16-2710>.
- [16] S. Kundu, S. Paul, and S. Pal, "A deep learning based approach to transliteration," *Proceedings of the Seventh Named Entities Workshop*, 2018, pp. 79–83, <https://doi.org/10.18653/v1/W18-2411>.
- [17] N. Dershowitz and O. Terner, *Transliteration of Judeo-Arabic Texts into Arabic Script Using Recurrent Neural Networks*, 2020, arXiv preprint arXiv:2004.11405 (2020). [Online]. Available: <https://arxiv.org/abs/2004.11405>.
- [18] V. Chavan, A. Malage, K. Mehrotra, and M. K. Gupta, "Printed text recognition using BLSTM and MDLSTM for Indian languages," *Proceedings of the 2017 4th International Conference on Image Information Processing, ICIP 2017*, vol. 2018, pp. 345–350, 2018, <https://doi.org/10.1109/ICIP.2017.8313738>.
- [19] T. Deselaers, S. Hasan, O. Bender, and H. Ney, "A deep learning approach to machine transliteration," *Proceedings of the Fourth Workshop on Statistical Machine Translation*, 2009, pp. 233–241, <https://doi.org/10.3115/1626431.1626476>.
- [20] N. Sankaran, A. Neelappa, and C. V. Jawahar, "Devanagari text recognition: A transcription based formulation," *Proceedings of the International Conference on Document Analysis and Recognition, ICDAR*, pp. 678–682, 2013, <https://doi.org/10.1109/ICDAR.2013.139>.
- [21] P. Keshri, P. Kumar, and R. Ghosh, "RNN based online handwritten word recognition in Devanagari script," *Proceedings of the International Conference on Frontiers in Handwriting Recognition, ICFHR*, vol. 2018, pp. 517–522, 2018, <https://doi.org/10.1109/ICFHR-2018.2018.00096>.
- [22] P. Krishnan, N. Sankaran, A. K. Singh, and C. V. Jawahar, "Towards a robust OCR system for Indic scripts," *Proceedings of the 11th IAPR International Workshop on Document Analysis Systems, DAS 2014*, pp. 141–145, 2014, <https://doi.org/10.1109/DAS.2014.74>.
- [23] C. Wick, C. Reul, and F. Puppe, *Calamari – A High-Performance Tensorflow-based Deep Learning Package for Optical Character Recognition*. 2018, arXiv preprint arXiv:1807.02004. [Online]. Available at: <http://arxiv.org/abs/1807.02004>.
- [24] C. Wick, C. Reul, and F. Puppe, "Comparison of OCR accuracy on early printed books using the open source engines Calamari and OCRopus," *J. Lang. Technol. Comput. Linguistics*, vol. 33, no. 1, pp. 79–96, 2018, <https://doi.org/10.21248/jlcl.33.2018.219>.
- [25] C. Reul, U. Springmann, C. Wick, and F. Puppe, *State of the Art Optical Character Recognition of 19th Century Fraktur Scripts Using Open Source Engines*, 2018, p. 1810.03436. [Online]. Available at: <https://doi.org/10.48550/arXiv.1810.03436>.
- [26] S. Kulkarni, P. Borde, R. Manza, and P. Yannawar, "Review on recent advances in automatic handwritten MODI script recognition," *Int J Comput Appl*, vol. 115, no. 19, pp. 975–8887, 2015, <https://doi.org/10.5120/20257-2636>.



Dr. SOLLEY JOSEPH, an Associate Professor, Carmel College for Women. Areas of Interest: Pattern Recognition, Image Processing.



DR. JOSSY GEORGE, Director & Dean – CHRIST (Deemed to be) University. Areas of Interest: Image Processing, Biometric, Human Resource.

...