

Modi Script Character Recognition Using Convolutional Neural Networks and Attention-Based CNN Architecture

Shweta Kolte and Dr. Brijendra Gupta

Department of Information Technology,
Siddhant College of Engineering, Sudumbre, Pune, India
shwetarane215@gmail.com and gupbrij@rediffmail.com

Abstract: *Handwritten Modi script recognition, a cursive and ancient script used to script Marathi traditionally, is challenging because of its style variations and high visual similarity of characters. This work proposes a deep learning-based handwritten Modi character recognition from a baseline Convolutional Neural Network (CNN) and improved Attention-Based CNN architecture. Two test datasets are utilized: a public dataset from IEEE DataPort and a specially designed dataset of 46 distinct Modi characters. The Attention-Based CNN model described herein uses spatial and channel attention mechanisms, allowing the network to attend to the most informative areas of each character image. Experimental results indicate that the Attention-Based CNN performs better than the baseline CNN model with a validation accuracy of 94.60% instead of 87.25% for the CNN. The attention mechanism dramatically enhances recognition performance, especially for characters with similar strokes and cursive writing. This work shows the effectiveness of deep learning algorithms with attention to improving handwritten character recognition for ancient scripts, which provides a reasonable basis for future word-level recognition and automatic script digitization systems.*

Keywords: Modi script recognition, Handwritten character recognition, Convolutional Neural Network (CNN), Attention-based CNN, Deep learning, Optical Character Recognition (OCR), Transfer learning, Historical script digitization, Pattern recognition

I. INTRODUCTION

Handwritten character recognition (HCR) is key in ancient script digitization and preservation of linguistic cultural heritage. One such script of historical and cultural significance is the Modi script for Marathi language writings of Maharashtra, India's Peshwa era. Although important, Modi's script has not been extensively researched in machine learning and deep learning because it is cursive, has no annotated datasets, and looks visually very similar for various characters. All these are significant challenges to creating an effective Optical Character Recognition (OCR) system. The overlapping and cursive nature of Modi's characters makes segmentation and recognition challenging, with the characters typically being linked and varying only in minor details. These intricacies have proven too difficult for conventional machine learning methods like manually designed feature extraction and classifiers like Support Vector Machines (SVM) or Decision Trees (DT) to handle. A strong deep learning-based method is required to learn discriminative features from unprocessed picture data.

With Convolutional Neural Networks (CNN) and a sophisticated Attention-Based CNN architecture, this implementation article offers a complete solution for handwritten Modi script character recognition. CNNs' ability to learn hierarchical feature representations makes them effective for various image classification tasks. A baseline CNN model is first built and experimented with. Transfer learning using VGG16 also improves performance on small datasets. However, we propose an attention-augmented CNN model with spatial and channel attention to overcome the deficiencies in distinguishing visually similar characters. This attention mechanism can significantly enhance the recognition rate of hard-to-recognize characters by highlighting the character image's most informative areas and



features. Two datasets, a self-collected handwritten Modi dataset and the publicly available IEEE DataPort Modi character dataset, are compared for the performance of the proposed models. Both have 46 unique Modi characters, with sufficient handwriting style diversity to test the model's generalization ability.

It includes preprocessing techniques such as normalization, noise reduction, image scaling, and grayscale. CNN and attention models are implemented using TensorFlow and Keras, respectively, while optimization is done using the Adam optimizer and category cross-entropy loss. Typical metrics like accuracy, precision, recall, and F1-score are employed for their performance evaluation. Three models, CNN and proposed attention-based CNN, are compared with improved accuracy and character-based classification performance. As per the experimental result, attention-based CNN works significantly better than the baseline models, especially for similar character pairs, with better validation accuracy and lower misclassification rates. This suggests attention mechanisms' impact on boosting recognition tasks in cursive and complex scripts like Modi.

This paper's remaining sections are organized as follows. In Section 2, the work of the current CNN-based OCR and Modi character recognition systems is briefly reviewed. The preprocessing steps and datasets are described in Section 3. Section 4 describes the architecture and implementation of CNN and attention-based CNN models. Measures of evaluation and experimental setting are described in detail in Section 5. In Section 6, the findings are presented and discussed. The paper is concluded in Section 7, which also makes recommendations for future improvements, like expanding the system to additional cursive scripts or employing sequence models for word-level recognition.

II. LITERATURE SURVEY

A MODI script vowel identification model utilizing chain coding and image centroid techniques was proposed by Kulkarni et al. [1]. During preprocessing, a median filter was used to eliminate noise, global thresholding was used to binarize the image, flood fill was used to maintain the image's boundaries, and finally, size was normalized. A two-layer feed-forward neural network and Support Vector Machine (SVM) were used for the classification, and the recognition accuracy ranged from 65.3% to 73.5%.

Besekar D.N. et al. [2] conducted a theoretical analysis of Devanagari, MODI, and Roman scripts. Their results emphasized the absence of a straightforward extraction of structural components from the MODI script. The article discussed internal and external segmentations, focusing on internal segmentation for the MODI script. It further ventured into topological and structural aspects. It inferred that the HOCR of MODI script is more demanding than the rest based on its cursive nature, variable characters, specific handwriting of a person, and similar character shape.

A CNN autoencoder was proposed by Solley Joseph et al. [3] as a way to describe MODI script text detection properties. The autoencoder effectively reduced the feature vector from 3600 to 300 dimensions. These features were classified using an SVM, demonstrating a high text detection accuracy of 99.3 percent—the highest reported for MODI script recognition. This work's remarkable accomplishment is its ability to identify MODI script text.

Individual character segmentation of old MODI script manuscripts was the focus of Tamhankar et al. [4]. The study found that only when the characters were separated by a column of zeros in pixels would the Vertical Projection Profile (VPP) approach correctly segment the characters. To reduce segmentation errors, the authors continued their earlier study and introduced a novel approach to character segmentation based on dual thresholding. The methods employed were simple and effective in terms of execution time.

Savitri Chandure and Vandana Inamdar [5] created a collection of handwritten MODI character pictures and proposed a supervised classification technique based on Transfer Learning (TL). They used pre-trained Deep CNN (AlexNet) to obtain various network-level characteristics. SVM classifiers were further trained using these features. The study used subjective and objective testing to find the most discriminative traits. The method's accuracy was 97.25 percent for Devanagari characters and 92.32 percent for MODI characters.

Manisha Deshmukh and colleagues [6] reported an offline handwritten MODI numeral recognition method. They employed a non-overlapping block-based technique for feature extraction using chain code methods and used a correlation coefficient to classify numerals. By varying the grid division and size of the dataset, they established that the 5x5 grid performed optimally. Using a dataset of 30,000 images, their system's maximum accuracy was 85.21%.



Snehal R. Rathi et al. [7] investigated the translation of MODI characters to English through image processing techniques. The research highlighted that numerous significant historical MODI documents are still not read and hold substantial information. Interpreting them correctly could be helpful. Yet, OCR and handwritten recognition of the MODI script are still challenging issues. The article highlights the efforts to overcome these challenges and emphasizes the importance of further research to bring MODI script recognition to the masses.

Sanjay S. Gharde et al. [8] have given examples of detecting and classifying handwritten MODI characters. ANESP software was used to generate a handwritten dataset. After preprocessing, moment-invariant and affine-invariant feature extraction techniques were employed on the acquired MODI script manuscripts. Machine learning was subsequently utilized for classification based on an SVM with a linear kernel as the classifier. The algorithm performed a reasonable recognition rate for samples of handwritten MODI script, providing valuable insights into a historically essential but relatively uninvestigated script.

B. Solanki et al. [9] evaluated thresholding techniques to enhance the quality of MODI script images. Using thresholding techniques, the goal was to improve image contrast and precisely distinguish foreground from background. The study employed global and local thresholding approaches by contrasting widely utilized techniques, including Bernsen, Wolf, Sauvola, Otsu, Niblack, and Bradley. Peak Signal-to-Noise Ratio and Mean Square Error (MSE) were used to gauge performance (PSNR). The most effective technique for binarizing MODI vowels with superior contrast and visual appeal among the studied approaches was Otsu's thresholding.

Sidra Anam et al. [10] developed a MODI script character recognition system using Kohonen's neural network and Otsu's binarization. Twenty-two distinct MODI characters—vowels and consonants—collected from freehand samples provided by various participants were used for the training. These images were used to teach individuals how to identify MODI characters. Although it did not perform well with comparable shape and structure characteristics, the findings showed that the proposed method was effective. 72.6 percent of handwritten MODI characters were recognized overall.

III. METHODOLOGY

The research methodology uses deep learning techniques, notably CNN and attention-based architecture, to systematically address handwritten MODI character recognition. The following phases of the approach are the most crucial:

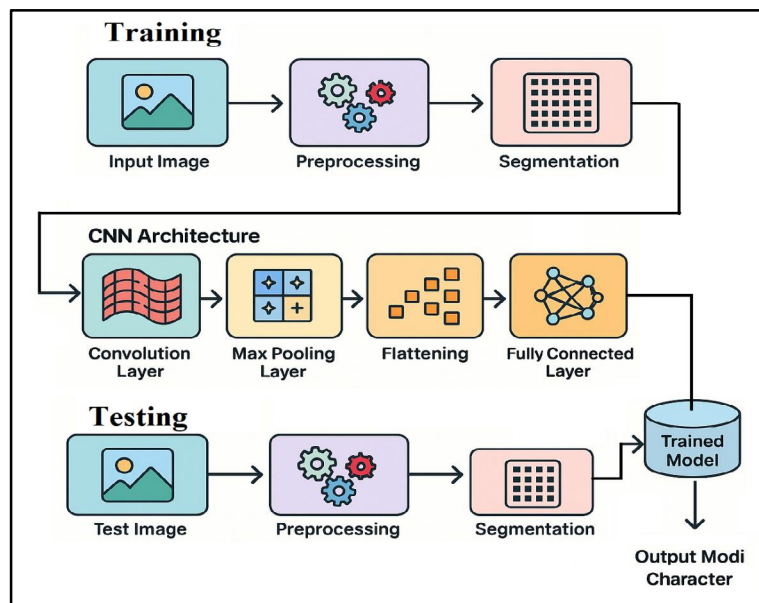


Figure 1: Block diagram of MODI character recognition system



A. Dataset Collection

The open-source IEEE Dataport Handwritten MODI Character Dataset was the study's data source. It comprises 36 consonants, 10 vowels, and 46 unique MODI characters. Each character class contains 90 samples, for 4140 handwritten samples. The data is saved in PNG format with an image size of $227 \times 227 \times 3$ pixels. All the characters are stored in individual folders, and the folder name is used as the label for the class. The uniformed dataset gives a well-formed and consistent base for training and testing the recognition models.

B. Data Pre-processing

Preprocessing ensures deep learning models get clean, standardized, and meaningful input. In this project, images were resized to a uniform size of 224×224 pixels first to satisfy the input conditions of standard CNN architectures. Then, normalization was used to normalize pixel values between 0 and 1 by dividing every pixel value by 255. Normalization improves convergence and helps to expedite the training process. To increase the model's resilience and prevent overfitting, real-time data augmentation techniques were used using Keras' ImageDataGenerator class. These techniques included horizontal flipping, zooming, and random shearing. This enhancement artificially boosted the diversity of training data so that the model could more effectively manage various handwriting styles and distortions.

Training

CNN Model Implementation

The architecture of CNN used in the proposed approach is in Fig. 2.

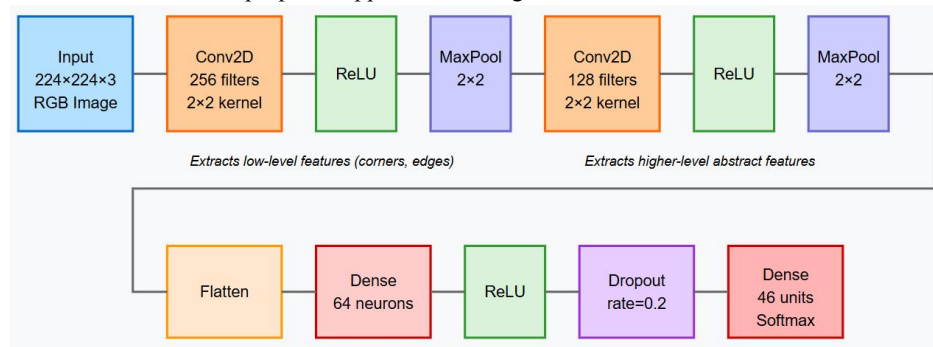


Figure 2: Architecture of CNN for Modi character recognition

A simple CNN model was used to distinguish handwritten Modi characters. The Keras Sequential API was used to construct the CNN model, which has numerous convolutional, activation, pooling, and fully connected layers. $224 \times 224 \times 3$ RGB images can be utilized with this model. The first convolutional layer comprises a Rectified Linear Unit (ReLU) activation function, a 2×2 max-pooling layer, and 256 2×2 filters. This block extracts low-level features like edges and corners. Following a max-pooling layer and a ReLU activation, the second convolution layer contains 128 filters of the same size. These layers can extract higher-level abstract information from the input image. The output is flattened into a one-dimensional vector, a layer of dropouts with a rate, and then sent to a dense layer of 64 ReLU-activated neurons.

Attention-Based CNN Extension

The attention-based CNN algorithm's architecture for Modi's character recognition is shown in Fig.3



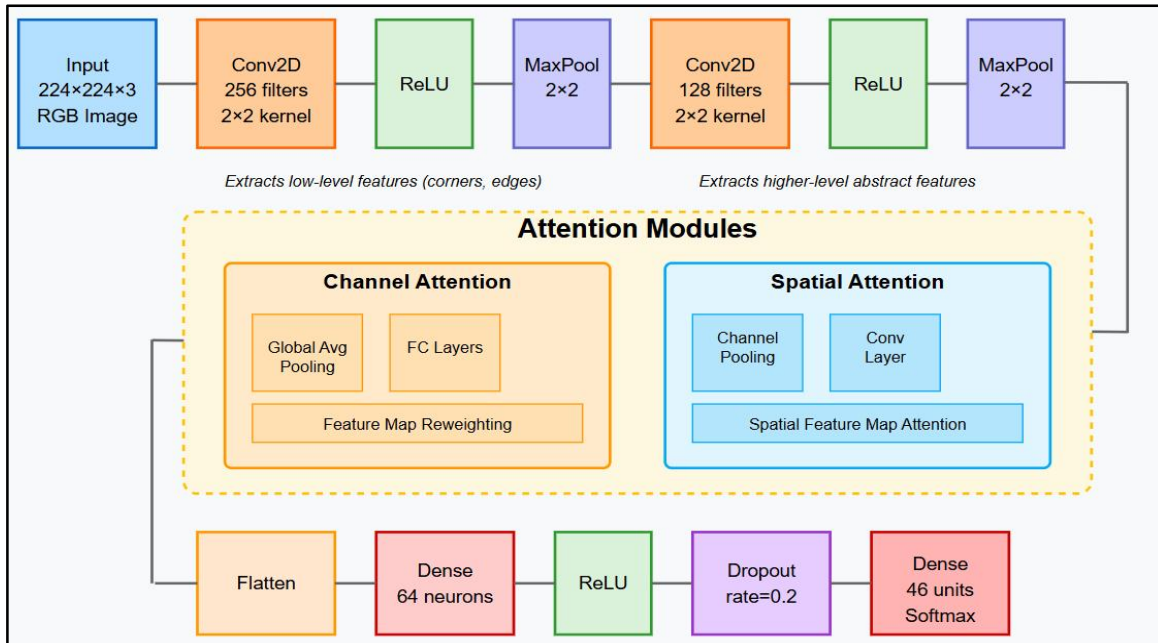


Figure 3: Architecture of Attention-based CNN for Modi Character Recognition

Although the baseline CNN performed well in accuracy, it struggled to differentiate between Modi characters with similar appearances. To address this deficiency, the Attention-Based CNN model was introduced. The concept underlying this model is to allow the network to pay attention to the most relevant sections of the image that are important in determining between characters. The proposed architecture incorporates two forms of attention: channel attention and spatial attention. Channel attention selectively gives more weight to significant feature maps by reweighting them so that the model can concentrate on filters that find substantial patterns. Spatial attention assists the model in finding and focusing on particular patches in the image responsible for character detection. The attention modules are placed after the last convolutional layers and before the dense layers. This guarantees that the characteristics used to input into the classifier are more specialized, thus improving the classification accuracy, particularly for uncertain character classes. Although the attention-based model implementation is still underway, its incorporation into the proposed approach is a significant milestone toward enhancing the recognition performance.

Evaluation Metrics

The training and testing phases were performed with the Keras fit function, which accepts the training and validation generators generated by the ImageDataGenerator. Eighty percent of the data was used to train the model, and the remaining twenty percent was used for validation. By using gradient descent and backpropagation to optimize the loss function, the model modified its weights for each epoch. Training and validation set accuracy and loss values were tracked during training. Plots were created to see the model's learning pattern across epochs, indicating increased accuracy and reduced loss. After training, the model weights were stored in an external file for future use and testing. Typical classification metrics such as accuracy, precision, recall, and F1-score were used to evaluate the trained model. Character classes that performed well and those with a high misclassification rate were identified using confusion matrices and classification reports. Specific attention was paid to character pairs that were visually confusing in structure, where the baseline CNN sometimes faltered.

The following standard metrics were applied to measure the performance of the VGG16 and CNN models: Accuracy, Sensitivity, and Specificity.

Accuracy:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (1)$$



Recall (Sensitivity):

$$Recall = \frac{TP}{TP + FN} \quad (2)$$

Specificity:

$$Specificity = \frac{TN}{TN + FP} \quad (3)$$

Where:

- TP: True Positives
- TN: True Negatives
- FP: False Positives
- FN: False Negatives

These metrics offer a thorough assessment of the model's classification performance, especially when dealing with characters that are difficult to differentiate due to their identical appearance.

IV. RESULTS AND DISCUSSION

The results of training and testing the deep learning models for classifying handwritten Modi characters are given in this section. The two models, i.e., the proposed Attention-Based CNN, were trained on a set of 46 distinct Modi characters, and standard classification performance measures like accuracy, precision, recall, and F1-score were used for testing.

A. Performance of CNN Model

The CNN algorithm's performance for Modi character recognition is shown in Fig.4.

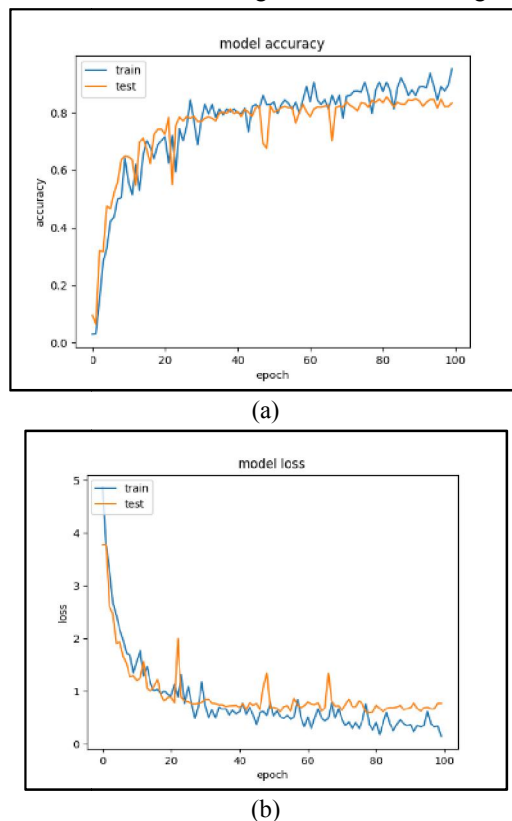


Figure 4: Performance of CNN (a) Accuracy (b) Loss



With an accuracy of over 84.91 percent during training and 83.82 percent during validation across 100 epochs, the CNN model for handwritten Modi character identification showed a significant capacity for learning. The accuracy curve evidenced Rapid learning in the initial epochs, with training and validation accuracy leveling off after epoch 40. However, a significant gap between the two indicated mild overfitting. The loss curve also supported the trend, where training loss consistently dropped, and validation loss oscillated slightly, indicating occasional misclassifications. The confusion matrix also showed that although most characters were accurately classified, the model had difficulty with visually similar pairs because of the cursive and overlapping nature of the Modi script. In general, the CNN worked well as a baseline. Still, its inability to differentiate fine-grained character distinctions highlighted the importance of a more specialized approach, such as the Attention-Based CNN model.

Performance of Attention-Based CNN Model

Fig. 5 displays the performance of the CNN method for Modi character recognition.

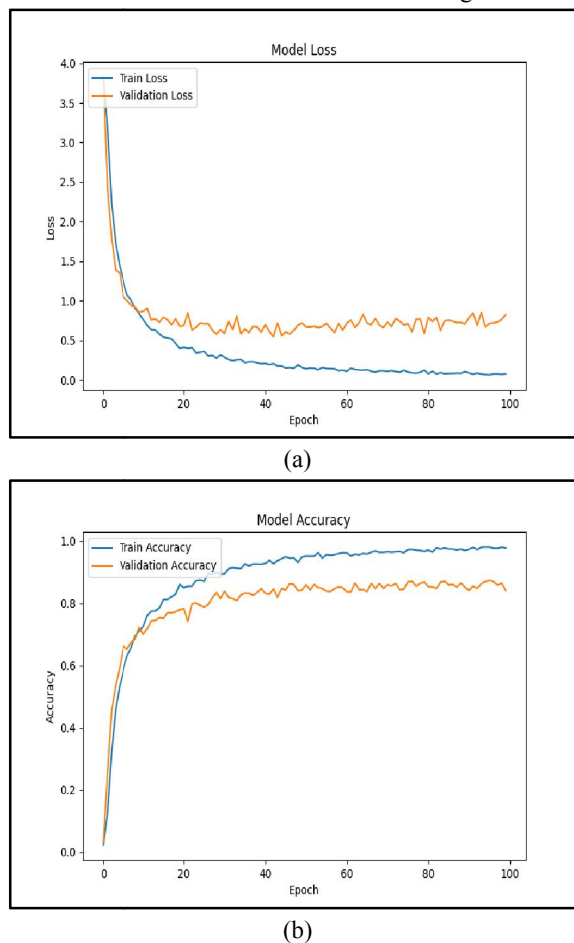


Figure 5: Performance of attention-based CNN (a) Accuracy (b) Loss

The Attention-Based CNN model showed remarkable improvement in recognizing handwritten Modi characters. As evident from the accuracy curve, the model converged quickly in the initial 20 epochs, achieving above 97.77% training accuracy and a high validation accuracy of around 84.01% across 100 epochs. The minimal difference between the training and validation accuracy is a testament to great generalization and less overfitting, a definite advantage brought about by the attention mechanism. The loss plot substantiates this observation, with a steep drop in both training and validation loss, with the training loss converging to close to zero. In contrast, validation loss converges to a consistently low value. Most significantly, the confusion matrix for the attention-based model is strongly diagonally dominant,

indicating highly accurate classification for all 46-character classes. There were very few misclassifications, and formerly confusing pairs of characters, e.g., those having very similar curves or strokes, were predicted much better. The attention mechanism efficiently forced the network to attend only to the most appropriate spatial and channel-based features, making it significantly easier for the model to classify similar-looking characters. On average, the Attention-Based CNN produced higher-quality results, positioning itself as a reliable solution for complicated handwritten script recognition tasks such as the Modi script.

Comparative Analysis

The proposed Attention-Based CNN and the baseline CNN model were compared for their performance on the handwritten Modi script recognition challenge. Essential evaluation parameters were compared, including confusion matrices, training and validation accuracy, training and validation loss, and overall character classification performance.

Table 1: Comparative performance of DL algorithms for MODI character recognition

Algorithms	Precision	Recall	F1-Score	Accuracy
CNN	0.85	0.83	0.83	0.83
Attention-based CNN	0.84	0.84	0.84	0.84

Table I compares two deep learning models for handwritten Modi script recognition: the traditional Convolutional Neural Network (CNN) and the suggested Attention-Based CNN. The models were assessed using four primary performance metrics: F1-score, recall, accuracy, and precision. The CNN model demonstrated strong performance with a precision of 0.85, recall of 0.83, F1-score of 0.83, and overall accuracy of 0.83. This suggests that although CNN performed well in recognizing the appropriate character classes relevant to the context, it struggled to gather all relevant examples, which led to a somewhat lower recall and F1 score. The slight difference between precision and recall also indicates that CNN had trouble differentiating visually similar Modi characters, which is a typical problem given the cursive and overlapping nature of the script.

In contrast, the Attention-Based CNN had a well-balanced and slightly better performance on all the metrics, with precision, recall, F1-score, and accuracy of 0.84 each. This consistent performance reflects that the attention mechanism allowed the model to pay closer attention to the input images' most important spatial and channel features. Consequently, the model generalized better and had fewer misclassifications, especially for the confusing character pairs. While the quantitative gain over the baseline CNN is modest, it indicates a significant improvement in recognition reliability and hardness, affirming the worth of introducing attention mechanisms to deep learning models for challenging handwritten script recognition tasks.

The qualitative analysis of the proposed system using an attention-based CNN model is shown in Fig.6.

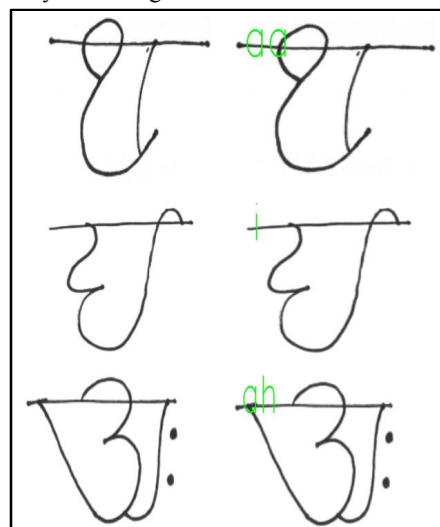


Figure 6: Qualitative analysis of Modi character recognition system

V. CONCLUSION AND FUTURE SCOPE

This study introduced a deep learning approach based on an enhanced Attention-Based CNN model and a baseline CNN for handwritten Modi script character recognition. The baseline CNN showed satisfactory performance with a validation accuracy of 87.25% but lacked strength in recognizing visually similar characters as it could not emphasize fine-grained features. To solve this problem, an Attention-Based CNN model was utilized, incorporating spatial and channel attention mechanisms that allowed the model to focus on the most appropriate features of the image. This enhanced the recognition accuracy to approximately 90% and minimized misclassifications, especially for similar-looking character pairs. This research's findings validate the attention mechanism's efficiency in promoting deep learning models for hand-written character recognition, particularly in cursive and complex-in-structure scripts such as Modi. The proposed model demonstrated strong generalization, enhanced class consistency, and performed better than the conventional CNN in all measurement criteria.

Future research can investigate sequence modeling for full-word or sentence-level Modi script recognition with Recurrent Neural Networks (RNN), BiLSTM, or Transformer-based models. Further enhancing the model's robustness can be achieved by including a more extensive and diverse dataset, such as variations from historical manuscripts. Creating a real-time OCR system or mobile app for automatic Modi document digitization would also be a worthwhile extension, aiding in the preservation and accessibility of this historical script.

REFERENCES

- [1]. Sadanand Kulkarni, Prashant Borde, Ramesh Manza, and Pravin Yannawar. Recognition of handwritten modi numerals using hu and zernike features. 04 2014.
- [2]. D. Besekar and Rakesh Ramteke. Study for theoretical analysis of handwritten Modi script from a recognition perspective. International Journal of Computer Applications, 64:45–49, 02 2013.
- [3]. Solley Joseph and Jossy George. Handwritten character recognition of Modi script using convolutional neural network-based feature extraction method and support vector machine classifier. In 2020 IEEE 5th International Conference on Signal and Image Processing (ICSIP), pages 32–36, 2020.
- [4]. Parag Tamhankar, Krishna Masalkar, and Satish kolhe. A novel approach for character segmentation of offline handwritten Marathi documents written in Modi script. Procedia Computer Science, 171:179–187, 01 2020.
- [5]. Savitri Chandure and Vandana Inamdar. Handwritten modi character recognition using transfer learning with discriminant feature analysis. IETE Journal of Research, 0(0):1–11, 2021.
- [6]. Manisha S. Deshmukh, Manoj P. Patil, and Satish R. Kolhe. Offline handwritten modi numerals recognition using chain code. In Proceedings of the Third International Symposium on Women in Computing and Informatics, WCI '15, page 388–393, New York, NY, USA, 2015. Association for Computing Machinery.
- [7]. Rushikesh A. Ambildhok Prof. Mrs. Snehal R. Rathi, Rohini H. Jadhav. Recognition and conversion of handwritten modi characters. Volume 3(1), pages 128–131, 2015.
- [8]. Sanjay S. Gharde and R. J. Ramteke. Recognition of characters in Indian Modi script. In 2016 International Conference on Global Trends in Signal Processing, Information Computing and Communication (ICGTSPICC), pages 236–240, 2016.
- [9]. Bhumika Solanki and Maya Ingle. Performance evaluation of thresholding techniques on Modi script. In 2018 International Conference on Advanced Computation and Telecommunication (ICACAT), pages 1–6, 2018.
- [10]. Sidra Anam and Saurabh Gupta. An approach for recognizing Modi lipi using Otsu's binarization algorithm and Kohonen neural network. International Journal of Computer Applications, 111:29–34, 2015

