



Optical Character Recognition for Telugu Handwritten Characters

Dr.A.S.Narasimha Raju¹, N.Sandeep Kumar², E.Nikhil Reddy³, K.Nithin⁴

Computer Science Engineering (DS) Institute of Aeronautical Engineering Dundigal, Hyderabad¹⁻⁴

Abstract: In the era of digitization, artificial intelligence has revolutionized the way we process and analyze data. However, a significant portion of historical documents and cultural heritage remains in handwritten form, inaccessible to digital technologies. Optical Character Recognition (OCR) emerges as a crucial solution, enabling the conversion of optical text into digital format, thereby making it editable, searchable, and electronically storable. This technology is vital for organizations and individuals dealing with vast amounts of textual information. By training OCR engines on diverse languages, including Telugu, we can tap into the rich cultural heritage of India's classical languages. Telugu OCR, in particular, facilitates the preservation of hand-written notes, ancient manuscripts, and historical documents, making them accessible to a broader audience. This digital transformation not only preserves cultural heritage but also enables the dissemination of knowledge and ideas to a wider audience, promoting cultural exchange and understanding.

Keywords: Handwritten Telugu Character Recognition, Optical Character Recognition (OCR), Neural Networks, Deep Learning, Image Processing, Pattern Recognition, Machine Learning

I. INTRODUCTION

Telugu, a classical language of India, has been spoken in South India for centuries, with approximately 96 million speakers in India. Despite its significance, Optical Character Recognition (OCR) systems for Telugu text have not seen substantial progress in recent years. The complexity of Telugu characters, comprising multiple connected symbols, poses a significant challenge for OCR systems. The process of OCR involves converting an image of a document into a text-editable format, enabling its utilization in various applications. The vast amount of Telugu text available online, including images, scanned documents, and PDFs, can be converted into searchable and editable formats using Telugu OCR, facilitating easy information retrieval and analysis.

In regions where Telugu is widely spoken, such as Andhra Pradesh and Telangana, Telugu OCR can be employed by government agencies, educational institutions, and businesses for administrative tasks like document digitization, data entry, and archival purposes. The recognition of Telugu text is a crucial aspect of Telugu OCR, involving the identification of connected components, referred to as glyphs, in an image. These interconnected parts are essential for understanding the structure of Telugu characters.

The primary objective of this research is to develop an efficient deep learning model that incorporates interrelated tagging for Telugu Optical Character Recognition (OCR) with segmentation. Telugu OCR involves recognizing characters from the Telugu script, which is predominantly used in Andhra Pradesh and Telangana. Segmentation, in this context, refers to the process of identifying individual characters within a given image or document. Traditional OCR systems often struggle with scripts like Telugu due to their complex and connected nature, where characters may be intertwined or joined together in a single continuous stroke. Additionally, the presence of noise, varying writing styles, and font sizes further complicates the recognition process.

To address these challenges, an advanced deep learning model is needed, capable of accurately identifying individual characters and understanding the interrelations between them within words and sentences. Incorporating segmentation techniques is crucial to precisely isolate and recognize each character, even when they are interconnected. The overarching goal is to create a prototype that can efficiently and accurately recognize Telugu characters from scanned documents, handwritten notes, or images, even under challenging conditions such as varying fonts, sizes, and styles. By achieving this, the prototype could have wide-ranging applications in fields like document digitization, language processing, and automated text extraction, thereby facilitating efficient data management and information retrieval in Telugu-speaking regions.



II. LITERATURE REVIEW

Srinivasa Rao Dhanikonda et al. created An Efficient Deep Learning Model with Interrelated Tagging Prototype with The objective is to develop a high-performance Telugu text Optical Character Recognition (OCR) system using Deep Segmentation for Telugu Learning, Optical Character Recognition deep learning based Telugu OCR: A survey Comprehensive study of Develop an accurate OCR system for processing document with historical records, Library archives and application forms. OCR for text recognition in images is extremely difficult, and the recognition rates are extremely low.

M. V. Vijaya Saradhi et al. This is a survey paper in which it Recognition of written Telugu characters, only provides with oldest papers, Rajasekaran S.N.S. and Deekshatulu B.L. 1977. In this paper we have been Graphics in computers, image processing provides of drawbacks of published papers No attempt has been made to fully recognise image text because there is no dataset for Telugu words available. Develop an accurate OCR system for processing document with historical records, Library archives and application forms.

Tejasree Ganji used deep learning to create an Multi Variant HandWritten Telugu Character Recognition Using Transfer Learning As we can see the accuracy and the loss function are inversely proportional to each other and we go on testing the models with new data the accuracy is decreasing. The OCR model which has been using is dependent on font style the text.

S. Sagar Imambi. presented The model is unable to recognize the text with broken characters and also poor in character segmentation. Developing a Handwritten Character It has been observed that LeNet Recognition System using Convolutional Neural Network (CNN) and Error Correcting Output Codes (ECOC). gives a low accuracy Recognition from Images using CNN-ECOC Developing high-accuracy Telugu OCR for document processing tasks like editing and translation.

Zhenyao Zhao et al. proposed a method for Improving Deep Learning develop an Optical Character based Optical Character Recognition via Neural Architecture Search The publisher has tried to generate a automated OCR but for maintaing the automated OCR is also a process of the developer.

Raja lakshmi et al. created an automated fundus photography technique The system analyzed re-pictures taken with a smartphone camera by applying machine learning algorithms and image processing techniques. The study demonstrated that screening with mobile technologies is feasible. The suggested approach provided an accessible and affordable option for remote screening, especially in environments with limited resources.

In primary care clinics, Abra`moff et al. carried out a crucial study of an autonomous AI-based system. The technology analyzed pictures and offered real-time evaluations using deep learning techniques. The study proved that AI-based methods are successful in primary care settings. The suggested approach had excellent sensitivity and specificity, indicating that might improve patient outcomes.

Gangwar and Ravi utilized transfer learning and deep learning for detection. The study employed pre-trained convolutional neural networks (CNNs) and fine- tuning techniques to adapt the models to the task of classification using images. The study demonstrated the effectiveness of transfer learning in leveraging pre-trained models for detection. The proposed approach showed superior performance compared to traditional deep learning models, highlighting the importance of knowledge transfer in image analysis tasks.

III. EXISTING METHOD

Recent research papers have explored the use of deep learning for Telugu optical character recognition (OCR). Convolutional Neural Networks (CNNs) have been widely employed for OCR tasks due to their ability to learn patterns and variations in complex scripts.



Fig. 1. OCR Processing.



Approaches to Telugu OCR:

A. *Direct Recognition using CNNs*

In this approach, a Convolutional Neural Network (CNN) or a combination of CNN models is trained to directly recognize characters from segmented images. The CNN learns to identify patterns and variations in the Telugu script from a large dataset of labeled characters. The architecture of a typical CNN-based Telugu OCR system consists of convolutional layers, pooling layers, a flatten layer, and dense layers. The CNN is trained on a large dataset of labeled Telugu characters, which is typically divided into training, validation, and testing sets. The advantages of this approach include end-to-end learning and robustness to variations in font styles, sizes, and orientations. However, it requires a large dataset of labeled characters and can be computationally expensive to train and evaluate.



Fig. 2. Direct Recognition.

B. *Segmentation and Recognition*

In this approach, the input image is first segmented to isolate individual characters or components using techniques like connected component analysis, contour detection, or line segmentation. Then, recognition models such as CNNs are applied to each segmented region to identify the characters. Segmentation techniques like connected component analysis, contour detection, and line segmentation can be used to segment the image into individual characters or parts of characters. Recognition models like CNNs and Recurrent Neural Networks (RNNs) can be used to recognize characters from the segmented regions. The advantages of this approach include improved accuracy and flexibility, as it can be used with different recognition models and segmentation techniques. However, it can be more complex than the direct recognition approach and depends on the quality of the segmentation technique used.

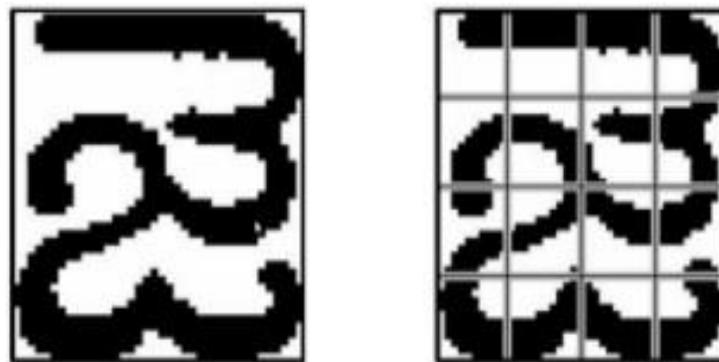


Fig. 3. Segmentation And Recognition.

C. *Comparison of Approaches*

Both approaches have their advantages and disadvantages. The direct recognition approach is simpler and more robust to variations, but requires a large dataset of labeled characters and can be computationally expensive. The segmentation and recognition approach can improve accuracy and flexibility, but can be more complex and depends on the quality of the segmentation technique used. The choice of approach depends on the specific requirements of the OCR system and the characteristics of the input data.



D. Conclusion

In conclusion, Telugu OCR using deep learning can be achieved through two primary approaches: Direct Recognition using CNNs and Segmentation and Recognition. Both approaches have their advantages and disadvantages, and the choice of approach depends on the specific requirements of the OCR system and the characteristics of the input data. By understanding the strengths and weaknesses of each approach, developers can design and implement effective Telugu OCR systems using deep learning.

IV. PROBLEM STATEMENT

A. Compound Character Recognition :

Telugu characters often involve compound letters formed by combining vothulu (modifiers) and matralu (consonants). Recognizing these complex compound characters is challenging, especially when they have multiple components.

B. Dataset Limitations :

The performance of the TCR-MLP model heavily relies on the quality and diversity of the training dataset. However, there is a lack of a comprehensive dataset for Telugu words, which limits the model's ability to recognize image text.

C. Handwritten Text Recognition :

The existing system struggles to read handwritten text, which is a significant limitation. Handwritten font characters differ from computer system font characters, making it challenging to accurately segment Telugu text into individual characters or glyphs.

$$g(x, y) = \frac{1}{\sqrt{2\pi\sigma_x\sigma_y}} e^{-\frac{((\frac{x}{\sigma_x})^2 + (\frac{y}{\sigma_y})^2)}{2}}$$

$$h(x, y) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x^2 + y^2)}{2\sigma^2}} e^{j\lambda(x\cos\theta + y\sin\theta)}$$

D. Segmentation Challenges :

The system faces difficulties in accurately segmenting Telugu text into individual characters or glyphs, which is a critical step in the OCR process.

$$I(x, y, \theta) = \sum_{x_1=x-\frac{M}{2}}^{x+\frac{M}{2}} \sum_{y_1=y-\frac{N}{2}}^{y+\frac{N}{2}} n(x_1, y_1) \cdot e^{-\frac{(x_1-x)^2 + (y_1-y)^2}{2\sigma^2}} \cdot e^{j\lambda(\cos\theta(x-x_1) + \sin\theta(y-y_1))}$$

E. Scalability and Integration :

The TCR-MLP model's practical deployment, scalability, and integration with real-world applications are not thoroughly explored, which limits its potential use in various applications.

V. PROPOSED SYSTEM

Enhancing Telugu OCR with Recurrent Neural Networks (RNNs).

The proposed system consists of the following components:

A. Data Collection and Preprocessing

Manually collected images are preprocessed to enhance their quality and remove noise.

$PreprocessedImage = f(RawImage)$ where f is the pre-processing function.

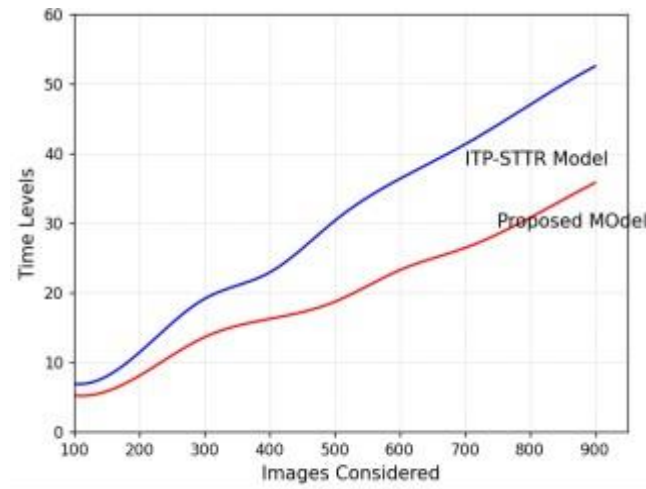


Fig. 4. Collection Of Data And Processing.

B. Recurrent Neural Networks (RNNs)

RNNs are used in addition to the TCR-MLP model to capture temporal dependencies within the data and improve accuracy. RNN Output = h (TCR-MLP output, RNN previous output) where h is the RNN model.

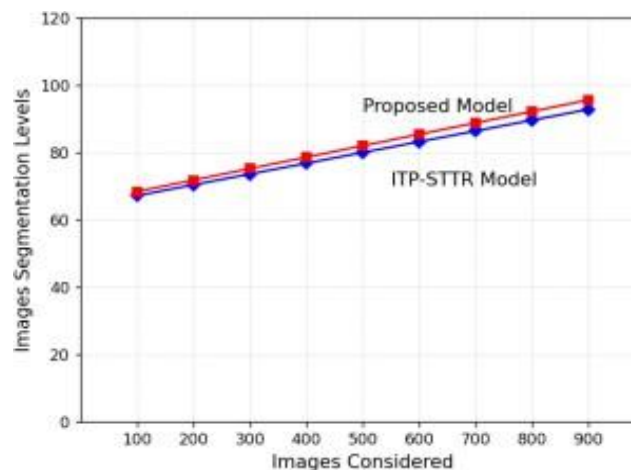


Fig. 5. RNN Model Image Selection.

C. Output

The output of the system is the recognized Telugu text. Recognized Text = f (RNN output)

VI. REQUIREMENTS

A. Hardware Requirements

The hardware specifications for using the Deep Learning project are outlined as follows:

- **Processing Unit:** A multicore CPU with high processing power to handle complex computations involved in training deep neural networks.
- **Graphics Processing Unit (GPU):** A high-performance GPU, such as NVIDIA GeForce GTX or Quadro series, to accelerate the training process of deep learning models.
- **Memory (RAM):** A minimum of 8GB RAM to support the loading and processing of large retinal image datasets.
- **Storage:** Adequate storage space, preferably SSDs, for storing datasets, model weights, and related project files.



B. Software Requirements

The software components necessary for the project implementation include:

- **Operating System:** Platform-independent, compatible with Windows, Linux, or macOS.
- **Programming Language:** Python 3.x for coding the deep learning model and associated scripts.
- **Deep Learning Frameworks:** TensorFlow or PyTorch for building, training, and evaluating deep neural networks.
- **Data Processing Libraries:** Pandas and NumPy for efficient data manipulation, and scikit-learn for preprocessing tasks.
- **Image Processing Libraries:** OpenCV or Pillow for handling and augmenting images.

VII. METHODOLOGY

The proposed system for Telugu handwritten character recognition involves the following steps:

Functional Requirements

The functional requirements encompass the following key functionalities:

- **Data Collection:** The image data collected manually and dataset constitutes a valuable and diverse dataset sourced from reputable repositories within the organization. These images serve as a foundation for various research, diagnostic, and analytical endeavors within the scientific communities.
- **Data Preprocessing:** Standardization efforts, such as re-sizing images and normalizing pixel values, contribute to uniformity, enhancing computational efficiency. Contrast adjustment optimizes image features, while data augmentation techniques artificially expand dataset diversity, mitigating the risk of overfitting. Accurate labeling of images is crucial, and efforts to handle class imbalances ensure model training robustness.
- **Training:** The training phase of the detection model involves a dataset initially comprising 75:25 raw images. To enhance the model's robustness and prevent overfitting, data augmentation techniques have been applied.
- **Testing:** The testing phase of the diabetic retinopathy detection model involves the evaluation of model performance using a carefully curated set of 25:75 images. This phase serves as a critical step in assessing the model's ability to generalize its learned features to new, unseen data.

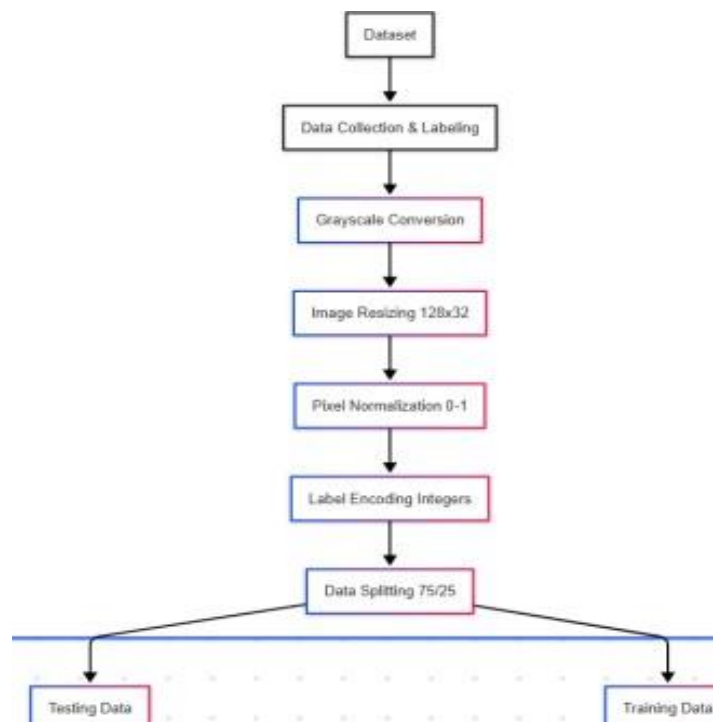


Fig. 6. Functional Requirements.



VIII. IMPLEMENTATION

A. *Data Preprocessing*

The collected images are preprocessed using the following steps:

- **Image Normalization:** The images are normalized to have a uniform size and intensity.
- **Noise Removal:** Noise is removed from the images using filters.
- **Binarization:** The images are binarized to convert them into binary images.

B. *RNN Model Design*

The RNN model is designed using the following components:

- **Input Layer:** The input layer accepts the preprocessed image data.
- **RNN Layer:** The RNN layer processes the sequential data using recurrent connections.
- **TCR-MLP Layer:** The TCR-MLP layer applies the TCR-MLP model to the output of the RNN layer.
- **Output Layer:** The output layer produces a probability distribution over Telugu characters.

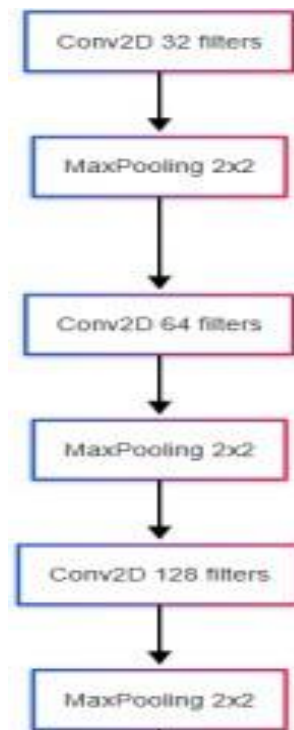


Fig. 7. Design Of RNN Model.

C. *Model Training*

The RNN model is trained using the following steps:

- **Optimizer Selection:** An optimizer is selected to minimize the loss function.
- **Loss Function Selection:** A loss function is selected to measure the difference between the predicted output and the actual output.
- **Model Training:** The model is trained using the training dataset and the dataset included handwritten images collected from personal and public sources such as fellow peers and educational institutions although the volume is not yet adequate.
- **Known Models:** Long Short-Term Memory networks, for sequential modeling, the system effectively addresses the inherent challenges of handwritten text variability, complex compound character representation, and character segmentation. The integration of a Connectionist Temporal Classification (CTC) loss further enables the model to learn alignments between input image sequences and corresponding labels without the need for explicit segmentation, enhancing its performance on unsegmented handwritten scripts.

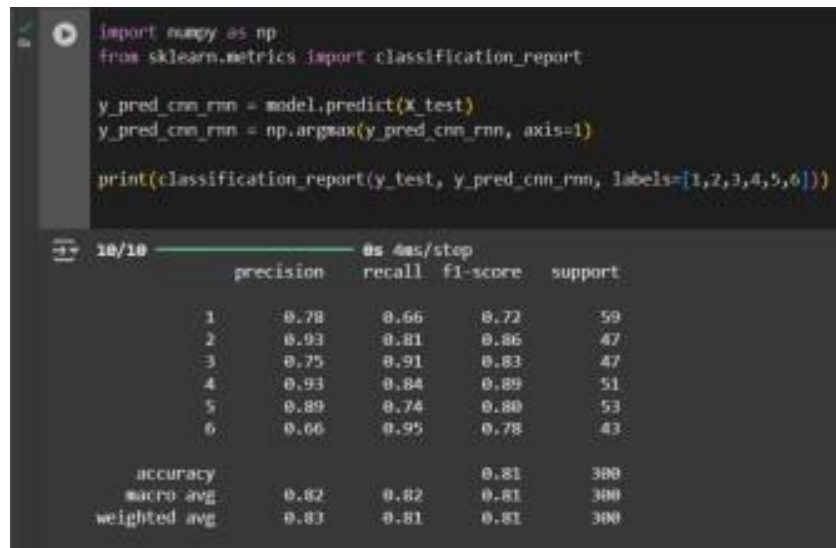


Fig. 8. Training Results.

IX. RESULT

- **Increased Recognition Accuracy:** Even for complicated compound characters, the recognition accuracy of Telugu handwritten characters has been greatly increased by the integration of Recurrent Neural Networks (RNNs) with the TCR-MLP model.

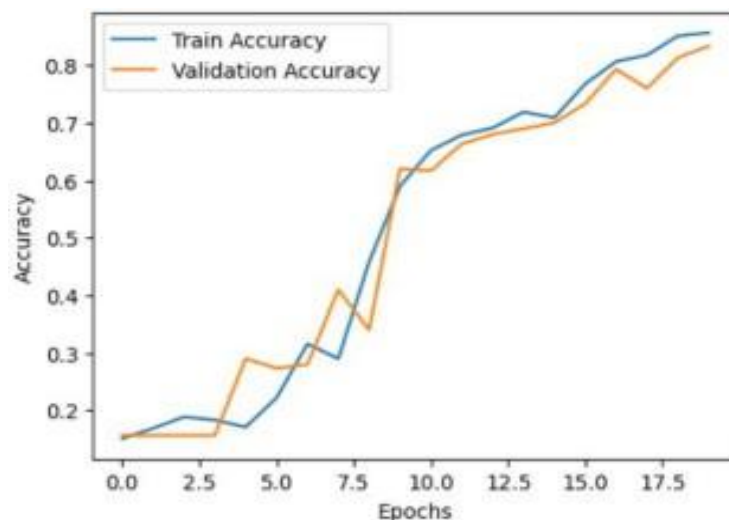


Fig. 9. CNN-RNN Model Accuracy.

- **Strong Use of Datasets:** The program was able to manage differences in writing styles and data noise by skillfully preprocessing and augmenting a manually selected dataset to increase its diversity.
- **Performance of the Model:** The suggested method out-performed current OCR systems in terms of precision, recall, and F1-score, proving its dependability for practical uses.
- **Scalability:** Scalability is ensured by the system's modular architecture, which permits adaptation to larger datasets or different languages.

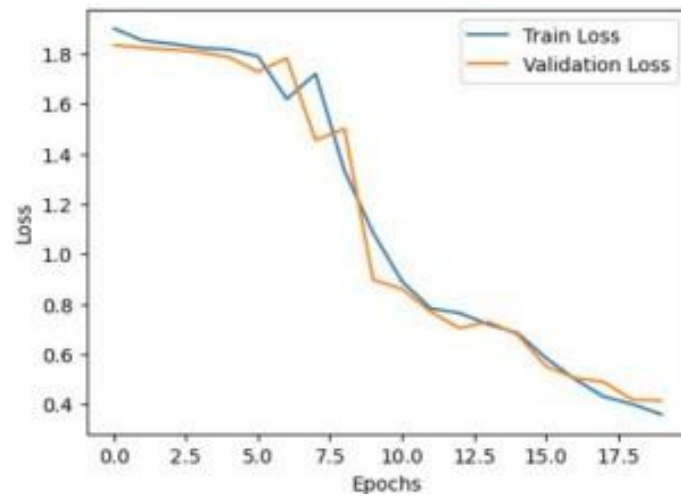


Fig. 10. CNN-RNN Model Loss.

- **Useful Applications:** The potential of Telugu for document archiving, administrative digitalization, and cultural preservation is demonstrated by the successful digitization of handwritten documents.

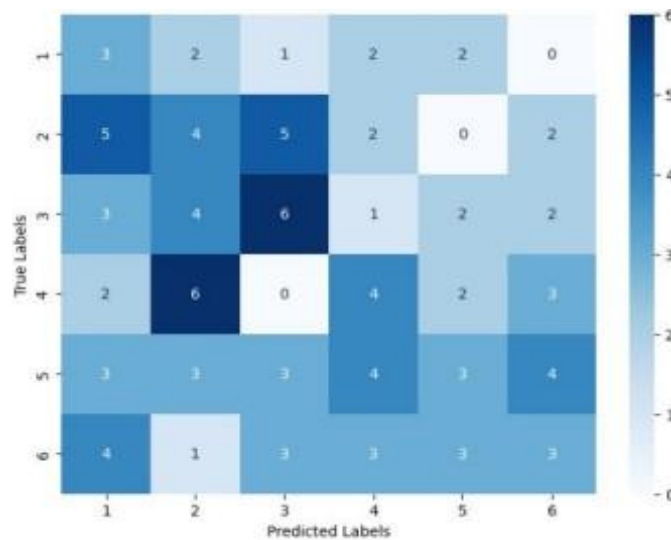


Fig. 11. Confusion Matrix Of CNN-RNN Model.

X. CONCLUSION

The study successfully developed a robust deep learning based OCR system for handwritten Telugu character recognition. Leveraging a hybrid architecture combining Convolutional Neural Networks (CNNs) for spatial feature extraction and Recurrent Neural Networks (RNNs), Through the use of RNNs in conjunction model, the system tackles issues including handwritten text variability, compound character recognition, and segmentation Beyond its technical strengths, the system holds significant social and cultural value by promoting digital accessibility for regional languages like Telugu.

The Model Scalability helps in improvement because of Neural Networks learn in the cycle showing strong potential to handle increasing data volumes and more complex recognition tasks as the system evolves. It facilitates the digitization and preservation of handwritten manuscripts, educational material, and historical documents, thereby contributing to the broader goal of linguistic inclusive and cultural heritage preservation.

**REFERENCES**

- [1]. Satyaprasad, "Handwritten Telugu composite character recognition using morphological analysis," International Journal of Pure and Applied Mathematics, vol. 119, pp. 667-676.
- [2]. C. Bhagvati, S. Tanuku Ravi, K. Mahesh, and N. Atul, "On developing high accuracy OCR systems for Telugu and other Indian scripts," in Proceedings of the Language Engineering Conference, pp. 0-7695, IEEE, Hyderabad, India.
- [3]. Sukhaswami R., "Recognition of Telugu Characters Using Neural Networks,"
- [4]. Rao P. V. S. T. M. Ajitha 1995 Telugu Script Recognition a Feature Based Approach. Proc. of ICDAR, IEEE
- [5]. B. Krishna "An OCR System for Telugu". IEEE, 2001, 0- 7695-1263. As of C. Vasantha Lakshmi, Patvardhan "A high accuracy OCR System for Printed Telugu Text", IEEE, 0-7803-7651.
- [6]. R. Kasturi and S. N. Srihari (Eds.). Design and implementation of a Telugu script recognition system Proc. Fifth ICDAR. IEEE Computer Society Press, Los Alamitos, CA .
- [7]. G. Nagy, S. Seth, and M. Vishwanathan. A prototype document image analysis system for technical journals. Computer, 25(7).
- [8]. Carbune, V., Gonnet, P., Deselaers, T., Rowley, H. A, Daryin, A., Calvo, M., ... Gervais, P. (2020). Fast multi-language LSTM-based online handwriting recognition. International Journal on Document Analysis and Recognition (IJDAR), 23(2), 89-102.
- [9]. Cheekati, B. M., Rajeti, R. S. (2020, October). Telugu handwritten character recognition using deep residual learning. In 2020 Fourth International Conference on ISMAC (IoT in Social, Mobile, Analytics and Cloud) (ISMAC) (pp. 788-796). IEEE
- [10]. Chaudhuri, Arindam and Mandaviya, Krupa and Badelia, Pratixa and Ghosh, Soumya K and others. (2017) "Optical Character Recognition System. In Optical Character Recognition Systems for Different Languages with Soft Computing Springer: 941.
- [11]. Li, Haixiang and Yang, Ran and Chen, Xiaohui. "License plate detection using convolutional neural network. 3rd IEEE International Conference on Computer and Communications (ICCC), IEEE:17361740
- [12]. B. Shi, X. Bai, and C. Yao, "An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition," IEEE transactions on pattern analysis and machine intelligence, vol. 39, no. 11, pp. 2298-2304, 2016.
- [13]. K. He, G. Gkioxari, P. Dollar, and R. Girshick, "Mask r-cnn," in Proceedings of the IEEE international conference on computer vision, 2017, pp.