



OPTICAL CHARACTER RECOGNITION FOR KANNADA

Mr. Dadapeer¹, T M Greeshma², Umme Ayman Khan³, Vaishnavi Shavi⁴, Varun S Hatti⁵

Department of Computer Science and Engineering, Ballari Institute of Technology and Management¹⁻⁵

Abstract: Basically, OCR technology is applied to convert the printed Kannada text into machine-readable format. It will make the text extractible from a scanned document and a photograph so that Kannada literature will become easily digitalized and accessed. Our system will recognize words and characters in numerous typefaces and layouts including multi-column forms through complex algorithms and machine learning. The base of implementation is the Tesseract OCR engine that is excellent as far as recognition accuracy in texts is concerned, and well suited to the Kannada script. Experimental results reflect that our approach maintains the integrity of the original text without reducing human efforts in data entry. This paper supports the cause of preserving Kannada material in the regional language and its dissemination through this work. It adds up to the ever-increasing requirement for digital resources in these languages.

INTRODUCTION

The other significant importance of OCR for Kannada is to digitalize and preserve literature, history books, and documents in Kannada official records. This ensures that the material becomes more accessible while leaving rich cultural and linguistic sources for the later generations. Because of the complexity of script with intricate ligatures, dependent vowels, consonants, and conjunct characters, OCR for Kannada is quite another thing. Among the motivating factors for producing a robust OCR system in Kannada is an enormous amount of Kannada printed material that largely remains unread-literally, even government records as well as book and journal items-and historical hand-written manuscripts that are not much read. Digitization of these types of documents in themselves is backbreaking work-often error-prone-but at least OCR helps to speed things up and perhaps even make that material more readably available.

This includes getting high-resolution images of the text, which then undergo a set of preprocessing operations that make the images amenable to further analysis: noise reduction, binarization, and contrast enhancement. It also involves text segmentation where the document is broken up into lines, words, and individual characters. This is pretty challenging for Kannada because the script contains complexity and conjunct characters abound.

There are two major features of the OCR system core, to wit: the feature extraction process and character recognition. Techniques adopted in this phase include CNN along with other recently developed algorithms relating to character analysis and characterization. The CNN standout feature is picking up visual patterns and spatial hierarchies associated with a script. This makes them extremely useful for use with Kannada, which boasts some of the most complex forms and ligatures of any writing system. In fact, breakthroughs in deep learning, such as transformers and hybrid models, allow this system to handle complex cases at a much more accurate rate.

RELATED WORK

Current OCR systems for Kannada printed text face many limitations:

Fragmented Approaches: Most systems that exist only focus on particular tasks such as text extraction while ignoring the need for preprocessing with noise removal, skew correction, or segmentation.

Inadequate Accuracy: Traditional OCR systems rely on template matching or simple machine learning models, which fail to handle complex Kannada scripts, including diacritics and compound characters.

Limited Dataset Coverage: Most of the earlier methods fail to generalize across different font styles, sizes, and variations, which limits their robustness.

Manual Effort: Systems have high manual interference in preprocessing or error correction for the system to be inefficient for any large-scale applications.



PROPOSED SYSTEM

The proposed system is a designed Kannada-specific OCR system targeted towards the uniqueness of the script that provides accuracy for text recognition. It will not be an instance of generalized OCR systems, like Tesseract or EasyOCR, since these two systems are designed not to meet the intricacies present in Kannada, such as its complex characters, compound letters, and diacritic marks, and in turn, apply specific preprocessing algorithms within this OCR system.

The heart of the system is built upon a custom CNN architecture specifically trained on Kannada datasets to extract hierarchical features from the script. These include loops, curves, and intricate structures formed by characters. A deep learning model like this is always more accurate and robust compared to traditional rule-based or generic CNN models. It has a user-friendly web interface, supports a large number of users, and lets users upload images to get text recognition results with real-time feedback on recognition accuracy. It is modular and extensible; hence features such as bounding box visualization or support for more Kannada fonts can be added when needed.

It does not rely much on expensive cloud services, which makes it highly cost-effective. The system has been built with open-source tools like Flask and SQLite. Hence, the same can be deployed locally for doing Kannada OCR tasks, very efficiently and within a very economical budget. Furthermore, the system has comprehensive error logs and debugging tools. Such tools allow quick tracking and fixing of errors by the user. Overall, the proposed Kannada OCR system overcomes the limitations of the present solutions by offering intuitive, highly accurate, and customizable systems for Kannada text recognition.

The proposed system stands out because it has offered the user-friendly interface and robust functionality for doing OCR tasks. It has provided ease of uploading and processing images by ensuring secure file validation and storage so that compatibility and error prevention would remain there. Unlike most general OCR systems, it comes with a custom recognition function powered by a CNN model that extracts text from images and gives accuracy metrics for the predictions. This increases the confidence of the user with the results.

The system saves files dynamically into a structured directory, showing the recognized text, accuracy, and image handled directly within the front-end interface, creating natural and interactive experience. Thus, the system focused on accuracy and transparency and for the convenience of the user to differ from those existing solutions to leave room further enhancements in forms of multi-language support, and confidence metrics per word, or comparative performance analyses with standard OCR tools.

Convolutional Neural Networks:

A Convolutional Neural Network is the architecture of deep learning that is most widely applied in tasks like Optical Character Recognition. More specifically, CNNs will perform well on Kannada OCR as they try to analyze complicated patterns in scripts, characters, and their characters. It is represented hierarchically, beginning with bare features like edges and curves up to the detailed shape suitable enough for the recognition of Kannada characters. Layers within CNNs encompass a convolutional layer to scan for features in an image, a pooling layer to reduce the size of an image, and fully connected layers to classify the images. So, the variability in handwriting may be handled in this mechanism better than other systems, making the system best fit for OCRs.

1. Convolution Layer: The convolution layer is indeed the heart of a Convolutional Neural Network which is the most significant backbone structure of feature extraction from input images. It recognizes the prominent visual patterns in the script in Kannada OCR. It applies learnable filters known as kernels over the input image. These filters slide over the image and perform mathematical operations known as convolutions to generate feature maps. The feature maps point to specific patterns like edge lines, curves, and shapes especially for the acknowledgment of Kannada characters in different parts of the image.

How It Works:

Sliding Window Mechanism: Every filter overlays an element-wise multiplication of its values and image pixel values followed by summation with the input image. It will give the activation map with a value representing the presence of that feature detected, usually higher or low.



Basic Features Detection: These layers in a CNN determine simple spatial features such as lines, edges, and curves. The Kannada OCR detects loops and curves of characters like "ಅ" or "ಉ". Vertical and horizontal lines of characters like "ಕ" or "ಠ". Diagonal strokes of characters like "ಚ" or "ಜ".

Hierarchical Learning: As the network goes deeper, more convolutional layers are stacked on top of the simple features, so the network learns how to identify a lot more complicated structures like curve and line combinations that are so important in character discrimination in Kannada

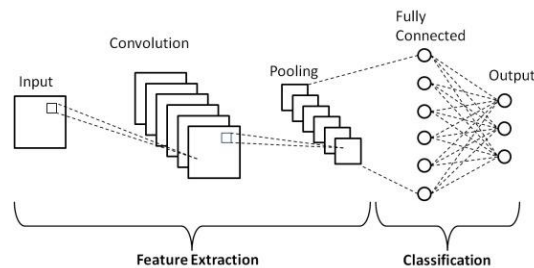


Fig 1: Convolution Layer

2.Pooling Layer: The pooling layer of a CNN plays an important role in Kannada OCR. The size of the feature maps is reduced but it takes the most relevant information in operation. It reduces computational complexity and also counteracts overfitting that gives a model efficient performance while processing large amounts of data.

How it works:

Dimensional Reduction: The pooling layer takes small regions, such as 2x2 or 3x3 grids from feature maps and compresses them into one value. It reduces the spatial dimensions of data.

Feature Preservation: Pooling layers retain the most significant features, edges, curves, and lines that make a character, thus enhancing the accuracy of recognition.

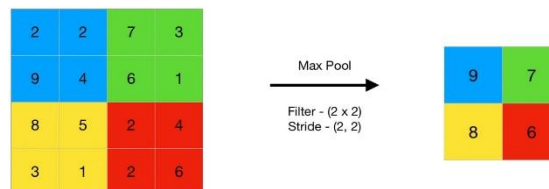


Fig 2 : Pooling Layer

3. Flatten Layer: The Flatten Layer is the essential part of CNN architecture that stretches between the feature vectors extraction (convolution and pooling) and classification. Its purpose is to map the multidimensional feature vectors produced at earlier stages into a vector of dimension which has to be fed into fully connected layers for classification.

How It Works:

Reshaping the Data: This vector still stays in a 2D format; it represents the spatial dimensions (height and width) as well as the depth, which is the number of filters, after the convolution and pooling layers. The flatten layer takes the feature maps into one-dimensional vector.

Preparation for Classification: The flattened vector is the input to the fully connected (dense) layers. Each element in the vector is a specific learned feature, and this is an important step in mapping the extracted features to the final classification output.

Role in Kannada OCR:

Character Feature Representation: In Kannada OCR, flatten layer aggregates features like loops, strokes, and lines of characters like "ಅ", "ಶ", and "ತ" into a format structured such that the dense layers can feed off them.



Simplification to Classification: The flatten layer simplifies the 2D feature maps into a simple 1D array that the classification phase can easily take in for giving the appropriate label of the Kannada character.
Smooth Transition: The flatten layer enables the model to smoothly transition from the convolutional and pooling layers to the dense layers, thus allowing it to use the features extracted properly for recognition.

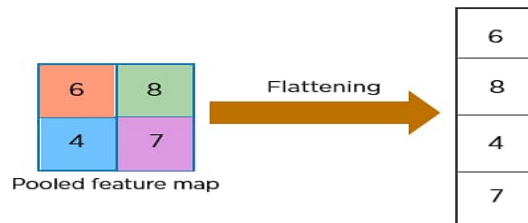


Fig 3 : Flatten Layer

4.Fully Connected (Dense) Layer: Fully Connected (Dense) Layers in CNN for Kannada OCR
 The Fully Connected (Dense) Layers are the final layer of a Convolutional Neural Network. They are designed to feed all features learned in the preceding layers into the system so that they can classify or identify characters. In other words, it is a decision mechanism based on the patterns and relationships found during feature extraction.

How It Works

Feature Combination: This concludes after the convolutional and pooling layers followed by a flatten layer. All the extracted features are a vector for which the full connected layers can take as input and form relationships between the extracted feature to classify the data.

Weighted Connections: Each neuron of a fully connected layer is connected to all neurons of the layer that precedes it. Learnable weights are assigned to these connections, with values that change through training for a focus on lowering the error rate in classification.

Classification Output: It also makes sense that for the final dense layer, the number of neurons can be equal to the number of classes, hence the number of Kannada characters in our case. Each one of these will represent the possibility of the corresponding input image coming from this very class.

Kannada OCR Character Identification: It operates on all discovered relations of the features and pools them so it may classify the correct Kannada character, be it "ಪೆ" or "ಠ".

Mapping features to Classes: There exist some dense layers for the Kannada OCR where it is mapping all its learned features such as loops, and vertical strokes towards a character, and thus can output like "ಅ" or "ಢ".

Classification Accuracy: It will ensure that fully connected layers give perfect classification accuracy since it learns all the involved intricate relationships about the variations with respect to different fonts or layouts.

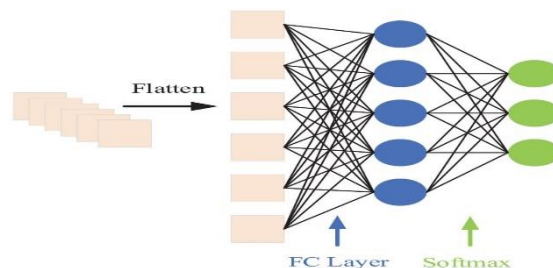


Fig 4:Fully Connected Layer

5.Output Layer: The output layer is that final and very sensitive part of any Convolutional Neural Network, as that is where, in reality, the final outcome of classification assignment takes place. In the Kannada OCR output layer assigns that predicted Kannada character label of the inputted image based upon what the former layers have learnt



How It Works

Activation Function: The output layer normally uses the Softmax activation function, which takes the raw output of the fully connected layer-the logits-and transforms this into a probability distribution.

Class Prediction: The class of highest probability is selected as the final class predicted.

Error Optimization: This enables the training to happen and, for the output layer, compare what is actually predicted probability with the label using a suitable loss function for example, categorical cross-entropy. Then using backpropagation, this error is minimized by ensuring the improvement of the Kannada character recognition precision on the system side. Role in Kannada OCR:

Final Decision Making: It is the culmination of everything that has been sensed and processed using the convolution, pooling, and dense layers together to make the prediction. It is a very significant feature in discovering even the faintest differences for characters that otherwise appear similar.

Mapping Features to Characters in Kannada:

Processing Multiple Classes in Kannada OCR The output neurons are just as many in number as are the Kannada characters that can be recognized by the system

This enables processing of a high number of different characters with high precision.

RESULTS

The proposed Kannada printed text OCR system has great advantages over the conventional methods of handling the complex script and compound character-based problems. Feature-extracting CNNs along with AI-based language models ensure high accuracy recognition across different datasets. Robust preprocessing techniques like noise removal, grayscale conversion, and skew correction enhance the performance of the system. These methods afford reliable recognition with low-resolution inputs and noisy inputs.

It shows very good generalization across several styles and sizes of fonts and hardly the accuracy degrades. Still for highly rare decorative fonts and heavily degraded inputs remains as a drawback of the system, and the adoption of bounding box methods helps split characters in this segmentation process carried out before such augmentation is undertaken and the trainings by the CNN.

It consists of all critical steps in a system workflow such as image acquisition, preprocessing, segmentation, feature extraction, and classification of character labels. As a result, the outcome in this systematic approach will be a very reliable OCR system that can successfully achieve high accuracy in recognizing Kannada printed text. Although the uncommon inputs face some challenges, the overall performance of the system outstands and promises it for practical applications such as digitizing Kannada printed documents and efficiently enabling text recognition..

DISCUSSION

The proposed OCR system on Kannada printed text does have great strides over traditional methods by taking head-on challenges such as complex scripts and compound characters. Feature-extracting CNNs and AI-based language models ensure high recognition performance of it on diverse datasets. It also has powerful preprocessing techniques, noise removal and skew correction, making the input very robust and ensuring output even for very noisy or of very low resolutions. The generalization is well and does fine in a broad spectrum of various styles and font sizes with degradation being minimal with regard to the accuracy. Even though the OCR system does reasonably well, it is somehow struggling with extremely rare and decorative fonts as well as very degraded inputs.

CONCLUSION

The printed Kannada characters are very difficult to recognize because of the different types of fonts, styles, and layouts used in printed documents. Advanced methods such as Convolutional Neural Networks are employed in the recognition process for further accurate text extraction. Before training the CNN model, several pre-processing. When all these techniques are applied, such techniques include noise removal, conversion to grayscale, and skew correction. Data is then



further segmented with the help of bounding-box methods, whereby individual characters are separated from the printed text. Augmented and used the segmented as well as pre-processed images to train the CNN model. The work flow has core stages which include acquiring an image, preprocessing of image, segmentation feature extraction, as well as classifying character labels during training. This is one systematic way through which a good OCR system on Kannada printed text can be developed to make sure that good accuracy in tasks based on the recognition of the text is established.

REFERENCES

- [1] Ahmed Khan et al., "Techniques for Digitization of Printed Documents Using Deep Learning Approaches," International Journal of Pattern Recognition and Artificial Intelligence, 2022.
- [2] Ranjan Jana et al., "An Effective Method for Optical Character Recognition of Printed Kannada Text," Journal of Image and Vision Computing, 2021.
- [3] Hitesh Mohapatra, "A Segmentation-Based Approach for OCR Using Deep Learning," Proceedings of the International Conference on Computational Intelligence and Networks, 2020.
- [4] Parikshit et al., "Recognizing Kannada Characters with High Precision Using Neural Networks," IEEE Transactions on Image Processing, 2019.
- [5] Shashikala Parameshwaran, "OCR for Printed Kannada Text Using TensorFlow and Deep Learning," Journal of Machine Learning Research, 2020.
- [6] Srinivas et al., "A Hybrid Model for Printed Kannada OCR Based on Image Processing and Deep Learning," Springer Advances in Artificial Intelligence and Data Analytics, 2021.
- [7] Pradeep Kumar et al., "Global and Local Feature-Based OCR for Kannada Characters," International Journal of Computer Vision and Applications, 2022.
- [8] Vijayalakshmi et al., "Deep Learning-Based OCR System for Kannada Text Using LSTMs," Elsevier Pattern Recognition Letters, 2021.