

Handwritten Tamil Character Recognition in Palm Leaf Manuscripts using BiLSTM Classifier

Dr. M. Mohamed Sathik¹, R. Spurgen Ratheash²

¹Principal and Head, Research Department of Computer Science,
Sadakathullah Appa College, Tirunelveli, Tamil Nadu, India.
Affiliated to Manonmaniam Sundaranar University, Abishekapatti, Tirunelveli-627012.
TamilNadu, India.

²Research Scholar, Reg.No.12334, Sadakathullah Appa College, Tirunelveli, Tamil Nadu, India
Affiliated to Manonmaniam Sundaranar University, Abishekapatti, Tirunelveli-627012.
TamilNadu, India.

E-mail :¹mmdsadiq@gmail.com, ²spurgen@gmail.com

Abstract

Recognizing and reading the Tamil characters which were written many centuries ago in the palm leaf manuscripts is a tough task. As the character has taken a different form over the centuries, the shape of the characters has not been known to contemporary Tamil readers. Neural network classifiers are the tool to unlock the treasure of knowledge paralyzed for such reasons. Only those who are familiar with the shape and strokes of Tamil characters can read palm leaf manuscripts. The knowledge gained from the Tamil literatures have learned about the palm leaf manuscripts has been prepared through computer so that all the people can know the written ideas. This method works in three stages like text line segmentation, character segmentation, character recognition. The final stage of recognition has done by the Bidirectional Long Short Term Memory (BiLSTM) classifier that produces a better result than other conventional CNN methods in Tamil character recognition.

Keywords: RNN, BiLSTM, LSTM, CNN, palm leaf manuscripts, Tamil character recognition, softmax layer, fully connected layer

1. Introduction

Initially, ancient people have written the Tamil scripts on the pots and then started writing on palm leaves. The palm leaf manuscripts existed as a medium to write and record incidents, events, innovations in medicine, astrology and literatures. The collection of palm leaf manuscripts is known as *suvaligal* which contains approximately 40 to 50 preserved leaves with 15 to 30 cm length and 3 to 12 cm width written on both the sides. Each side of the palm leaf has 5 to 6 written lines with sequence of characters from left to right. When the leaves were in dilapidated condition, the scribes copied the contents of one leaf to another new leaf [1]. While copying, an inherited knowledge about shapes and writing style of locality influences to define the shapes of the character. Tamil language has 247 independent characters with the combination of 12 vowels and 18 consonants and one special character ‘ / ’ as in *Figure 1*. Most of the recent Tamil characters were not written in palm leaf manuscripts and the written characters are also different in forms. As a result, very much vital information in siddha medicine, astrology, and literatures could not be recognized [2] because of the shapes of old characters have been forgotten by the new shapes. Even though digitization is the key solution to preserving and storing the

palm leaf manuscripts, there has been no progress in publishing them as books as the characters are unrecognizable. This article discusses the highly productive way of Tamil character recognition by matching the early day character shapes with present-day

| | | VOWELS | | | | | | | | | | | |
|------------|----------|--------|----|----|----|----|----|----|----|----|----|----|----|
| | | அ | ஆ | இ | ஈ | உ | ஊ | ஓ | ஔ | எ | ஏ | ஐ | ஔ |
| CONSONANTS | VALLINAM | க | கா | கி | கீ | கு | கூ | கை | கே | கொ | கோ | கோ | கோ |
| | | ச | சா | சி | சீ | சு | சூ | சை | சே | சொ | சொ | சொ | சொ |
| | | ட | டா | டி | டீ | டு | டூ | டை | டே | டொ | டொ | டொ | டொ |
| | | த | தா | தி | தீ | து | தூ | தெ | தே | தொ | தொ | தொ | தொ |
| | | ப | பா | பி | பீ | பு | பூ | பை | பே | பொ | பொ | பொ | பொ |
| | MELLINAM | ங | ஙா | ஙி | ஙீ | ஙு | ஙூ | ஙை | ஙே | ஙொ | ஙொ | ஙொ | ஙொ |
| | | ஞ | ஞா | ஞி | ஞீ | ஞு | ஞூ | ஞை | ஞே | ஞொ | ஞொ | ஞொ | ஞொ |
| | | ண | ணா | ணி | ணீ | ணு | ணூ | ணை | ணே | ணொ | ணொ | ணொ | ணொ |
| | | ந | நா | நி | நீ | நு | நூ | நை | நே | நொ | நொ | நொ | நொ |
| | | ம | மா | மி | மீ | மு | மூ | மை | மே | மொ | மொ | மொ | மொ |
| | IDAYINAM | ய | யா | யி | யீ | யு | யூ | யை | யே | யொ | யொ | யொ | யொ |
| | | ர் | ரா | ரி | ரீ | ரு | ரூ | ரை | ரே | ரொ | ரொ | ரொ | ரொ |
| | | ல் | லா | லி | லீ | லு | லூ | லை | லே | லொ | லொ | லொ | லொ |
| | | வ | வா | வி | வீ | வு | வூ | வை | வே | வொ | வொ | வொ | வொ |
| | | ள் | ளா | ளி | ளீ | ளு | ளூ | ளை | ளே | ளொ | ளொ | ளொ | ளொ |

Figure 1. Tamil Characters

characters. The result produces better performance on Tamil character recognition in palm leaf manuscripts than other conventional methods.

The remaining part of the article has been arranged as follows. Section 2 discusses the related work in Tamil character recognition. Section 3 explores the general BiLSTM layer framework. Section 4 discusses the process of Tamil character recognition. Section 5 describes dataset and training as an experimental setup. Section 6 presents the result and discussions. The last section gives the conclusion of the paper.

2. Related work

The character recognition process offered in many languages such as Thai, Khmer, Arabic, English, and Tamil are in different techniques. The Tamil character recognition has been taken into consideration for the analysis. The Kohonen Self-Organizing Map (SOM) tuned by global feature technique in the type of Artificial Neural Network used to classify handwritten Tamil characters [3]. The symbols, numerals and Tamil characters are recognized by the techniques of Gabor Filter and Support Vector Machines (SVM) [4]. Hilditch's Algorithm is used in Neural Network to recognize typed Tamil characters by passing the Horizontal histogram, Vertical histogram, radial, input and output features in minimum number of classes [5]. The features of character height, width, number of vertical and horizontal lines, curves, circles, slope lines, dots are extracted and processed by SVM and Kohonen SOM in Artificial Neural Network to recognize offline handwritten Tamil characters for eight classes[6]. A method to extract the feature of the characters is Hu's invariant and Zernike movements and classifies the characters using Feed Forward Neural Network [7]. A survey in Tamil character recognition explained deep belief network method to extract the features, Restricted Boltzmann Machines model to train the character using deep learning in large data [8]. An ideal edge identification method used in palm leaf Tamil characters using Canny Edge Detection by three ways such as great discovery, great confinement, and negligible reaction in Artificial Neural Network. Finally, their method enhances the character with the binarization technique to provide a remarkable result in recognition [9]. A survey provides the method about character recognition in palm leaf manuscripts of Southeast Asia languages like Balinese, Khmer, Sundanese using CNN [10]. CNN is used to recognize the character of Tamil palm leaf manuscripts by five layers such as convolution, pooling, activation, fully connected layers and softmax classifiers [11]. Nine layers including five convolution layers and each

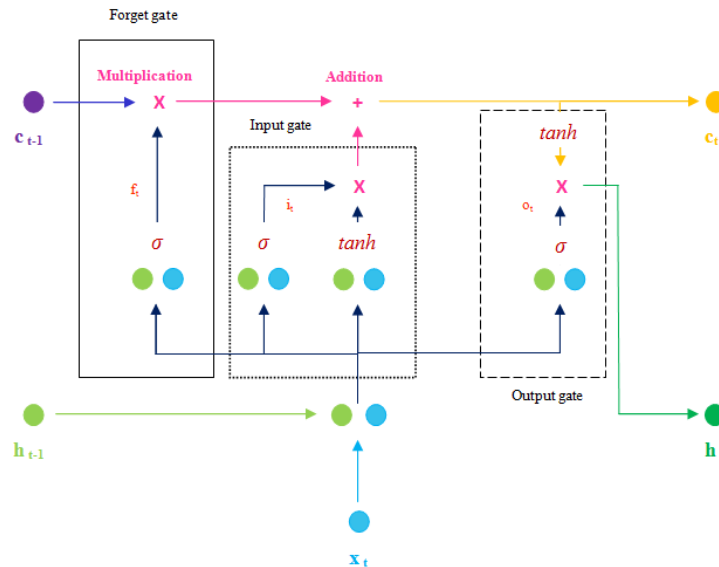


Figure 3. LSTM block

functions are applied on input such as sigmoid and tanh. The first decides which input value gets 0 or 1. The tanh adds weight to the value from input ranges from -1 to 1. The results of both functions are multiplied and update the previous cell state which has multiplied by the forget gate. The third gate is output gates that apply sigmoid function on the combined value of hidden state and input values to know which gets 0 or 1. The tanh function applied on the cell state updated by previous gates. The multiplied results of both produce an input to the next LSTM block [15].

3.2 Recurrent Neural Network (RNN)

In RNN, the image vectors are taken as input for sequence input layer to extract the features and passed to the Bidirectional Long Short Term Memory (BiLSTM) layer. The sum and multiplying options are applied on features in forward and backward LSTM cells. The output of LSTM blocks passed to multiply the weight matrix and to add the bias vector by fully connected layer [16]. The activation function is added to the weighted features in softmax layer followed by classification layer used to compute loss and finally the output produced.

3.2.1 Sequence Input Layer: Sequence Input Layer is the first layer in PLTCR architecture. The two dimensional vector sequence input image considers input size as a scalar for the count of features. The vector has three elements such as height (h) and width (w), and number of channels (c) of an image.

3.2.2 BiLSTM Layer: BiLSTM is a wrap that conjoins two parallel LSTM layer. One of the two with input processed forward and other one with output processed backwards. Merger mode of the bidirectional layer combined the forward and backward outputs and passed to the subsequent layer. The merge being with the options of sum, multiplication to add and multiply the output together, concatenation and average are used to produce output for the next layer, the default option is concatenation. The LSTM is also known as memory blocks when they are connected recurrently. The memory blocks have three multiplicative units such as input gate (i_t), output gate (o_t) and forget gate (f_t). The memory cells update by hidden layer content (h_{t-1}), input (x) with the current time step (t) and add bias (b) value [17]. The sigmoid (σ) function makes the decision to retain the values in each gates as in the following mathematical representations.

$$i_t = \sigma(C_i[h_{t-1}, X_t] + b_i) \quad (1)$$

$$f_t = \sigma(C_f[h_{t-1}, X_t] + b_f) \quad (2)$$

$$o_t = \sigma(C_o[h_{t-1}, X_t] + b_o) \quad (3)$$

The tanh function supports to distribute the gradient longer by vector cell state (c_t) to memory cell to evade the vanishing or exploding gradient problem. The tanh function adds weight to the input with bias value and updates the previous cell state as in relationship 4.

$$c_t = \tanh(C_c[h_{t-1}, X_t] + b_c) \quad (4)$$

The benefit of sequence modelling to access both the past and future contents in BiLSTM can be achieved by forward and backward LSTM layers (Alex Graves et al., 2005). In character recognition, the output of two LSTM blocks for the character (h_i) is sum of features in forward and backward block output as in the following representation.

$$h_i = [\vec{h}_t \oplus \overleftarrow{h}_t] \quad (5)$$

3.2.3 Fully Connected Layer: The weight and bias add to all neurons in this layer produced by the previous BiLSTM Layer. The patterns can be identified by combining all the features that learned from the succeeding layer. The process is used to classify the image. The layer is independent at each time step when the sequence input [18]. The bias adds with the weight of an input (x) with current time stamp (t).

3.2.4 Softmax Layer: The softmax is activation function to compute the probability distribution for the list of classes (x_i) in the range between 0 and 1 with the sum of probability is equivalent to 1. The softmax can be calculated by the following representation [19].

$$S(x_i) = \frac{\exp^{x_i}}{\sum_j \exp^{x_j}} \quad (6)$$

The layer is unlike sigmoid rather performs multi class classification task. In loads of architecture, the layer more or less exists at the end so it also known as output layer of deep learning architecture.

3.2.5 Classification Layer: The layer calculates the cross entropy by assign mutually exclusive classes to each input values taken from the softmax function by the following mathematical relation. The number of samples (s), number of classes (c), an output (o), and an indicator (t) of i^{th} sample fit in with j^{th} class as in the following representation [20].

$$C = -\frac{1}{n} \sum_{i=1}^s \sum_{j=1}^c t_{ij} \ln o_{ij} \quad (7)$$

4. Proposed method

The Tamil character recognition in palm leaf manuscript using RNN has two phases and each has two processes such as normalization and labelling for first phase, training and testing for second phase as in *Figure 4*. In the first phase, the character segmentation by HorVer method produces different size of images. Each image must be normalized in equal aspect ratio 30 x 30 by the centre point of the strokes. The measurement of above and below pixels from the centre decides to acquire the complete shape of the character. In the second phase, the normalized images are labelled individually to train the character. The background of an input image is black in colour and foreground character is in white colour with the value of 0 and 1 respectively.

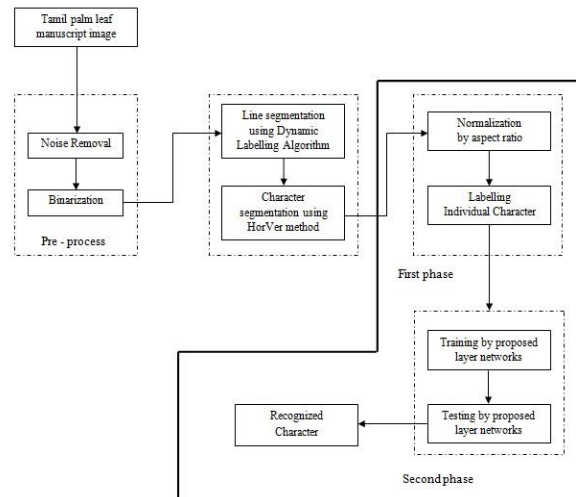


Figure 4. Process of Tamil Character Recognition

4.1 Architecture

The Tamil character recognition in palm leaf manuscripts have the hierarchical order of sequence input layer, BiLSTM layer, Fully connected layer, Softmax layer and Classification layer with suitable functionalities to implement RNN as in *Figure 5*.

4.2 BiLSTM Layer in Tamil Character Recognition

The RNN model has created as specified in the previous section for Tamil character recognition in palm leaf manuscripts. BiLSTM layer update the features with ‘sigmoid’ gate activation function, ‘tanh’ state activation function, and 200 x 1 hidden states. 100 epochs of training cycle and maximum 97 iterations per epoch with the learning rate of 0.001. The third layer, fully connected layer update the weights by total number of characters x 200. The loss calculated by crossentropyex loss function in softmax layer as fourth layer end with the allocation of labels for classes in fifth layer as classification layer. The second phase of testing, a new data set contains 2,640 Tamil characters in palm leaf manuscripts used to recognize 88 characters which was not trained. In training the layer used several hyper parameters to recognize Tamil characters in palm leaf manuscripts as in *Table 1*.

Table 1. Hyper Parameters for Tamil Character Recognition

| Hyper Parameters | Values |
|------------------|--------|
| Initialization | Glorot |
| Batch Size | 27 |
| Interpreter | Adam |
| Epochs | 100 |
| Learning rate | 0.001 |

5. Experimental setup

5.1 Dataset

The palm leaf Tamil character recognition has a unique dataset created by the researchers. The different styles of vowel and consonant character images are in white strokes in black background. The dataset has IWFHR2010Tamil vowel characters collected from HP Labs India available at free of cost. The consonant character images are handwritten characters collected from 270 members and digitalized by the scanner.

Table 2. Dataset of Palm Leaf Tamil Characters

| Contents | Counts |
|-------------------------------|----------|
| Total no of palm leaves | 950 |
| No of text lines per leaf | 5 |
| Total no of lines | 4750 |
| No of characters in each line | 50 or 45 |
| Total no characters | 213750 |

The scanned image has been segmented by the researchers and classified as single. The detail of the dataset has shown in *Table 2*. In the collection of data, 10457 individual character mages have been selected for training and 2640 for testing in neural network.

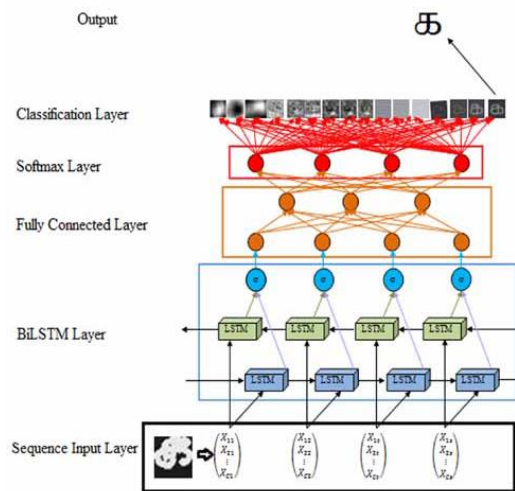


Figure 5. Layer Architecture for Tamil Character Recognition

5.2 Training and Testing

Different batch size were trailed and fixed the batch size as 27. After the layer architecture is compiled, the specific processed character dataset is loaded to train the model. The training is done by 100 epochs as in *Figure 6*. The tuning is made in hidden states and activation function to achieve an optimum in accuracy.

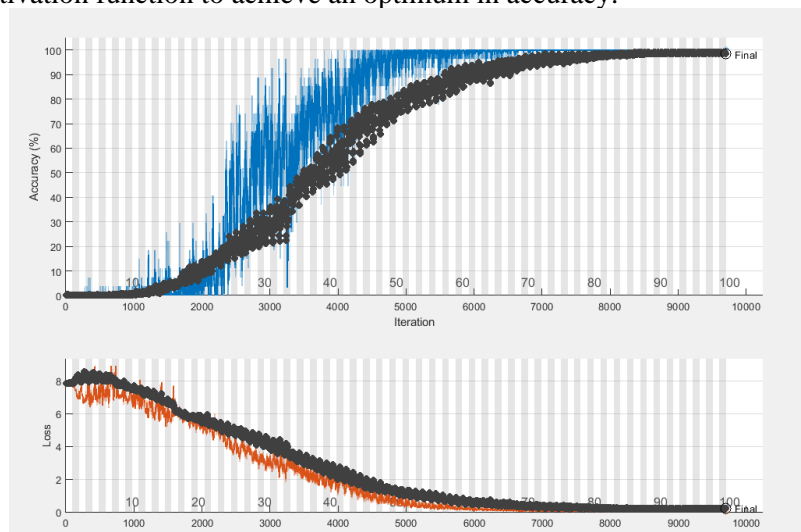


Figure 6. BiLSTM Training and Validation

6. Results and discussions

The RNN layer architecture recognizes the character written by different scribes. This work provides 1 % of wrong prediction rate in testing with the Tamil character in palm leaf manuscripts. The dataset is trained separately for CNN and LSTM classifiers and compared the Recognition Accuracy with proposed layer architecture using BiLSTM classifier as in *Figure 7*. The BiLSTM classifier has the benefit to defeat over-fitting

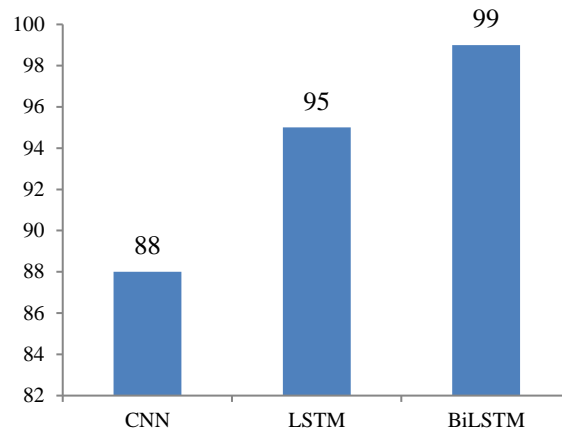


Figure 7. Performance of Palm Leaf Tamil Character Recognition

problem without using dropout techniques in CNN. The LSTM does not have backward propagation that reduces the performance in palm leaf Tamil character recognition. The recognition accuracy of proposed work is compared with other previous works as in *Table 3*.

Table 3. Comparison of Tamil Character Recognition

| Years | Dataset | Language | Document | No. of classes | Method | Training Accuracy | Recognition Accuracy |
|----------|-----------|----------|-------------|----------------|-------------|-------------------|----------------------|
| 2016 | HP Labs | Tamil | Handwritten | 35 | CNN | 99% | 94.40% |
| 2018 | HP Labs | Tamil | Handwritten | 146 | CNN | - | 88.86% |
| 2019 | Own | Tamil | Palm Leaf | 60 | CNN | - | 96.21% |
| 2019 | HP Labs | Tamil | Handwritten | 156 | CNN | 95.16% | 97.70% |
| Proposed | Palm leaf | Tamil | Palm Leaf | 88 | BiLSTM -RNN | 92.31% | 99.57% |

The CNN classifier combined with principal component analysis and trained 50 epochs to get maximum accuracy [21]. The over-fitting problem between training and validation and different kind of regularization methods were applied to solve that produced 89.3% of test accuracy initially in CNN. Stochastic pooling, probabilistic weighting and dropout techniques were used to get marginal changes in result [22]. The single scribe character set 3.79% of erroneous prediction and 0.64sec time has taken to predict one character [11]. The dropout regularization technique were used in every convolution layer with an initial probability 0.1 and increased by the same to overcome the over-fitting [12].



Figure 8. Two Characters in Single Image

7. Conclusion

The research work used RNN classifier to recognize Tamil characters in palm leaf manuscripts. The work has recognized all characters used in Tamil palm leaf manuscripts long before 300 to 400 years. The RNN classifier also predicts exact character when two characters have joined together in a single image *Figure 8*. The recognition accuracy in this work provides much better result than the previous conventional handwritten Tamil character recognition methods. This work is benchmarking for Tamil character recognition in palm leaf manuscripts that will extend to recognize the stone epigraphs in future.

References

- [1] D. Udaya Kumar, G.V Sreekumar and U.A Athvankar, "Traditional Writing System in Southern India Palm Leaf Manuscripts", Design thoughts, IDCITB, (2009), pp. 1–7.
- [2] S. Rajkumar, N. Srinivasan, T. Thirunarayanan and R. Sangeetha, "Survey of Tamil Siddha manuscripts in possession of Traditional Healers in Northern Tamil Nadu", International Journal of Pharmacology and Clinical Sciences, vol. 1, no. 3, (2012), pp. 68-73.
- [3] K. Sarveswaran and D. Ratnaweera, "An Adaptive Technique for Handwritten Tamil Character Recognition", Proceedings of International Conference on Intelligent and Advanced Systems, Kuala Lumpur, Malaysia, (2007) November 25 – 28.
- [4] P. Soman, R. Ramanathan, S. Ponmathavan, N. Valliappan, L. Thaneshwaran and S. Arun.S.Nair, "Optical Character Recognition for English and Tamil Using Support Vector Machines", Proceedings of International Conference on Advances in Computing, Control, and Telecommunication Technologies, Trivandrum, India, (2009) December 28 -29.
- [5] V. Karthikeyan, "Hilditch's Algorithm Based Tamil Character Recognition", International Journal of Computer Science & Engineering Technology, vol. 4, no. 3, (2013), pp. 268 - 273.
- [6] P. Banumathi and M. Nasira, "Handwritten Tamil Character Recognition using Artificial Neural Networks", Proceedings of International Conference on Process Automation, Control and Computing, Tamilnadu, India, (2011) November 20 - 22.
- [7] Amitabh Wahi, S. Sundaramurthy and P. Poovizhi, "Handwritten Tamil Character Recognition" Proceedings of Fifth International Conference on Advanced Computing, Chennai, India, (2013) December 18 - 20.
- [8] R. Jagadeesh Kannan and S. Subramanian, "An Adaptive Approach of Tamil Character Recognition Using Deep Learning with Big Data-A Survey", Proceedings of the 49th Annual Convention of the Computer Society of India, Telangana, India, (2014) December 12 - 14.
- [9] P. Selvakumar and S. Hari Ganesh, "Tamil Character Recognition using Canny Edge Detection Algorithm", Proceedings of World Congress on Computing and Communication Technologies, Tamil Nadu, India, (2017) February 2 – 4.
- [10] MadeWinduAntaraKesiman, Dona Vally, Jean-Christophe Burie, Erick Paulus, Mira Suryani, SetiawanHadi, Michel Verleysen, SopheaChhun and Jean-Marc Ogier, "Benchmarking of Document Image Analysis Tasks for Palm Leaf Manuscripts from Southeast Asia", J. Imaging, vol. 4, no. 2, (2018), pp 1- 27.
- [11] R. S. Sabeenian, M. E. Paramasivam and P. M. Dinesh, "Palm-Leaf Manuscript Character Recognition and Classification Using Convolutional Neural Networks", Computing and Network Sustainability, Springer, (2019), pp. 397-404.

- [12] B. R. Kavitha and C. Srimathi, "Benchmarking on offline Handwritten Tamil Character Recognition using Convolutional Neural Networks", *Journal of King Saud University – Computer and Information Sciences*, (2019), pp. 1–8.
- [13] R. Spurgen Ratheash and M. Mohamed Sathik, "Line Segmentation Challenges in Tamil Language Palm Leaf Manuscripts", *International Journal of Innovative Technology and Exploring Engineering*, Vol. 9, No.1, (2019), pp 2363 - 2367.
- [14] M. Mohamed Sathik and R. Spurgen Ratheash, "Optimal Character Segmentation for Touching Characters in Tamil Language Palm Leaf Manuscripts using Horver Method", *International Journal of Innovative Technology and Exploring Engineering*, Vol.9, no. 6, (2020), pp 1010 - 1015.
- [15] SeppHochreiter and JurgenSchmidhuber, "Long Short-Term Memory", *Neural Computation*, Vol. 9, No. 8, (1997), pp 1735 – 1780.
- [16] Alex Graves, Santiago Fernandez and Jurgen Schmidhuber, "Bidirectional LSTM Networks for Improved Phoneme Classification and Recognition", *Proceedings of International Conference on Artificial Neural Networks*, Hamburg, Germany, (2014) September 15-19.
- [17] Rania Maalej and Monji Kherallah, "Convolutional Neural Network and BLSTM for offline Arabic handwriting recognition", *Proceedings of International Arab Conference on Information Technology*, Lebanon, (2018) November 28-30.
- [18] J. Wu, "Compression of fully-connected layer in neural network by Kronecker product," *Proceedings of Eighth International Conference on Advanced Computational Intelligence*, Chiang Mai, (2016) February 14 - 16.
- [19] R. Hu, B. Tian, S. Yin and S. Wei, "Efficient Hardware Architecture of Softmax Layer in Deep Neural Network," *Proceedings of International Conference on Digital Signal Processing*, Shanghai, China, (2018) November 19 - 21.
- [20] S. Albawi, T. A. Mohammed and S. Al-Zawi, "Understanding of a convolutional neural network,"*Proceedings of International Conference on Engineering and Technology*, Antalya, Turkey,(2017) August 21 - 23.
- [21] Prashanth Vijayaraghavan and Misha Sra, "Handwritten Tamil Recognition using a Convolutional Neural Network", *Massachusetts Institute of Technology*, (2016), pp 1- 7.
- [22] M. Sornam and C. Vishnu Priya, "Deep Convolutional Neural Network for Handwritten Tamil Character Recognition Using Principal Component Analysis", *Springer Nature*, (2018), pp. 778–787.

Authors



Dr. M. Mohamed Sathik is the Principal of Sadakathullah Appa College, Tirunelveli, India. He received two Ph.Ds majored in Computer Science and Computer Science & Information Technology in Manonmaniam Sundaranar University, Tirunelveli, India. He has many more feathers in his cap by degrees such as M. Tech, MS (Psychology) and MBA. He is pursuing post Doctoral degree in Computer Science. Known for his active involvement in various academic activities, he has attended many national and international seminars, conferences and presented numerous research papers. With publications in many international journals, he has published two books besides having guided more than 40 research scholars. The prolific academician is a member of curriculum development committee of various universities and autonomous colleges in Tamil Nadu, India. His areas of specialization are Virtual Reality, Image Processing and Sensor Networks.

R. Spurgen Ratheash is a research scholar of Sadakathullah Appa College, affiliated to Manonmaniam Sundaranar University, Tirunelveli, India. He is an Assistant Professor of Information Technology, received his MCA degree in Computer Applications from Bishop Heber College, Bharathidasan University, Trichy, India in 2007. He received his



MPhil degree in Computer Science and an M. Tech in Information Technology from Manonmaniam Sundaranar University, Tirunelveli, India in 2012 and 2014 respectively. His major research interests include Digital Image Processing, Document Image Analysis and Character Recognition of Tamil language palm leaf manuscripts.