

Recognition of Ancient Tamil Handwritten Characters in Palm Manuscripts Using Genetic Algorithm

E. K. Vellingiriraj, P. Balasubramanie

Department of C.S.E., Kongu Engineering College, Perundurai Tamilnadu India
girirajek@rediffmail.com

Abstract—The objective of this research is to develop computer software that can recognize the Ancient Tamil handwritten characters by using the genetic algorithm technique (RATHCPM). The system consists of 5 main modules, which are: 1) image acquisition module, 2) image preprocessing module, 3) feature extraction module, 4) character recognition module, and 5) display result module. Each module has the following details. First, the image acquisition module collects an unknown input character from a user. Second, the input image is transformed into a suitable image for the feature extraction module. Third, the system extracts character features from the image. There are 3 main features of Tamil characters which are stroke, loop and location of loop and stroke connection. Fourth, the extracted character information is kept in the form of bits string chromosome in a genetic algorithm. Finally, the system displays the best fitness chromosome for the recognition result.

Keywords: RATHCPM, vowels and consonant

I. INTRODUCTION

Handwritten character recognition is one of the most difficult tasks in the pattern recognition system. There are a lot of difficult things that need many image processing techniques to solve, for examples: 1) how to separate cursive characters into an individual character, 2) how to recognize unlimited character fonts and written styles, and 3) how to distinguish characters that have the same shape but different meaning such as the character o and number 0. Many researchers try to apply many techniques for breaking through the complex problems of handwritten character recognition. There are many applications that need to take advantage of the handwritten character recognition system, namely, 1) automatic reading machine, 2) non-keyboard computer system, and 3) automatic mailing classification system. [1] The objective of this research is to try to help researchers to recognize Tamil handwritten characters by using the genetic algorithm technique. The Tamil alphabet has 12 vowels, 18 consonant, combination vowels and consonant 216, and one Ayutha letter, totally 247 letters in Tamil 10 numerical symbols, as shown in Figure 1, respectively. Normally, Tamil characters consist of small circles or loops, which are connected to circular zigzag lines and straight lines. Most Tamil characters are written by using a single stroke. The structure of Tamil words is written in a four-line level style, which is shown in Figure 2.

அ ஆ இ ஈ உ ஊ எ ஏ ஐ ஓ ஔ ஶ - உயிர் எழுத்துக்கள்
ஃ - ஆயுத எழுத்து
௦ ௧ ௨ ௩ ௪ ௫ ௬ ௭ ௮ ௯ - மெய் எழுத்துக்கள்

0	1	2	3	4	5	6	7	8	9
௦	௧	௨	௩	௪	௫	௬	௭	௮	௯

Fig 1. Tamil Characters and Numbers

அம்மா
உயிர்
தமிழ் எழுத்து
கிணற்றுத் தண்ணீர்

Fig 2. The Tamil words in Four line level

II. LITERATURE REVIEWS

Historically, handwritten character recognition applications used three major approaches; the statistical approach, the structural or syntactic approach, and the neural network-based approach. This section reviews handwritten character recognition applications based on these three approaches.

A. Statistical Analysis Approach

Statistical Pattern Recognition uses statistical and/or probabilities functions for building a recognition algorithm. The input features are extracted from a set of characteristic pattern measurements. A limitation of this approach is the difficulty to express pattern classification in terms of structural information. [2, 3, 4, 5, 6, 7 and 8].

B. Structural or Syntactic Analysis Approach

Syntactic Pattern Recognition uses syntactic or structural information of patterns to generate knowledge that is related to patterns. This approach extracts the similarity of patterns and builds pattern syntax or structural rules. The information of pattern syntax rules is used to explain, classify and recognize unknown patterns. This approach is suitable for building a handwritten character recognition system because it uses a

structural approach to build unlimited handwritten character patterns syntax. A limitation of this approach is the difficulty to build learning structural rules. [9, 10, 11, 12 and 13].

C. Neural Network Based Approach

Neural Pattern Recognition emulates knowledge of how a biological neural system stores and manipulates information. This artificial neural system is called “neural networks”. The notion is that an artificial neural network can solve all problems in automatic reasoning, including a pattern recognition problem. This approach classifies patterns through predictable properties of neural networks. A limitation of this approach is a little amount of semantic information from a network. [14, 15 and 16].

III. METHODOLOGY

In this section, we present all details of our system design. First, we start with the overall framework of the Tamil handwritten character recognition system. Then, we give each component detail. Finally, we present the user interface.

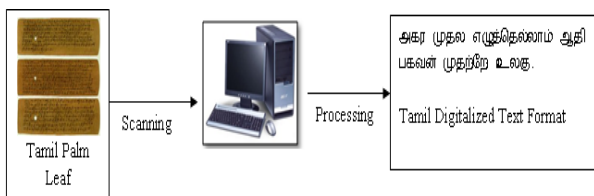


Fig 3. Framework of Ancient Tamil Handwritten Characters Recognition System

A. System Architecture Overview

Based on the Tamil handwritten character recognition system in the previous section, we design the RATHCPM system framework as shown in Figure 3. The RATHCPM has the following workflow. First, the system captures the Tamil written character images and stores them in a computer system. Second, the system extracts several features from the character images such as the number of circles, number of lines and connection locations between line and circle in each character by using image processing techniques. Third, the system uses all features of a character to generate a genetic chromosome. Fourth, the system recognizes Tamil characters by comparing genetic chromosome between unknown characters and the training character set in a database. Finally, the RATHCPM displays the best fitness genetic chromosome for the output result.

B. System Structure Chart

Based on the system framework in the previous section, we convert the RATHCPM framework to the system structure chart as depicted in Figure 4. The RATHCPM system consists of five main modules, which are 1) image acquisition, 2) image preprocessing, 3) feature extraction, 4) character recognition, and 5) display result, as shown in Figure 4. The detail of each module is described as the following:

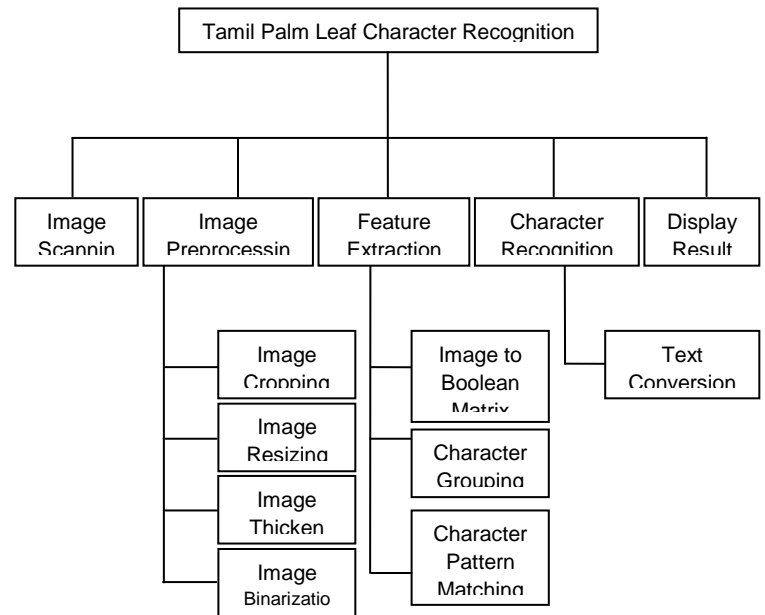


Fig 4: Structure chart of Ancient Tamil handwritten characters recognition by image zoning using the Boolean matrix

1) Image Acquisition

In the first stage, the image acquisition module captures an unknown input character from the system interface. The RATHCPM system provides a workspace and lets the user to draw an online Tamil handwritten character on the workspace. After the user finishes drawing, all the details of a written character is saved into a bitmap image and the process passes the next step.

2) Image Preprocessing

In the image preprocessing module, the system prepares a it able handwritten character image for the feature extraction module. The image preprocessing stage consists of four sub-processes, which are: 1) image cropping, 2) image resizing, 3) image binarization, and 4) image Thicken. Each sub-process has the following details.

a) Image Cropping Sub-process

The character image from the image acquisition stage has the white space that is not necessary in the recognition process. Moreover, the white space needs more CPU power for the recognition process and may cause an erroneous result. Therefore, the system needs to crop only the written character boundary. The example of the cropped character image is shown in Figure 5.



Fig 5. Input Image to Cropping Image

b) Image Resizing Sub-process

The input image may have different size, which will affect the recognition results. Therefore, every input image will be resized to 100 x 100 pixels image. The example of the resized character image is shown in Figure 6.

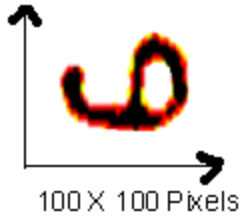


Fig 6. Resized Character by 100 X 100

c) Image Binarization Sub-process

The binarization sub-process will change an input character image into a binary image (0s and 1s only). A binary image helps the feature extraction module to extract character features easily. The example of the binarized character image is shown in Figure 7.



Fig 7. Binarized Character Image

d) Image Thicken Sub-process

The Thicken sub-process reduces a written character of a thick line into a thin character or skeleton character. A thin character is easier to extract its feature than a thick character. The example of the thinned character image is shown in Figure 8.



Fig 8. Thicken Process

3) Feature Extraction

The feature extraction module extracts the basic components of Tamil characters, such as loops, straight lines, zigzag lines and the position of the connection between the loop and the straight line. Normally, a Tamil character consists of none to three loops, none to three zigzag lines and none to five straight lines. After this module extracts the basic components from the Tamil character image, then it translates all basic components into chromosome bits string in a genetic algorithm. There are two main basic components that the RATHCPM extracts, which are: 1) a stroke extraction, and 2) a loop extraction. Both of them have the following details.

a) Stroke Detection

The stroke detection sub-process extracts the type of a line by separating a character into 5 X 5 pixels image blocks and finding the slope of a line in each block. Finally, it analyses

and classifies a line type of a character in four categories, which are 1) vertical line, 2) horizontal line, 3) zigzag line, and 4) the tail line. Each line category has the following details.

(1) Vertical Stroke Analysis

The vertical stroke analysis is a function that extracts the vertical stroke lines in the Tamil character. The RATHCPM separates one character into three vertical regions, namely: 1) Left Region, 2) Middle Region, and 3) Right Region, as shown in Figure 9 (a). The output of this function is produced in 4-bit chromosome. The first three chromosome bits string represents the vertical line in left, middle and right regions, respectively. And the fourth chromosome bit string represents the second vertical line in the right region.

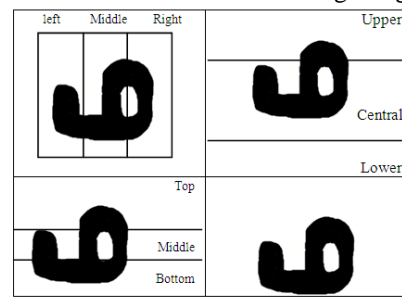


Fig 9. The Tamil Character Stroke analysis region

(2) Horizontal Stroke Analysis

The horizontal stroke analysis is a function for extracting the horizontal stroke line in the Tamil character. The RATHCPM separates a character into three regions, which are 1) Upper Region, 2) Central Region, and 3) Lower Region, as shown in Figure 9 (b). The output of this function generates 3-bit chromosome, which are 1) Upper Region, 2) Central Region and 3) Lower Region.

(3) Zigzag Stroke Analysis

The zigzag stroke analysis in RATHCPM means the stroke that contains a turning point. The zigzag line in a Tamil character is divided into 3 zones, namely; 1) Top Zone, Middle Left Zone, and Bottom Zone, as shown in Figure 9 (c). The output of this function gives 3-bit chromosome, which are 1) Top zigzag line, 2) Middle left zigzag line, and 3) Bottom zigzag line. (4) Tail Stroke Analysis. The tail stroke analysis is a function to extract the tail stroke of the Tamil character. Some Tamil characters have a long tail stroke. RATHCPM focuses on the location of a tail.

There are 2 positions, which are Upper tail and Lower tail, as shown in Figure 9 (d). The output of this function produces 2-bit string chromosome.

b) Loop Detection

The loop detection sub-process extracts three characteristics of the loop, which are 1) number of loops in a character image, 2) position of each loop in a character image, and 3) type of

the loop. The loop detection sub-process has the following details.

(1) Number of Loops Analysis

Based on the observation, every Tamil character consists of zero to three loops, for example, the Tamil character “**f**” has no loop, character “**n**” has one loop, character “**d**” has two loops and character “**z**” has three loops. The RATHCPM applies a color filling algorithm to find the loops in a character image. The concept of a color filling algorithm is filling a white color pixel in an image with a black-color background. A background will change to a white-color if there is no completed loop in the image. The color filling algorithm process is shown in Figure 10. The output of this process generates 3-bit string chromosome to flag number of loops in a character.

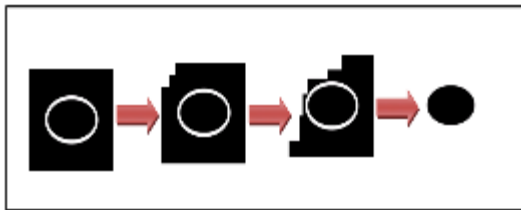


Fig 10. The Color filling algorithm to find a loop

(2) Loop Location Analysis

Each loop in a Tamil character is located in nine locations, which are 1) Top-Left, for example, a character “**t**”, 2) Top-Middle, for example, a character “**fp**”, 3) Top-Right, for example, a character “**k**”, 4) Middle-Left, for example, a character “**s**”, 5) Center, for example, a character “**b**”, 6) Middle-Right, for example, a character “**n**”, 7) Bottom-Left, for example, a character “**n**”, 8) Bottom-Middle, for example, a character “**s**” and 9) Bottom-Right, for example, a character “**x**”. Each character has maximum 3 loops and each loop has 9 locations; therefore, there are 3X9 equal to a 27-bit chromosome.

(3) Loop Type Analysis

There are eight points in a loop that can connect to a line which are 1) Top-Left, 2) Top-Middle 3) Top-Right, 4) Middle-Left 5) Middle-Right 6) Bottom-Left 7) Bottom-Right and 8) Bottom-Middle. Each character has maximum 3 loops and each loop has 8 connection points to a straight line; therefore, there are 3X8 equal to a 24-bit chromosome.

4) Character Recognition

The image recognition module uses information from feature extraction to form a genetic chromosome, after that it uses chromosome string to recognize the Tamil character mage. This module consists of 2 functions, namely, 1)

chromosome generation function, and 2) chromosome evaluation function. Each function has the following details.

(1) Chromosome Generation Function

The chromosome generation function produces the Tamil character chromosome by combining all features extracted in he previous module together, namely:

- 3-bit for number of loops
- 27-bit for location of each loop
- 24-bit for loop connected with a straight line
- 12-bit for location of the lines

There are 66-bit chromosomes in the Tamil character. The sample of a Tamil character “**t**” chromosome bits string are “11100 00100 00110 00010 00000 00000 01000 00000 00000 00001 00101 00000 00000 0”, where 0s is none of the features are shown in a character, and 1s represent a character image shows that feature. The meaning of each chromosome bit string..

(2) Chromosome Evaluation Function

The chromosome bit string from the chromosome generation function is used to recognize a character by comparing the fitness value of an unknown character with all Tamil characters in the database. The highest fitness value is the recognition result. The fitness value is calculated by using Equation 1 as the following:

$$\text{Fitness Value} = \sum_{i=1}^{66} |(S_i+1.0) - (L_i+1.0)| * W_i \quad (1)$$

Where S is a chromosome bit string in database L is a chromosome bit string of an unknown character W is weight of each chromosome bit string.

5) Design Results

The RATHCPM system interface has 2 panels. The first one is a workspace (label number 1). This space is provided for drawing a Tamil character to test the system. The next one is a result panel (label number 2). This panel is used to show the picture of a character which is a result of recognition. Moreover, there is a label that shows the recognition process time (label number 3) and fitness value of the input character (label number 4). There are 2 types of tools on the RATHCPM interface. There are drawing and recognition tools. In a drawing tool, there are write button (label number 5) and clear button (label number 6). The write button will call the drawing function that lets the user draw on the workspace. The clear button is used for clearing workspace. For a recognition tool, there is a recognition button (label number 7). The recognition button accesses a recognition method to identify an image from the workspace and shows the result in the result panel. There is a group of radio boxes that are used to specify the line level of written characters. A user need to tell the system whether the input character is consonant (label number 8), ordinary vowel (label number 9), vowel in upper level (label number 10), or vowel in lower level (label number 11).

IV. CONCLUSIONS

In this paper, we fulfill our research objective by applying the genetic algorithm technique for recognizing Ancient Tamil handwritten characters, based on the basic features of handwritten characters namely, 1) loop, 2) line, and 3) location of loop and line connection. The system generates 66-bit string chromosome to represent a handwritten character. Then the system uses the 66-bit string chromosome to identify each handwritten character. The Tamil handwritten character recognition system still has many difficult problems that need more complex algorithms to solve. For example, to recognize a cursive handwritten character, the recognition system needs a more effective character segmentation algorithm to separate a cursive character into individual characters and needs a more efficient character recognition algorithm to identify uncertain handwritten shapes. Because every handwritten character has unlimited shapes, patterns and styles even when it is written by the same or a different person. So it is difficult to define a standard structure, for a general algorithm to explain uncertain handwritten characters.

REFERENCES

- I. Chomtip Pornpanomchai, Dentcho N. Batanov and Nicholas Dimmitt, "Recognizing Thai handwritten characters and words for human-computer interaction", *International Journal of Human-Computer Studies*, pp. 259-279, (2001)
- II. Chomtip Pornpanomchai, Pattara Panyasrivarom, Nuttakit Pisitviroj and Piyaphume Prutkraiwat, "Thai Handwritten Character Recognition by Euclidean Distance", *The 2nd International Conference on Digital Image Processing (ICDIP 2010)*, pp. 53-58' (2010)
- III. Chomtip Pornpanomchai and Montri Daveloh, "Printed Thai Character Recognition by Genetic Algorithm", *The International Conference on Machine Learning and Cybernetics*, pp. 3354-3359, (2007)
- IV. Boontee Kruatrachue, Nattachat Pantrakarn and Kritawan Siriboon, "State Machine Induction with Positive and Negative for Thai Character Recognition", *The International Conference on Communications, Circuits and Systems*, pp. 971-975, (2007)
- V. Parinya Sanguansat, Widhyakorn Asdornwised and Somchai Jitapunkul, "Online Thai Handwritten Character Recognition Using Hidden Markov Models and Support Vector Machines", *The International Symposium on Communications and Information Technologies*, pp.492-497, (2004)
- VI. Rud Budsayaplakorn, Widhayakorn Asdornwised and Somchai Jitapunkul, "On-line Thai handwritten character recognition using hidden Markov model and fuzzy logic", *The IEEE 13th Workshop on Neural Networks for Signal Processing*, pp. 537-546, (2003)
- VII. Kritawan Siriboon, Apirak Jirayusakul and Boontee Kruatrachue, "HMM topology selection for on-line Thai handwriting recognition", *The First International Symposium on Cyber Worlds*, pp. 142-145, (2002).
- IX. Arrak Pornchaikajornsak and Arit Thammano, "Handwritten Thai character recognition using fuzzy membership function and fuzzy ARTMAP", *The IEEE International Symposium on Computational Intelligence in Robotics and Automation*, pp. 40-44, (2003)
- X. Supachai Tangwongsan and Orawan Jungthanawong, "Refinement of stroke structure for printed Thai character recognition", *The 9th International Conference on Signal Processing*, pp. 1504-1507, (2008)
- XI. Khampheth Bounnady, Boontee Kruatrachue and Takenobu Matsuura, "Online Unconstrained Handwritten Thai Character Recognition Using Multiple Representations", *The International Symposium on Communications and Information Technologies*, pp. 135-140, (2008)
- XII. Jarernsri L. Mitranont and Urairat Limkonglap, "Using Contour Analysis to Improve Feature Extraction in Thai Handwritten Character Recognition Systems", *The 7th IEEE International Conference on Computer and Information Technology*, pp.668-673, (2007)
- XIII. Hyung Il Koo and Nam Ik Cho, "Text Line Extraction Chinese Documents Based on an Energy Minimization Framework", *IEEE Trans. On Image Processing*, Vol.21, no.3, pp 1169-1175, Mar 2012.
- XIV. G. Nagy, S. Seth, and M. Viswanathan, "A Prototype Document Image Analysis System for Technical Journals," *Computer*, vol. 25, no. 7, pp. 10-22, Jul. 1992.
- XV. F. Shafait, D. Keysers, and T. M. Breuel, "Performance Evaluation and Benchmarking of Six-Page Segmentation Algorithms," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 30, no. 6, pp. 941-954, Jun. 2008.
- XVI. Y.Liang, M.C.Fairhurst & R.M.Guest, "A Synthesis Word Approach to Word Retrieval in Handwritten Documents", *Elsevier Pattern Recognition*, Vol.45, PP 4225-4236, June 2012.
- XVII. Giuseppe Pirlo, Donato Impedovo, "Adaptive Membership Functions for Handwritten Character Recognition by Voronoi-Based Image Zoning", *IEEE Trans on Image Processing*, Vol 21, No 9, PP 3827-3836, Sep 2012.
- XVIII. Chomtip Pornpanomchai, Verachag Wongsawangtham, Satheanpong Jeungdomporn, and Nannaphat Chatsumpun, "Thai Handwritten Character Recognition by Genetic Algorithm (THCRGA)", *IACSIT Journal of Engineering and Technology*, Vol 3, No 2, Apr 2011.
- XIX. Qiu-Fend Wang, Fei Yin, and Cheng-Lin Liu, "Handwritten Chinese Text Recognition by Integrating Multiple Contexts", *IEEE Trans on Pattern Analysis and Machine Intelligence*, Vol 34, No 8, Aug 2012.
- XX. Tiji M Jose and Amitabh Wahi, "Recognition of Tamil Handwritten Characters using Daubechies Wavelet Transforms and Feed-forward Back Propagation Network", *IJCA*, Vol 64, No 8, PP 0975-8887, Feb 2013.
- XXI. Jin Chen, Daniel Lopresti, "Model Based Ruling Line Detection in Noisy Handwritten Documents", *Pattern Recognition Letters*, Elsevier, 2012.
- XXII. A Bharath and Sriganesh Madhvanath, "HMM-Based Lexicon-Driven and Lexicon-Free Word Recognition for Online Handwritten Indic Scripts", *IEEE Trans on Pattern Analysis and Machine Intelligence*, Vol 34, No 4, Apr 2012.