# Off-line Handwritten Thai Name Recognition for Student Identification in an Automated Assessment System

Hemmaphan Suwanwiwat, Vu Nguyen, and
Michael Blumenstein
School of Information and Communication Technology,
Griffith University, Australia

Umapada Pal
Computer Vision and Pattern Recognition Unit
Indian Statistical Institute
India

*Abstract*—In the field of pattern recognition, off-line handwriting recognition is one of the most intensive areas of study. This paper proposes an automatic off-line Thai language student name identification system which was built as a part of a completed off-line automated assessment system. There is limited work undertaken in developing off-line automatic assessment systems using handwriting recognition. To the authors' knowledge, none of the work on the proposed system has been performed on the Thai language. In addition the proposed system recognises each Thai name by using an approach for whole word recognition, which is different from the work found in the literature as most perform character-based recognition. In this proposed system, the Gaussian Grid Feature (GGF) and the Modified Direction Feature (MDF) extraction techniques are investigated on upper and lower contours, loops from full word contour images of each name sample, and artificial neural networks and support vector machine are used as classifiers. The encouraging recognition rates for both feature extraction techniques were achieved when applied on loop, upper and lower contour images (99.27% accuracy rate was achieved using MDF on artificial neural networks and 99.27% using GGF with a support vector machine classifier).

*Keywords*—*off-line handwriting recognition; automated assessment system; student identification system; modified directional feature; Gaussian grid feature*

## I. INTRODUCTION

Handwriting recognition is divided into 2 types, which are off-line and on-line handwriting recognition. Off-line handwriting recognition is one of the most difficult and challenging tasks in pattern recognition. It performs recognition of written documents by using a scanner. The hard copy document is commonly transformed into a binary pattern [1] which allows the recognition system to process the binarised handwritten image.

Off-line handwriting recognition is considered more difficult than its on-line counterpart as it cannot capture the written information whilst performing the writing as is the case in on-line handwriting recognition [2]. Nevertheless, many applications benefit from off-line handwriting recognition techniques, for example, postal address interpretation, signature verification, and bank cheque verification. However, there is only a small amount of research focusing on off-line assessment systems [3], [4], [5], [6]. To the best of the authors' knowledge, there is no off-line Thai language automated assessment system proposed in the literature.

Manual assessment of handwritten examinations is a complex task; it requires the marker's attentiveness, correctness and it is time consuming. An important part of the examination paper, besides the exam questions and answers themselves, are student name components, which are the name and last name, and student number. Commonly for manual assessment, when the marking of examination papers is concluded, the marker has to rewrite each student's mark into a report marking sheet. One problem of transcribing the mark of each student is that it could be error prone as the assessor may mistakenly ascribe the examination mark against the wrong student name.

This paper proposes a sub-system of an Off-line Automatic Assessment System (OFLAAS) called Student Identification System (SIS) similar to previous work the authors have proposed earlier [6]. However, in the present work, the experiments were performed on the Thai language to recognise student name components. To the authors' knowledge, there has not been much work done on Thai whole word recognition, and especially not on writer identification as per the work that is proposed here. This is due to the nature of the Thai language (please refer to sub-section II-B). Features used in the proposed system are different from the previous work as in the present work, the features were extracted from upper, lower, and loop images rather than full contour images.

Also in this proposed system, Artificial Neural Networks (ANNs) and Support Vector Machines (SVMs) were used as classifiers to compare the recognition rates on both techniques, rather than only those applied to ANNs as in the previous work. It must be noted that the system proposed here has no intention to verify student identities, only to identify them. The SIS with the ability to verify students may be proposed in future work.

It should be noted that the student number has not been used in this research as a student identification and verification system could be developed on student name components in the future. That is, the future system would be able to verify if the person who sat an examination is the same person who owns the name recognised using the student handwritten name and

the last name. Also as normally students do not sign examination papers, this research proposes the use of name components in developing the SIS but not the student signature or number.

For the proposed system, once the marking process (main process of OFLAAS) is completed, the report on each student's mark is automatically produced. Having such a system would reduce the chance of mismatching between the student's names and their marks. Another advantage of having the proposed system is that the list of the students who are absent from the examination can be produced automatically.

As stated above, the proposed system intends to identify students from their handwritten name and last names but not verifying if the written names have been forged. This research investigates and compares the performance of the Modified Direction Feature (MDF) and Gaussian Grid Feature (GGF) extraction technique. Since the proposed system is used to recognise Thai words, which have different characteristics compared to English, different input images rather than using only boundary images, were employed with ANNs and SVMs to obtain the highest recognition rates possible.

As there is no suitable database of Thai name components available, a new database was created to be used for experimentation of the proposed system. The database used in the present research consists of 2,060 handwritten name components from a total number of 103 writers.

The remainder of this paper is organised into three sections. Section II describes the methodology employed in this research. Section III details the results obtained and puts forward a discussion and analysis. Finally, conclusions are drawn in Section IV and future work is also described.

## II. METHODOLOGY

This section discusses the methodology and techniques used in conducting the research. The topics in this section include the proposed system (block diagrams), data collection, nature of Thai language, proposed methodology, and the experimental setup.

Significant research has been undertaken in the area of off-line character and handwriting recognition. Nevertheless, to the authors' knowledge, there has only been a limited amount of work in the literature reporting the development and investigation off-line automatic assessment systems. Also there has not been any research undertaken for student identification using Thai handwritten name components written on examination papers.

The proposed methodology includes data collection, image processing, effects of different input images to each feature extraction technique, and the investigation of the MDF and GGF techniques in conjunction with classifiers employing different parameters in order to achieve the optimum results. Classifiers used in conducting the proposed SIS are ANNs and SVMs. Fig. 1a. illustrates a block diagram of a complete OFLAAS which consists of a main component including a short answer question automatic marking module and a SIS, which is a sub-system of OFLAAS. The OFLAAS is used to mark each student's examination paper, and is also used to identify students from their name components. Once both processes are completed, the full report containing a list of students who attended the examination along with the marks they achieved is produced.

The process of SIS begins with the data collection of the students' name components. The scanning process is used to transform raw data into digitised patterns. Binarisation and preprocessing, including line and word segmentation, noise removal, filling and skew correction are then applied to the images.

The feature extraction techniques which were selected in the proposed system are the MDF and GGF. The MDF and GGF extraction techniques were chosen due to their ability to successfully extract those important features from images, which have enabled accurate recognition rates to be attained in a number of applications [6], [7], [8]. After the feature vectors are generated by employing each technique, the features are then applied to the ANNs and SVMs for training, and testing for the recognition/identification process. The SIS recognition accuracy rate was evaluated once the results were obtained. A proposed SIS block diagram can be found in Fig. 1b.

### A. Data Collection

There is no publicly available dataset of Thai language handwritten name components from examination papers; as a result, a data collection process was performed to create a custom dataset. The dataset collected for the proposed system is the first database of its type in the Thai language.

In the research proposed here, the recognition of words was based on one writer per name components. Although in some cases, the writers (students) may have had the same name components (such as their first or last name); the system will identify the students by using the rest of their name components, as well as by using other criteria (please refer to Section III). Upon request, the database is available for download to the research community.
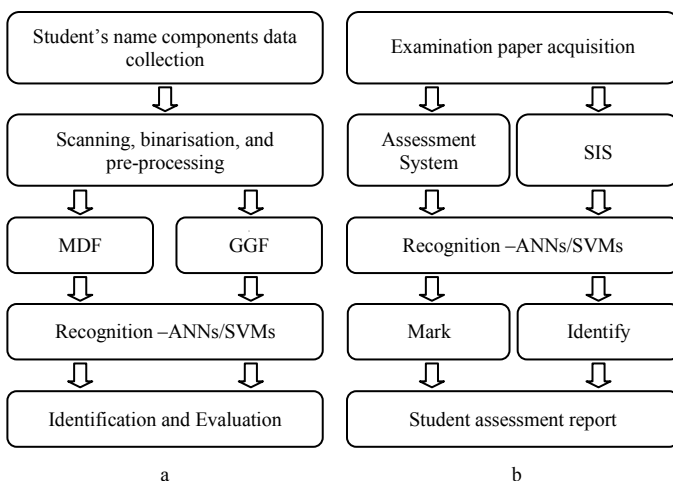
Fig. 1a Block diagram illustrating a complete Off-Line Automatic Assessment System (OFLAAS). Fig. 1b. A block diagram illustrating the proposed Student Identification System (SIS)

Fig. 2. Name component examples. Intra-class differences, various lengthsofname components, and duplication of name components written by different writers can be found.

The restricted number of one hundred and three volunteers was determined based on an assumption that for most classes, the number of students per class is less than one hundred, except for larger classes such as some lectures in universities. Therefore, one hundred and three volunteers, both male and female, provided handwritten Thai language first and last names in order to be used for training and testing the proposed SIS.

In total, there were 2,060 (206 name components x 10 samples of each name component) samples obtained. All samples were written with minimum constraints e.g. writing instruments and handwriting styles were not restricted within the given space. All the name components were written in Thai.

Each volunteer wrote their name components one after another ten times. Even though the volunteer wrote their name components one after the other, intra-class differences were still found (Fig. 2). Some of the Thai name components are particularly long, for example, สุวรรณวิวัฒน์, สวรรค์บรรจบริการ, and เอก ลักษณานนท์. One of the Thai name characteristics was that the name components both first and last names may begin, end, or include some common words for name components such as "wat", "chai", "ya", "kit", and therefore can be quite confusing for automatic classification because of the common characteristics, which were shared by the volunteers. In conclusion, the dataset contains names with quite varied word lengths (from 2 - 5 syllables which can be up to 14 characters), some duplicated last names, and some of the name components containing common words within each component.

Name component examples can be found in Fig. 2. From those samples, different lengths of words were observed. Intra-class differences were found and name component duplication of different writers existed (the name component of "สุวรรณ วิวัฒน์" was written by two different writers).

### B. Thai Language

Thai language is very different from the English language. It consists of 44 consonants, 18 vowels, 4 voice tones, and 3 special symbols, which can be seen in Table I. Altogether, there are up to 69 characters (excluding Thai numerals) in the Thai language.



Fig. 3. Thai Sentence Structure and its Zones

TABLE I. THAI CHARACTERS

| Type | Type Members |
|---|---|
| Consonant | กขฃคฅฆงจฉชซฌญฎฏฐฑฒณดตถทธนบปผฝพฟ ภมยรลวศษสหฬอฮ |
| Vowels | ะาอิอีอึอือุอูเแโไใอะอ็อ็อฤฦ (where อ can be any other consonant) |
| Tones | อ่อ้อ๊อ๋ (where อ can be any other consonant) |
| Punctuation Marks | อ์ๆฯ (where อ can be any other consonant) |

Heads of Thai characters, which are small loops, play an important role because many of the characters look the same. Only the position or whether or not the head is present will tell what character it is. The heads of a character can be found in various locations. When classifying Thai characters using the characters heads, three categories are obtained. These three categories are:

*1) No-head character*, for example, "ก" for a consonant, "า" for a vowel, and "I" for a voice tone mark.

*2) One-head character*, for example, "ข", "ผ", "ล", "อ" for consonants, and "l" for a vowel.

*3) Two-head character*, for example, "ฌ", "ฮ", "ฬ", "ญ" for consonants, and "ะ" for a vowel.

An obvious example that shows the importance of the head is these three different characters that are only differentiated because of their heads: "ก", "ภ" and "ถ". Loops (heads of characters) can be found in many positions, including upper-left part of a character ("บ"), upper-right part of a character ("ห"), middle part of a character ("ฆ"), lower left part of a character ("ม"), lower right part of character ("น").

There is no space between words for the Thai language, which makes it hard to segment words from a sentence. For example "ตากลม" can mean sitting in the wind (ตาก ลม) or sitting with your eyes wide open (ตากลม). However this is not the problem in the research proposed here as the names were recognised as a whole (whole word recognition).

The last aspect to mention here is zoning. The Thai language structure can be classified into four main levels which are upper zone 2, upper zone 1, middle zone, and lower zone (see Fig. 3 for details).

### C. Datasets

For both GGF and MDF extraction techniques, a total number of 2,060 word samples were used to facilitate ANN training and testing. 80% of the samples of each name were

used as the training dataset, and 20% were used as the testing dataset. Altogether there are 206 x 8 = 1,648 samples used for training and 206 x 2 = 412 samples used for testing. All 2,060 samples were used in 10 fold cross-validation SVM.

### D. Image Acquisition and Preprocessing

Image acquisition was performed by using a scanner. The resolution of the images was 300 dpi; the scanned images were stored in a grey-level format. After that, binarisation was performed on each image using Otsu's thresholding algorithm [9].

Image preprocessing is a necessary operation which has to be performed before the recognition process. The main purpose of image preprocessing is to develop useful canonical descriptions of shapes and surfaces in the given image [10]. Outputs from these preprocessing steps are used as input for the training and recognition phase; hence, some preprocessing techniques were applied in order to achieve the most valid output.

Automatic line and word segmentation were performed using histogram projection. Line segmentation was performed first, and then word segmentation was performed in order to obtain each of the name components (first and last name). Salt and pepper noise removal and filling techniques were also performed. Skew normalisation [11] was finally applied to each name image. Slant correction was not performed on any of the images to preserve the unique characteristics of each student name.

For each image, boundary extraction was performed to isolate connected components. A boundary extraction algorithm was employed [12]. After that, loop as well as upper and lower contour extraction are performed. Examples of preprocessing successive images are presented in Fig. 4.

### E. Feature Extraction

Feature extraction is one of the most crucial components in a handwriting recognition system. The objective of feature extraction is to extract the salient information that needs to be applied in the recognition process. It reduces data by determining certain feature properties that distinguish input patterns [13].

Feature vector sizes of the two techniques are quite different, whilst the MDF vector size is 121, the GGF vector size is 864. It is noted that the 2 feature extraction techniques were applied separately and were not combined in any way. Originally, the MDF was created to extract direction and location information at background to foreground transitions from handwritten characters. Hence, the technique was developed to analyse information at the character level. However, in the proposed system, the MDF will be implemented to extract information from the whole word (each name component) image. Likewise, the recognition process will utilise holistic name information rather than recognising the handwriting at the character level.

As the nature of Thai language may contain loops in each character and that the loops play an important role in
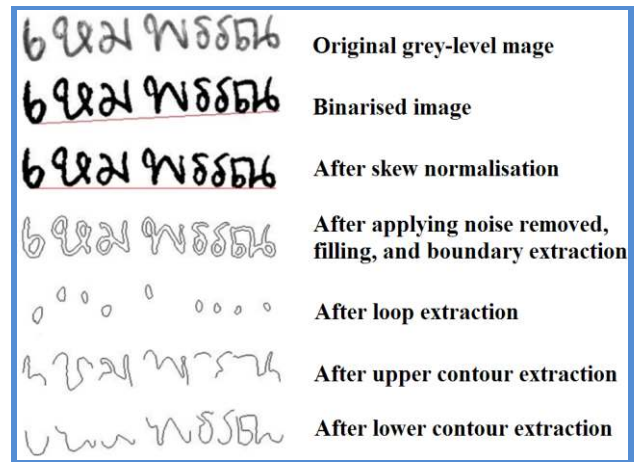


Fig. 4. Example images after each preprocessing step

distinguishing characters (please refer to Section B), the research proposed here uses loop image information extracted from full boundary images (please refer to Fig. 4). Also, because the MDF was initially created to be used at the character level, in this proposed research we also use upper and lower contours extracted from full boundary images. This is to reduce the number of transitions in each direction when applied to the MDF. It must be noted that at the word level, the number of transitions is generally larger than at the character level; therefore, to simplify transitions in each direction, upper and lower contour images are used.

As illustrated in Fig.5b, when finding transitions from background to foreground (1, 3, 5) and from foreground to background (2, 4), transitions of each direction can be more satisfactorily covered when applied on upper or lower contour images.

1) *The Modified Direction Feature Extraction Technique (MDF):* The MDF builds upon the direction feature [14]. The main difference in MDF is the way the feature vector is created. For MDF, feature vector creation is based on the calculation of transition features from background to foreground pixels in the vertical and horizontal directions.

Both the location transitions (LTs) calculated, and the direction value at that location are stored. The feature extraction processes includes 1) Determining the LT values, and 2) Determining the Direction Transitions (DTs).
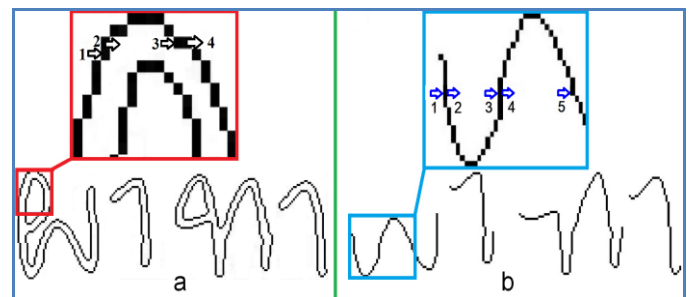


Fig. 5a. Example of transitions on full boundary image and Fig. 5b. example of transitions on lower contour .
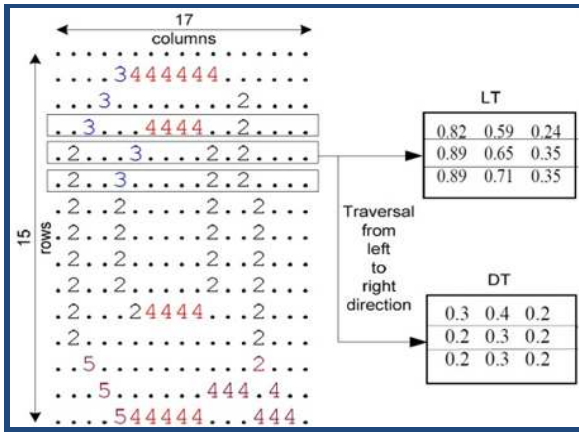
Fig. 6. Processing of DT and L T values in the left-to-right direction [14].

To determine the LT values, scanning the image in each row from left to right and right to left is necessary. Each column in the image also must be scanned from top to bottom and bottom to top. A fraction value of the distance traversed across the image is attained by computing the LT values in each direction. A maximum value was defined to be the largest number of transitions that may be recorded in each direction. DT values can be found after a transition in a particular direction is determined. The direction value at the position is stored along with its corresponding LT value. DT values are obtained by dividing the direction value by a predetermined number.

As a result, four vectors would be present for each set of feature values. Another process which has to be performed is re-sampling. Re-sampling of the vector is necessary to ensure that the dimensions are normalised in size. An example of processing of DT and LT values in the left-to-right direction can be found in Fig. 6. Further information regarding MDF can be found in [14].

In the research proposed here, when applying MDF to the three images, which are loops, upper, and lower contours, the vector size is increased to 361 ((120 × 3) + 1 = 361, where 1 is the ratio between width and height). The vector size remains 121 when applying MDF to the full boundary image.

*2)    The Gaussian Grid Feature Extraction Technique (GGF):* The Gaussian Grid Feature [7] is a relatively new feature extraction technique. Originally, it was developed for the signature verification problem. The GGF employs pattern contours as its input. From the contour representation of a name component image, the contour image is divided into 12 × 12 zones of equal size. By tracing the contours in each block, the 4-direction chain code histogram of each block is created. Every step from a pixel to its adjacent one of the four directions, which is either horizontal, vertical, left diagonal, or right diagonal, are tallied.

A Gaussian smoothing filter ($\sigma$ = 1.2) is then applied to each directional 12 × 12 matrix. The value of each element of each matrix obtained in the previous step is then normalised by dividing its value by the maximum value of the four matrices.

From the two-matrix pairs, horizontal (H) and vertical (V) matrices, left-diagonal (L) and right-diagonal (R) matrices, two new matrices ($\oplus$ and $\otimes$) are established by manipulating pairs of matrices of the perpendicular directions. The feature vector is formed by merging the six matrices H, V, L, R, $\oplus$ and $\otimes$. The dimension of the output feature vector is 12 × 12 × 6 = 864.

In the research proposed here, when applying GGF to the three images, which are loops, upper, and lower contours, the vector size is increased to 2592 (864 × 3 = 2592). The vector size remains 864 when applying GGF to the full boundary image.

*F.   Classification*

The ANNs were trained with the resilient backpropagation algorithm, which was selected above all others to address the problem of magnitudes of the partial derivative effects when using the sigmoid function (details about ANN settings can be found in sub-section G). For both MDF and GGF features, the neural networks were trained using 206 × 8 = 1,648 samples, and tested using 206 × 2 = 412 samples.

For the experiments using SVM, libsvm [15] was employed in conjunction with the WEKA toolkit [16]. For training the SVMs, ten-fold cross validation was used across all 2,060 handwriting samples.

*G.   Experimental Settings*

For both ANN and SVM settings and structure, there were 206 outputs for the 206 first and last names. The duplicated name components from different writers, for example "สุวรรณวิวัฒน์" from "เหมพรรณ สุวรรณวิวัฒน์" and "สุวรรณวิวัฒน์" from "ธนวัฒน์ สุวรรณวิวัฒน์" were classified into 2 different outputs. However, in the recognition phase, "สุวรรณวิวัฒน์" can be recognised as either "สุวรรณวิวัฒน์" of เหมพรรณ's output or of ธนวัฒน์'s output. The SIS will identify who the name component belongs to. Classification into 2 different outputs was carried out because in the future it is believed that this will be useful in developing the SIS that can identify and verify students from their name components.

For ANNs, the number of hidden units investigated during training was experimentally set from 30 up to 120 hidden units. The number of iterations set for training increased from 1000 up to 20000. All ANNs were trained with a learning rate of 0.l, and a momentum rate of 0.1. For the SVM settings, 10-fold cross validation was performed to get statistically meaningful results. The Gaussian kernel was used and the C parameter of the SVM was set to be 100.

*H.   SIS Evaluations*

Classification accuracy rates were evaluated by using the unseen testing dataset before being applied to the SIS. The unseen testing dataset included 412 samples (20% of the 2,060 samples). Only the best classification outcomes using each feature extraction technique (GGF/MDF) from both classifiers (ANN and SVM) were applied to the proposed system.

## III. EXPERIMENTAL RESULTS AND DISCUSSION

This section presents recognition rates obtained from the ANNs and SVMs classifiers trained using either GGF or MDF extraction techniques. Recognition rates were obtained through experimentation with the testing datasets (20% from the total number of samples available). ANNs and SVMs were trained using either the GGF or the MDF extraction technique, and tested with the remaining 412 unseen samples. The best recognition rates, together with their settings are displayed in Table II.

The best recognition rates (RR) obtained from applying MDF on upper, lower contours and hole images, contain somewhat less information about the handwriting than full boundary images (refer to Table II). This can be explained by the fact that the number of transitions in each direction was significantly reduced, so that the MDF, which was originally created to be applied at the character level, can extract more details from the image. Also with fewer transitions, less noise can be found.

For GGF, which was originally designed for the signature verification problem, attained the highest recognition rate of 99.27% when features were extracted from upper, lower contour and loop images using SVM as the classifier. The reason that the result is better than the recognition rate (98.40%) where features were extracted from the full boundary/contour images is that the full contours may have been too complex in the case of some classes and has introduced uncertainty in terms of additional, non-salient information.

For both ANNs and SVMs, it is also noted that the training time required when applying the MDF feature is shorter than the GGF due to its smaller feature vector size. After each of the name components were recognised, the SIS was used to map them against each student's name components. In order for the system to identify the student which the name components belong to, it employs the criteria described in Fig. 7.

When comparing the recognition rates of classifiers employing GGF or MDF extraction techniques and the identification accuracy of the SIS with the human marker together with the above criteria, a 100% accuracy was attained. The 100% accuracy was obtained, mainly because of the efficiency of the feature extraction techniques used, combined with an ability of the SIS and its criteria to identify who the name components belong to, and its ability to reject some name components for manual identification (see Fig. 7).

There is one work related to student identification systems available in the literature. However, the related work was performed in the context of assessment in the English language [6]. The comparison of the existing system and the proposed system can be seen in Table III. As can be seen, slightly different recognition rates were obtained for both languages. Slight improvement in recognition rates were found when extracting features from less detailed images (i.e. loops, upper and lower contours). From the results in Tables II and III, it can be seen that MDF yielded better recognition rates when applied to ANNs, and GGF yielded better recognition rates when applied to the SVM.

TABLE II.  RAW RECOGNITION RATES ATTAINED EMPLOYING GGF OR MDF FEATURE EXTRACTION TECHNIQUES IN CONJUNCTION WITH THE ANN AND SVM.

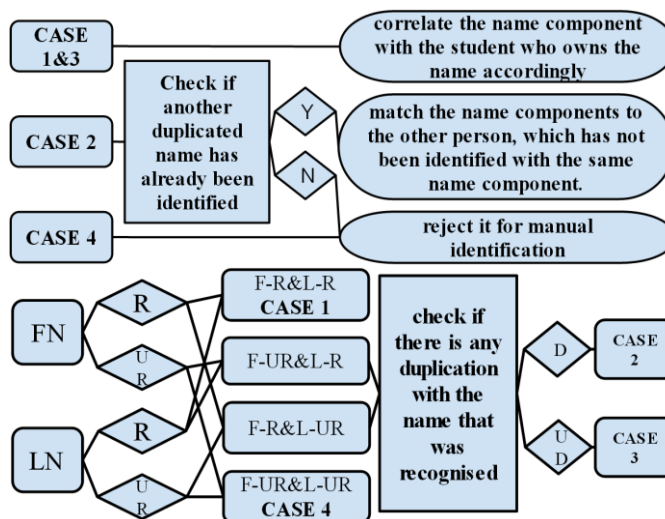| Feature Extraction Technique | ANN | | | SVM |
|---|---|---|---|---|
| | RR (%) | HU | Iterations | RR (%) |
| MDF on full boundary image | 99.03% | 60 | 15000 | 98.35% |
| GGF on full boundary image | 97.33% | 74 | 4000 | 98.40% |
| MDF on upper, lower contour, and loop images | 99.27% | 70 | 4000 | 97.52% |
| GGF on upper, lower contour, and loop images | 80.58% | 58 | 11000 | 99.27% |
| MDF on upper and lower contour images | 97.57% | 63 | 19000 | 95.73% |



Fig. 7. SIS criteria decisions based on whether or not first (FN) or last (LN) name components are recognised (R) or unrecognised (UR), and if just one of the components is recognised, or whether there is any duplication (D) of first or last names.

As mentioned earlier, most of the Thai language recognition techniques found in the literature were performed based on character recognition [17]. Feature extraction techniques such as the transition profile, Fourier descriptor, and edge direction, etc. were applied in that work; they were reported to yield the recognition rate of 92.94%. As the work proposed here is holistic recognition, MDF and GGF were chosen to be applied to the less detailed components of whole word images (upper, lower contours, and loops). MDF and GGF were chosen because of their abilities to extract details efficiently from whole words as well as at the character level. A recognition rate of 99.27% is encouraging for both MDF and GGF.

## IV. CONCLUSIONS AND FUTURE RESEARCH

This research presented a comparative performance analysis of two distinct feature extraction techniques within the context of a whole-word recognition system. The feature extraction techniques investigated were originally designed for different purposes, character recognition and signature verification, rather than handwritten whole-word recognition.

The experimental results obtained in this research indicated that the performance of both techniques were encouraging and comparable. The highest accuracy obtained by applying the

TABLE III. COMPARISON BETWEEN THE RESULTS OF THE PROPOSED SIS RECOGNITION RATES AND PREVIOUS WORK

| System – Technique – Classifier | Recognition Rate (%) |
|---|---|
| English SIS – MDF – ANN [6] | 99.02% |
| English SIS – GGF – ANN [6] | 98.28% |
| Thai SIS – MDF – ANN | 99.27% |
| Thai SIS – GGF – SVM | 99.27% |

GGF technique is 99.27% which is similar to that which was obtained when the MDF extraction technique was employed, though the results were attained by using different classifiers.

Despite its superiority for the signature verification problem (in previous work [7]), the GGF did not outperform the MDF for the problem of handwritten Thai word recognition. This can be explained by the higher intra-class variation observed in handwritten words compared to signatures.

What is interesting to observe is that comparable recognition results can be attained by using full contour information from the word or targeted features from the word contour. This demonstrates the usefulness of the proposed feature extraction approaches, which reduce the amount of information used for training (upper, lower contours and loop information) as well as testing the classifier for the first time on a Thai handwritten dataset.

It is encouraging to find that both MDF and GGF can work better with fewer details when applied to full contour images, though with different classifiers. This is a novel finding, however more experiments will be carried out on a larger database and also on English words to observe the results. More experiments with other baseline feature extraction techniques, such as local binary patterns and the histogram of oriented gradients will be carried out for comparison purposes.

In addition, the work will be extended to student verification so that the system can detect if the students who sat an examination were really the persons who own the name and not someone else.

REFERENCES

[1] E. Anquetil, and H. Bouchereau, "Integration of an on-line handwriting recognition system in a smart phone device," In proc. of the 16th International on Pattern Recognition, 2002, vo,3, pp. 192- 195 vol.3, doi: 10.1109/ICPR.2002. 1047827

[2] R. Plamondon and S. N. Srihari. Online and offline handwriting recognition: a comprehensive survey, IEEE Trans. on PAMI, 22(1):63-84, 2000.

[3] J. Allan. Automated Assessment Of Handwritten Scripts. PhD thesis, Nottingham Trent University, 2004.

[4] H. Suwanwiwat and M. Blumenstein. Short answer question examination using an automatic off-line handwriting recognition system and the modified direction feature. In Proc. of the 3rd International Conference on Machine Version (ICMV 2010), pages 476-480, 2010.

[5] S. Srihari, J. Collins, R. Srihari, H. Srinivasan, S. Shetty, and J. Brutt-Griffler. Automatic scoring of short handwritten essays in reading comprehension tests. Artificial Intelligence, 172:300-324, 2008.

[6] H. Suwanwiwat, V. Nguyen, M. Blumenstein. Off-line restricted-set handwritten word recognition for student identification in a short answer question automated assessment system, In Proc. 12th International Conference on , vol., no., Hybrid Intelligent Systems (HIS), pp.167,172, 4-7 Dec. 2012.

[7] V. M. Nguyen and M. Blumenstein. An application of the 2D Gaussian filter for enhancing feature extraction in off-line signature verification. In Proc. 11th International Conference on Document Analysis and Recognition (ICDAR 2011), pages 339-343. IEEE, 2011.

[8] V. Nguyen, M. Blumenstein, V. Muthukkumarasamy, and G. Leedham. Off-line signature verification using enhanced modified direction features in conjunction with neural classifiers and support vector machines. In Proc. 9th International Conference on Document Analysis and Recognition (ICDAR 2007), pages 734-738, 2007.

[9] N. Otsu. A threshold selection method from gray level histograms. IEEE Transactions on Systems, Man and Cybernetics, 9(1):62-66, 1979.

[10] A D. Kulkarni. Computer vision and fuzzy-neural systems. Prentice Hall PTR, Upper Saddle River, NJ, 2001.

[11] M. Blumenstein, C. K. Cheng, X. Y. Liu (2002), New Preprocessing techniques for Handwritten Word Recognition, Proc. of the 2nd lASTED conference on visualization, Imaging and Image Processing, pp. 480-484.

[12] J. R. Parker, 1994, Practical Computer Vision using C, John Wiley and Sons, New York, NY.

[13] C.C. Tappert, C. Y. Suen, T. Wakahara, "The State of the Art in Online Handwriting Recognition," IEEE Transactions on Pattern Analysis and Machine Intelligence, August 1990, pp. 787-808, doi: 10.1 109/34.57669

[14] M. Blumenstein, X. Y. Liu, and B. Verma. A modified direction feature for cursive character recognition. In Proc. Intl. Joint Conference on Neural networks, volume 4, pages 2983-2987, 2004.

[15] Chang, Chih-Chung and Lin, Chih-Jen, LIBSVM: A library for support vector Machines, ACM Transactions on Intelligent Systems and Technology, vol:2(3), 2011, pp27:1-27:27

[16] Mark Hall, Eibe Frank, Geoffrey Holmes, Bernhard Pfahringer, Peter Reutemann, Ian H. Witten (2009); The WEKA Data Mining Software: An Update; SIGKDD Explorations, Volume 11, Issue 1.

[17] W. Chatwiriya. Off-line Thai Handwriting Recognition in Legal Amount. PhD Dissertation, 2002.