



Arabic font recognition based on diacritics features



Mohammed Lutf^a, Xinge You^{a,*}, Yiu-ming Cheung^b, C.L. Philip Chen^{a,c}

^a Department of Electronics and Information Engineering, Huazhong University of Science and Technology, Wuhan, China

^b Department of Computer Science, Hong Kong Baptist University, Hong Kong SAR, China

^c Faculty of Science of Technology, University of Macau, China

ARTICLE INFO

Article history:

Received 7 November 2012

Received in revised form

17 July 2013

Accepted 21 July 2013

Available online 3 August 2013

Keywords:

Font recognition

Arabic diacritics

Composite of central and ring projection

ABSTRACT

Many methods have been proposed for Arabic font recognition, but none of them has considered the specialty of the Arabic writing system. Most of these methods are either general pattern recognition approaches or application of other methods which have been developed for languages other than Arabic. Therefore, this paper is the first attempt to present an alternative method for Arabic font recognition based on diacritics. It presents the diacritics as the thumb of Arabic fonts which can be used individually to identify and recognize the font type. Diacritics are the marks and strokes which have been added to the original Arabic alphabet. Though they are the smallest regions in the Arabic script, with today technology it is very easy to get a high resolution image with a very low cost. In this kind of images, the diacritics can reveal very useful information about the font type. In this study, two algorithms for diacritics segmentation have been developed, namely flood-fill based and clustering based algorithm. The experiments conducted proved that our approach can achieve an average recognition rate of 98.73% on a typical database that contains 10 of the most popular Arabic fonts. Compared with existing methods, our approach has the minimum computation cost and it can be integrated with OCR systems very easily. Moreover, it could recognize the font type regardless of the amount of the input data since five diacritics, which in most cases can be found in only one word, are enough for font recognition.

© 2013 Elsevier Ltd. All rights reserved.

1. Introduction

Font recognition is a fundamental issue in the identification and analysis of documents. In the past this task was considered highly demanding on computer hardware. The OCR techniques are making great successes in commercial software, but the possibility of increasing the efficiency of the OCR system can be guaranteed by taking font type into account. Automatic document processing (ADP) techniques tackle font recognition on the basis of two main aspects. First, it generalizes the font type for all characters in the document. The use of this approach only enables us to reduce the number of alternative forms of each class of a font family. This clearly leads to the recognition of only one kind of font. The second aspect that should be considered in ADP techniques is the identification of the font types used within the document, which is usually neglected in spite of its importance [1].

Different Optical Font Recognition (OFR) methods have been successfully applied in many languages except Arabic, because these methods always fail to accord with the characteristics of the

Arabic writing system. Arabic is an alphabetic language written in a cursive way and this is not only in handwriting but also in machine typing characters. It is something between Latin and Chinese. In Latin, each word consists of letters separated from each other and in Chinese each character consists of strokes connected with each other to build one block representing a word. In Arabic, however, the word consists of one or more subwords and each subword may have one or more letters connected with each other, which makes working with an Arabic text a very challenging task. This explains why these approaches are successful with most languages but not with Arabic.

The most common approaches toward font recognition are font recognition based on typographical features [2–4] and font recognition using textural features [5–8]. It has been reported that the second approach is more efficient than the former one [9]. However, solving the problem of font recognition is still just a small task in OCR or ADP system. Most of the existing OFR methods are full pattern recognition systems; they have their own processing blocks starting from reading the image to the final step which is the recognition of the font class. When one of these types of OFR systems is added to OCR or ADP system, it increases the computational cost and mainly reduces the performance of the system especially when it is used in real time applications, for example, using OCR system to read the road signs in an autonomous car. Therefore, the speed and the possibility of integrating

* Corresponding author. Tel.: +86 2787544817; fax: +86 2787544823.

E-mail addresses: Mohammed.lutf@gmail.com (M. Lutf), you1231cncn@gmail.com, yousg@mail.hust.edu.cn, XingeYou@hust.edu.cn (X. You), ymc@comp.hkbu.edu.hk (Y.-m. Cheung), philip.Chen@ieee.org (C.L. Philip Chen).

the OFR with other systems through sharing most of the processing blocks is a critical need. Another drawback in the currently proposed methods is that when identifying the font type, the features depend on the whole given text image and it is always preceded by a preprocessing step. This fact makes the font recognition even more costly. The study reported in this paper addresses the problems discussed above and proposes instead an Arabic font recognition system which can easily be integrated with OCR system by sharing all its preprocessing blocks. At the same time, the features are extracted only from the diacritics which are very easy to be segmented.

Over the last two decades, a few approaches were proposed for Arabic font recognition. It can be classified into two main categories: segmentation-free and segmentation-based approaches. Segmentation-free approaches intend to recognize the font by extracting the texture features globally from a predefined text entity, it could be a whole image, a text block or a text line. Also, the texture features extraction algorithm is always a very well known algorithm. Systems based on this approach differ only on the type of the text entity and the algorithm used for extracting the texture features.

In [10,11], Silmane et al. used Gaussian mixture model to estimate the font category likelihoods in a word images. It was the first attempt to evaluate Arabic OFR on a publicly available database(APTI), but due to the database limitation, this approach can be used only with already segmented word images. Imani et al. [12] proposed the use of wavelet to extract features from a 128×128 text block. In this approach, most of the training dataset was labeled using a learning algorithm which may produce wrong labeled data. As a result, this will reduce the final recognition rate. In [13], Bataineh et al. used Gray-Level Co-occurrence Matrix (GLCM) features extracted from the edges of a 512×512 text block. The text block was created by removing the spaces between words and lines, but multi-size words are not normalized which will produce different edge features for the same word written by the same font. Pourasad et al. [14] proposed the use of holes of letters and text line horizontal projection profile for Farsi font recognition. This approach will fail with fonts contain no holes and also will fail with the skewness text. Khosravi and Kabir [15] proposed the use of Sobel and Roberts gradients in 16 directions to extract the texture features form a 128×128 text block for Farsi font. This approach is flexible regarding the size of the input text block, but this is also a problem because the text block is processed without normalization and the measurements depend on pixels which make this approach fail with high resolution or fewer word input text. In [16], Ben Moussa et al. demonstrated how we can use a combination of BCD (box counting dimension) and DCD (dilation counting dimension) techniques to obtain features from text paragraphs, accuracy rate of about 98% has obtained. Abuhaiba [17] proposed the use of horizontal and vertical projection profiles, Walsh coefficients, invariant moments, geometrical features for Arabic OFR using word level images. The features in this approach are extracted from common words, and the author assumes that at least the whole paragraph is written using the same font. The problem with this approach is that the segmentation algorithm of the common words during the test process may filter out all the input text. Zramdini and Ingold [3] proposed the use of scale invariant feature transform (SIFT) with 128×128 text block. The authors claim to reach a recognition rate of 100%, but the computational cost is very high especially when using big database.

For segmentation-based approach, only one method was proposed by Abuhaiba [18]. This approach segments each word into symbols then creates a template for each symbol. The problem with this approach is that it is based on segmenting the individual characters in each word which is the most complicated problem in Arabic text. The method proposed in this work which uses the

word height for symbol segmentation is not valid with the majority of Arabic fonts and it works only under ideal conditions.

In addition, all techniques used in these approaches are general techniques which can be applied to any texture analysis or recognition problem. It does not have any specific treatment for Arabic text where a simple modification or direct application of these techniques may not solve the problem of Arabic font recognition. And the absence of commercial products for Arabic OFR is an evidence.

Although the texture features approaches are robust to noisy and low resolution text images, it is reliable only with uniform and homogeneous text block where all words have the same font which is not always the case. So, to take the advantages of texture features robustness and to overcome the complexity of Arabic character segmentation, we propose a novel method¹ for Arabic font recognition using diacritics-based rotation invariant features with a low computational cost. Diacritics are not connected to each other nor to the text body which makes the diacritics very easy to be segmented. Two efficient algorithms have been developed for diacritics segmentation which are the main contribution in our method. Besides font recognition, we also address the computation simplification and font recognition of irregular text like skewness text lines, multi-font formats, and multi-language text image.

Diacritics are the most common shapes that appear in any Arabic text. Unlike the normal characters' shapes, where some of them may not be found at all; it is very easy to find hundreds of diacritics in only one page of a text. The dots diacritics, for instant, are shared by many characters and the same vowel diacritic can be attached to almost all characters. Therefore, the fact that our recognition system is based mainly on diacritics allows us to ensure getting sufficient information from any input image even if it contains only few words. Thus, our focus on diacritics does not mean that normal characters are useless; it is similar to identifying a person using only his fingerprint. If we consider the retina, face, voice, DNA and many other biometrics, we will have more discriminatory information; but as long as the fingerprint gets the task done, there is no need to include other factors. The same thing is applied to diacritics.

The rest of this paper is organized as follows. In Section 2, we introduce Arabic diacritics in detail. In Section 3, we introduce our proposed method. In Section 4, the experimental results are given. Finally, Section 5 concludes the paper and gives some perspectives of future work.

2. Arabic diacritics

Arabic is a widely used alphabetic writing system in the world [20], and it has 28 basic letters. The alphabet was first used to write texts in Arabic, most notably the Qur'an, the holy book of Islam as shown in Fig. 1. With the spread of Islam, it came to be used to write many languages like, at various times, Urdu, Pashto, Uyghur (in China), Ottoman Turkish and Spanish (in Western Europe) [21]. To accommodate the needs of these languages, new letters and symbols were added to the original alphabet. This process is known as the *Ajami* transcription system, which is different from the original Arabic alphabet. Then many modifications and improvements have been made to the Arabic writing script, which results in additional letters and strokes. The new strokes are called diacritics, and the purpose of adding these diacritics was to

1. Distinguish between letters of the same or similar shape.
2. Indicate sounds (vowels and tones) that are not conveyed by the basic alphabet.

¹ This work is an extension to our conference paper [19].

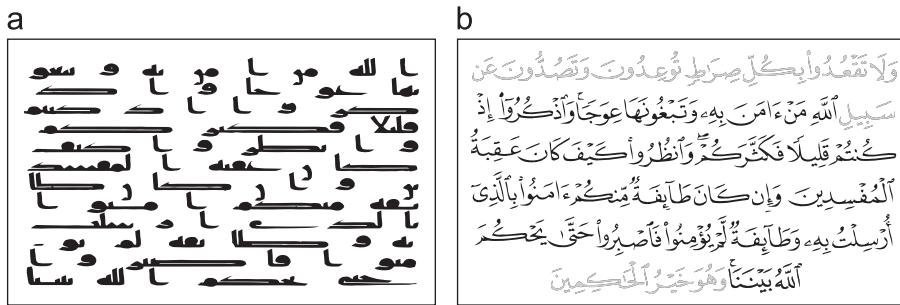


Fig. 1. Seventh-century kufic script (a), with its equivalent today script (b), diacritics are the only difference between them.

Table 1

Arabic alphabetic examples, four shapes for each letter, beginning, middle, end and in isolation. The associated numbers are the order of the letters in the alphabetic list.

No.	Beginning	Middle	End	Isolated
1			↑	-
2			↑	↑
3			↑	↑
8			↑	↑
12		↑	↑	↑
17		↑	↑	↑
18	↑	↑	↑	↑
27			↑	↑
28	↑	↑	↑	↑

3. Indicate the absence of a vowel.
 4. Clarify the difference in the meaning between words consisting of the same letters.

Nowadays, the diacritics are not something additional or optional to the language; rather they are a crucial part of it. In order to master the writing style of the Arabic alphabet, the writer needs to master the writing and the correct placing of the diacritics. The text written in modern Arabic may contain all diacritics called the *fully diacritic form* or contain only the necessary diacritics called the *basic diacritic form*. Therefore, modern Arabic writing system is composed of two main parts: letters and diacritics.

Unlike cursive writing which the Latin alphabet is based on, the standard Arabic style possesses substantially different shapes depending on whether a letter is connected with a preceding and/or a succeeding letter, hence each letter has either two or four shapes [22]. As shown in Table 1, the shapes correspond to the four positions: beginning of a subword, middle of a subword, end of a subword, and in isolation [23]. The writing style can be summarized into four steps: (1) defining the word letters, (2) selecting the appropriate letter shape with respect to its position, (3) connecting each letter to its succeeding letter unless the letter has no middle shape or it is the last letter in the word and (4) adding the corresponding diacritics to the word. Fig. 2 shows an example of this process. In this paper, we reverse the 4th step to separate the diacritics then use it for font recognition.

When typing in Arabic, all diacritics need to be entered separately except for the “dots” diacritics (diacritics 9, 10 and 11 in [Table 2](#)). This demonstrates the fact that the diacritics need to be concerned for new font design, meaning that the shape and the style of the diacritics have to follow the style and the appearance of the font. [Fig. 3](#) shows the “Hamza” diacritic written in 10 different fonts, and [Table 2](#) lists all Arabic diacritics used in modern Arabic writing.

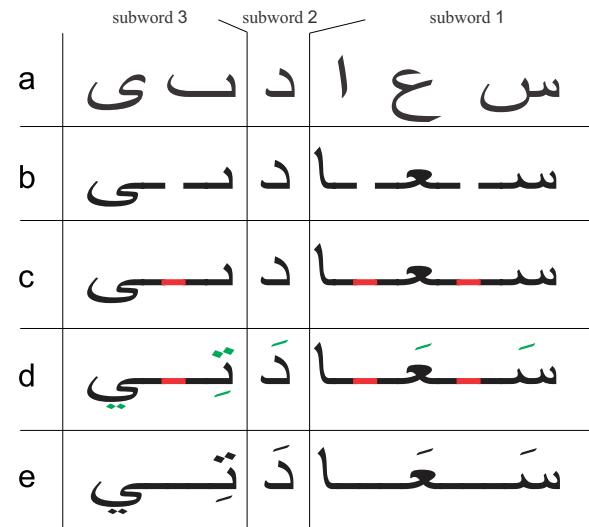


Fig. 2. Arabic writing style four steps: (a) step 1: defining the letters composing the word, (b) step 2: selecting the equivalent letter's shape with respect to its position, (c) step 3: connecting letters, (d) step 4: adding the diacritics, (e) final word. Our approach is to reverse the 4th step and separate the diacritics then use it for font recognition

Table 2
Arabic diacritics

No	Diacritic name	Shape
1	Fatha and Kasra	/
2	Tanwin (Fathateen and Kasrateen)	=
3	Damma	ؚ
4	Hamza	ؑ
5	Madda	~
6	Shadda	ؚ
7	Tanwin (Dammateen)	ؒ
8	Sukun	o
9	One dot	.
10	Two dots	..
11	Three dots	...

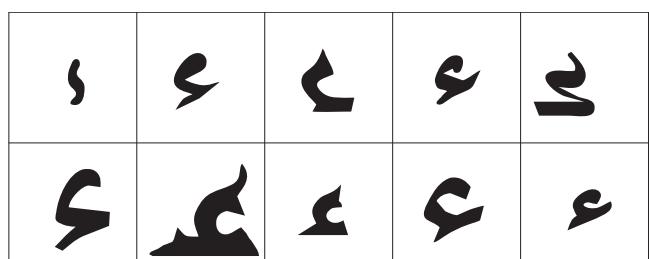


Fig. 3. The Diacritic “Hamza” written in 10 different fonts

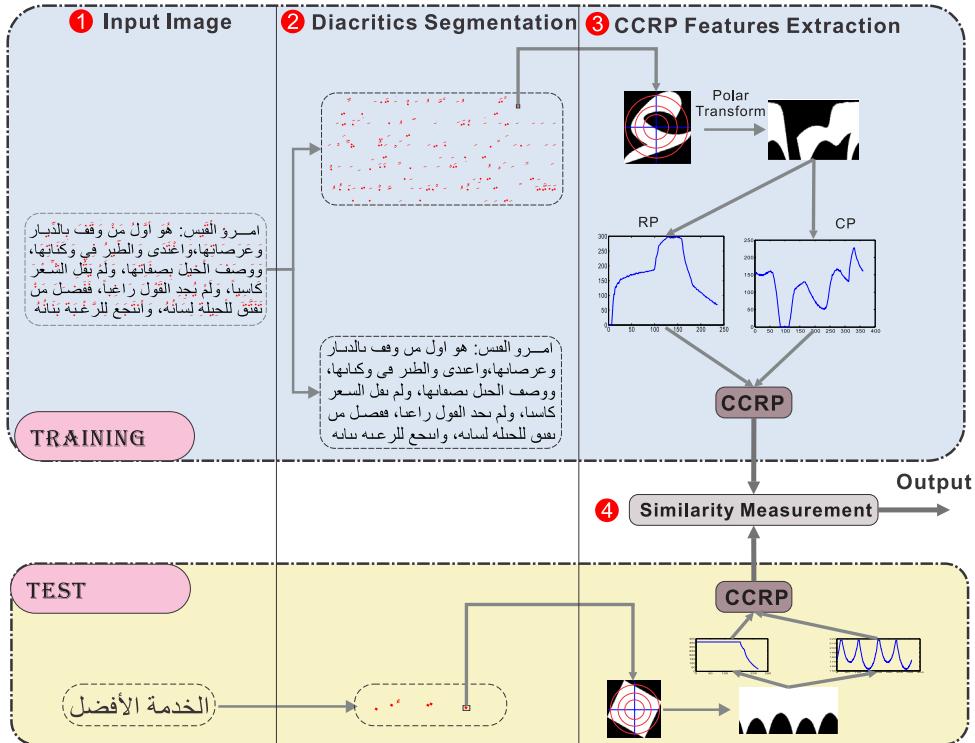


Fig. 4. The proposed system architecture.

3. Proposed method

The proposed method for Arabic font recognition is illustrated in Fig. 4 for training and testing. Initially, for training the input image is split into its original two parts, one part contains letters and the other part contains diacritics. Features are extracted then from each diacritic to generate a feature vector which is stored in the database. Since the representative features for each font are determined by the features extracted from training diacritics, a comparison is needed between the test diacritic features and the training diacritic features in order to estimate the similarity. In our approach, a composite of central and ring projection (CCRP) features proposed in [24] are extracted from each diacritic, and the normalized cross correlation is used to measure the similarity between the testing and training diacritics.

3.1. Preprocessing

The preprocessing mainly focuses on diacritics segmentation, all other preprocessing tasks like noise removing, orientation correction and text localization are assumed to be done by OCR or ADP system.

3.1.1. Flood fill-based diacritic segmentation

Flood-fill (sometimes called seed-fill) is a simple algorithm used to determine the region connected to a given pixel (the seed) and fill it with specific color. In Arabic text, the letters in a subword are connected with each other forming one region, while the diacritics associated with each subword are located above or below it. Each diacritic can be seen as an independent region, because they are not connected to each other nor to the text body. Consequently, if we apply a flood-fill operation on all text regions of an input document image; the result will be the same image containing only the diacritics. In our case, we need to determine

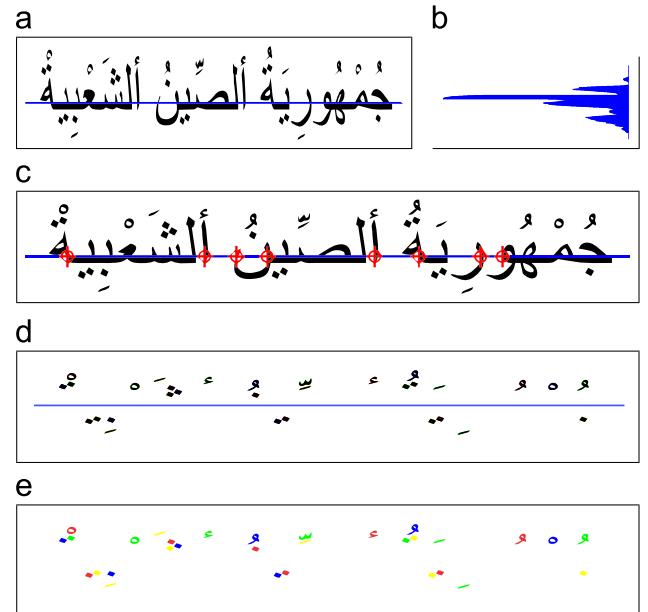


Fig. 5. Flood-fill based diacritics segmentation: (a) the input image, (b) the horizontal projection profile for the image in "a", (c) flood-fill seeds detection, (d) the image "a" after applying the flood-fill operations, and (e) the final result.

the seed of the flood-fill operation which should be a pixel located in the text region as follows.

Using the text horizontal projection profile, the maximum value in this profile indicates the location of the baseline as shown in Fig. 5(a) and (b). Let G represent the input document binary image, then

$$H_{profile}(i) = \sum_{j=1}^N G(i,j) \quad (1)$$



Fig. 6. Diacritics segmentation using flood-fill algorithm: (a) the input image with baselines (blue lines) and the flood-fill seeds shown in red circles while the actually used seeds are filled with yellow and (b) the resultant image. (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this article.)

where N is the number of columns in the image. The index of the baseline $I_{baseline}$ equals to the index of the maximum value in the profile:

$$I_{baseline} = \arg \max_i (H_{profile}(i)) \quad (2)$$

The pixels located on the baseline intersect only with the main text. Along the baseline, starting with $j=0$, if any pixel $p(j)$ with a value equal to 1 followed by a pixel $p(j+1)$ with a value equal to 0, then, $p(j+1)$ is a seed for a flood-fill operation as shown in Fig. 5(c). Some text regions may have more than one seed, but only one flood-fill operation will be performed because the value of any extra seeds will be set to 1 in the first flood-fill operation executed on this region. After clearing the whole text, we can easily detect diacritics (Fig. 5(d)) and obtain their labels (Fig. 5(e)). Providing that the input image contains more than one text line, we can detect the lines boundaries using the horizontal projection profile. While the maximum values in the profile indicate the baselines, the minimum value between two adjacent baselines indicates the boundary between them. Thus, we can easily segment each text line and apply the flood-fill segmentation algorithm which was mentioned above. Fig. 6 shows an example from our page level test dataset.

3.1.2. Clustering-based diacritics segmentation

When the text lines are not well separated or the text direction is not horizontal, the flood-fill based diacritic segmentation algorithm may result in a poor segmentation. Thus, to segment the diacritics, we use a K-means cluster-based method. A major advantage of the k-means algorithm is its computational simplicity. Furthermore, its conceptual simplicity has been a source of

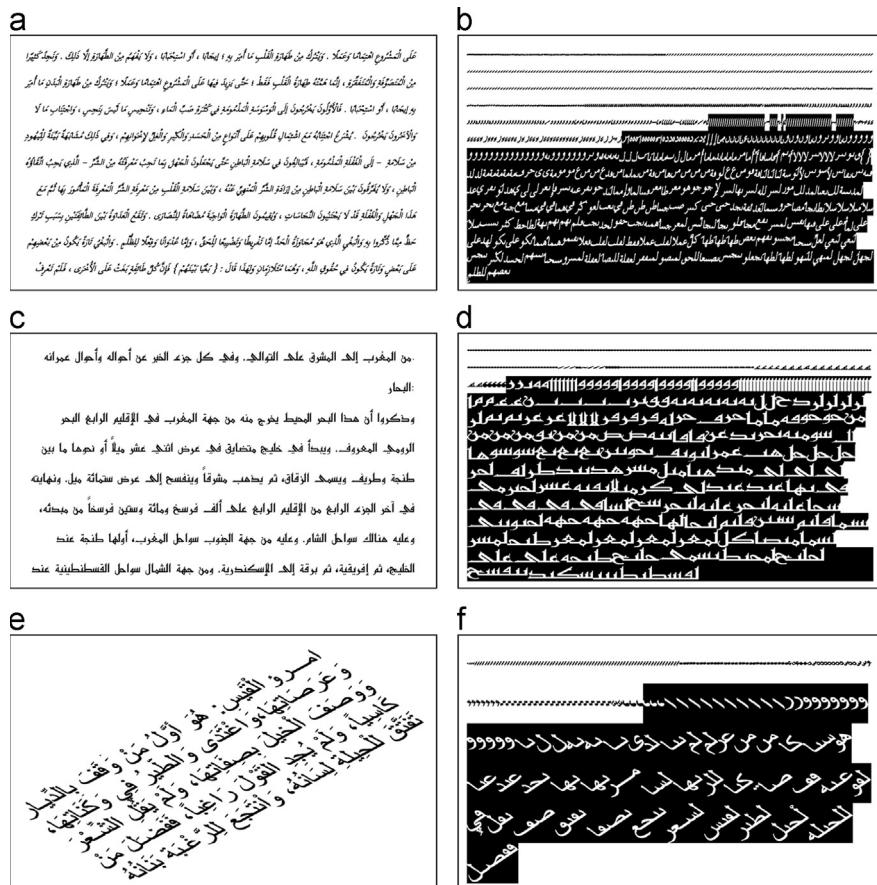


Fig. 7. Diacritics segmentation using clustering-based algorithm: (a) an example from our test dataset "Traditional Arabic font", (b) image "a" after sorting its components "text and diacritics" with respect to region size , the diacritics background shown in white and text background shown in black, (c) another example written in modern style where the vowels diacritics are not included "Kufi font", (d) the output of the image "c", (e) input image rotated 40°, (f) the output of the image "e".

inspiration to many researchers and a lot of methods have been proposed to improve its performance [25], especially the clustering initialization and the instance order. In our implementation, we only run one iteration of K-means followed by a validation process.

Given an input image G , let cluster C_1 represents the text, cluster C_2 represents the diacritics and vector $S = [S_1, S_2, \dots, S_n]$ represents the size of all regions R , where n is the total number of regions, then

Step1 : Initialize the clusters representatives as

$$\mu_1 = \text{mean}(S), \quad \mu_2 = \mu_1/4.$$

Step2 :

```

for  $i=1$  to  $n$ 
    if  $d(S_i, \mu_2) < d(S_i, \mu_1)$  then
         $C_2 \leftarrow R_i$ 
    end{if}
end{for}
```

where d is the Euclidian distance.

3.1.3. Diacritics validation

All diacritics can be segmented successfully, but the letters written in an isolated form may be misclassified as diacritics especially the letter “Alef” . Therefore, the diacritic D will be labeled as an isolated letter if the height m and width n of D satisfy $m > 2n$.

As shown in Fig. 7, if we reorder all regions (text and diacritics) in the image G according to the region size S , we can see that the diacritics are grouped together especially when G is written with the basic diacritics (which is the most widely used writing style) as shown in Fig. 7(d). Some diacritics may be classified as normal characters. This will not effect the recognition accuracy, because it only reduces the number of diacritics extracted from the input image which can be seen as an image contains only the successfully extracted diacritics.

3.2. Features extraction and classification

Rotation invariant is the most critical problem with text image processing, because the scanning procedure always remain difficult to be controlled perfectly. To avoid this issue, we use a composite of central and ring projection (CCRP) proposed in [24] for feature extraction. The CCRP is very easy to be implemented, and its very appropriate to be applied for diacritics, because each diacritic compose of one region only. Ring projection (RP) and central projection (CP) are methods for transforming 2-D patterns into 1-D pattern. RP reduces the dimensionality by performing projection along circles with different radius [26], while CP reduces the dimensionality by performing projection along lines with different polar angles [27]. As shown in Fig. 8, to calculate the CP and RP, we first find the centroid (x_0, y_0) of the diacritic region D then translate the origin of the diacritic image I to this centroid. Let

$$R = \max_{(x,y) \in D} \sqrt{(x-x_0)^2 + (y-y_0)^2} \quad (3)$$

be the longest distance from (x_0, y_0) to a point (x, y) on D . The Cartesian coordinate system then should be converted to polar coordinate system. The conversion is based on the following relations:

$$x = \gamma \cos \theta, \quad y = \gamma \sin \theta \quad (4)$$

hence,

$$I(x, y) = I(\gamma \cos \theta, \gamma \sin \theta)$$

where, $\gamma \in [0, R]$, $\theta \in [0, 2\pi]$.

The RP and CP methods are performed by computing the following integrals, respectively:

$$f(\gamma) = \int_0^{2\pi} I(\gamma \cos \theta, \gamma \sin \theta) d\theta \quad (5)$$

$$f(\theta) = \int_0^R I(\gamma \cos \theta, \gamma \sin \theta) d\gamma \quad (6)$$

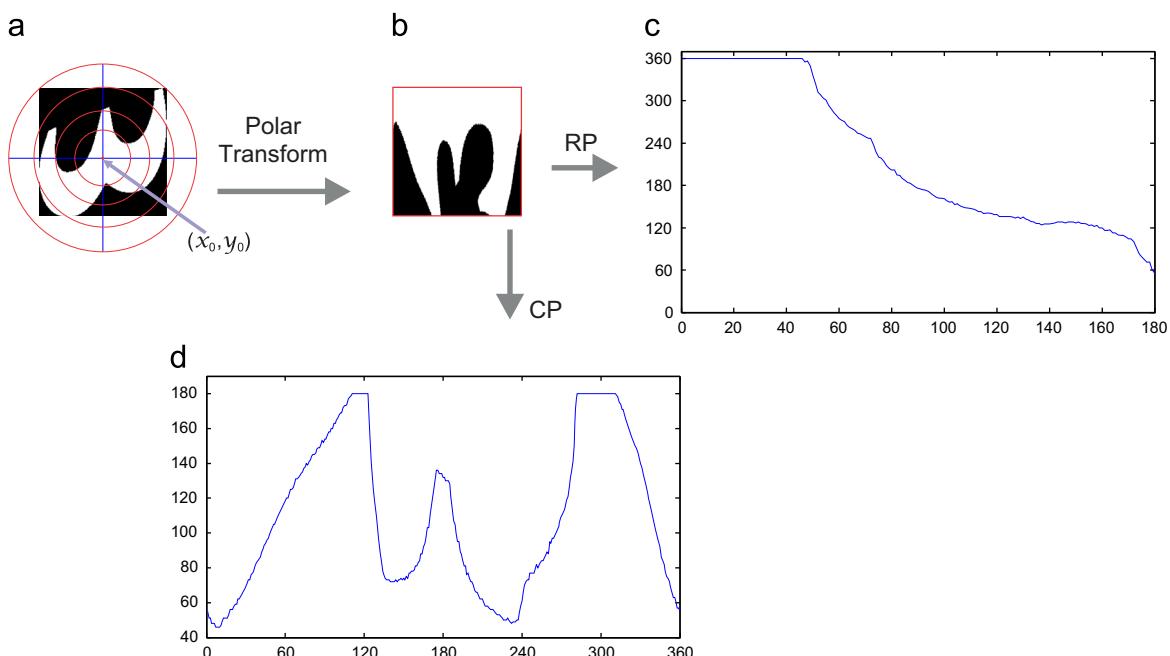


Fig. 8. Features extraction, (a) the input diacritic image I , (b) the polar transform of I , (c) the ring projection curve of I and (d) the central projection curve of I .

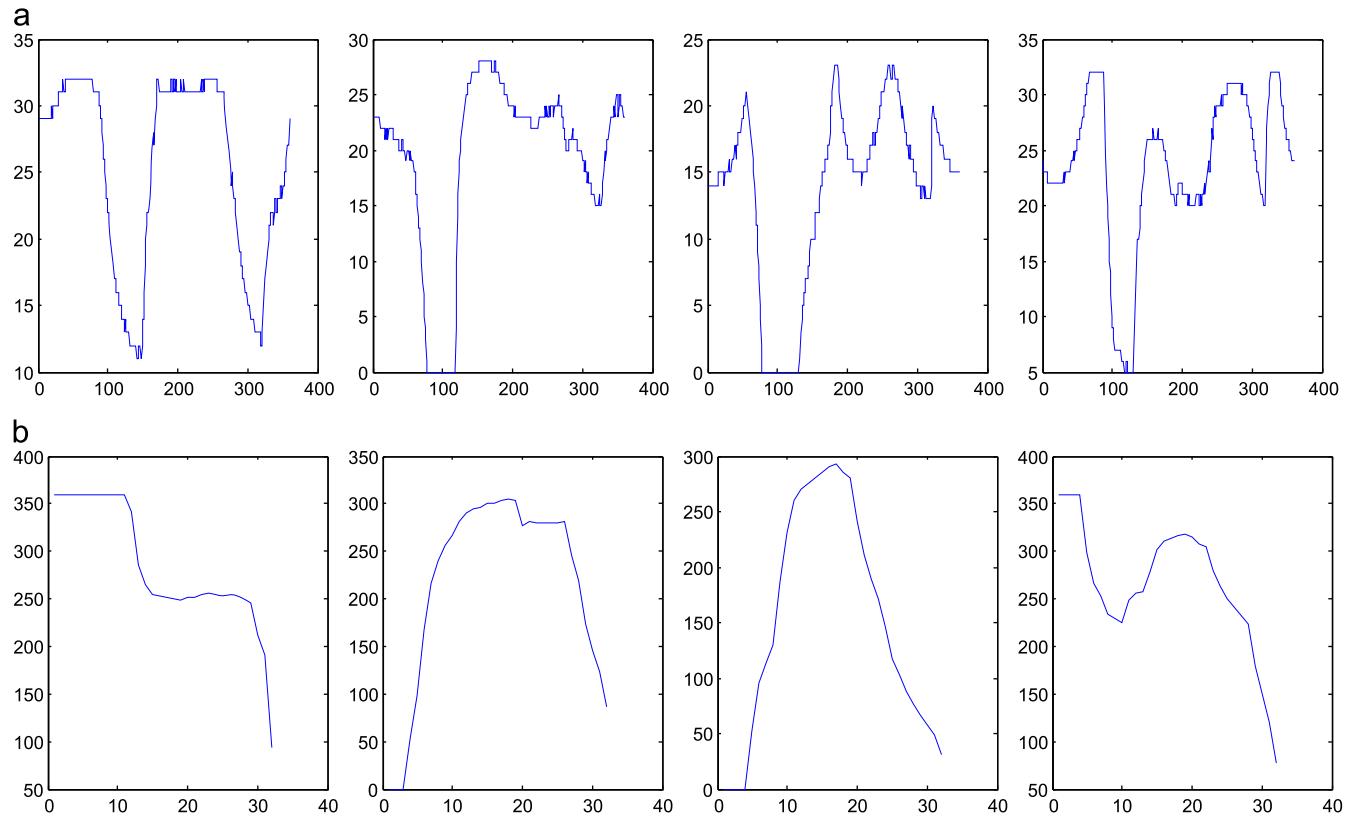


Fig. 9. CP and RP feature curves for the first four diacritics from left shown in Fig. 3. (a) The CP curves and (b) the RP curves.

Since the diacritic image is in a discrete format, Eqs. (5) and (6) can be approximated as follows:

$$CP : f(\gamma_l) = \sum_{\theta=0}^{2\pi} I(\gamma_l \cos \theta, \gamma_l \sin \theta) \quad (7)$$

$$RP : f(\theta_k) = \sum_{\gamma=0}^R I(\gamma \cos \theta_k, \gamma \sin \theta_k) \quad (8)$$

where, $\gamma_l \in [0, R]$, $l = 0, 1, \dots, M$ and $\theta_k \in [0, 2\pi]$, $k = 0, 1, \dots, N$. An illustration of CP and RP is shown in Fig. 9 for 4 different diacritics.

In this implementation, we group each class of diacritics in each font, then the mean CCRP for each diacritic class is calculated. The result is 110 CCRP features ($10 \text{ fonts} \times 11 \text{ diacritic class}$).

For classification we use the normalized cross correlation. Given two sequences a and b , the normalized cross correlation is

$$NCC_{(a,b)}(l) = \frac{\sum_k a_k b_{k-l}}{\sqrt{\sum_k a_k^2 \sum_k b_k^2}} \quad (9)$$

For test diacritic X and training diacritic class Y , let X_r and Y_r represent the curves obtained by RP, X_c and Y_c represent the curves obtained by CP for diacritics X and Y respectively, then the similarity is defined as follows:

$$Sim(X, Y) = \lambda_1 \cdot NCC_{(X_r, Y_r)} + \lambda_2 \cdot NCC_{(X_c, Y_c)} \quad (10)$$

where $\lambda_1 + \lambda_2 = 1$.

For M fonts, the diacritic X is of font K when $Sim(X, Y)$ is maximum at K :

$$X^K = \arg \max_K Sim(X, Y_k), \quad 1 \leq K \leq M \quad (11)$$

where Y is the training diacritic class which has the maximum similarity with X compared to all training diacritic classes belonging to the font K .

When the image G is a text of more than one character, it will contain more than one diacritic, say N diacritics. In this situation, a max voting strategy can be applied. If the diacritic X is classified as font K , then the vote V_K is added by one. The font with the maximum vote is selected as the font type for G :

$$G^K = \arg \max_K V_K, \quad 1 \leq K \leq M \quad (12)$$

4. Experimental results

The proposed font recognition method is evaluated by using our database. The experiments were conducted through eight experiments in which the first two experiments are for font recognition using both the page level dataset and the line level dataset. Experiment 3 is for determining the minimum number of diacritics required for font recognition. In experiment 4 we address the influence of the diacritic class in the recognition result. Experiment 5 is for font style recognition, and experiment 6 is for recognizing the font of an Arabic text written with non-Arabic text. In experiment 7 we show how our method works with text written in different orientation. Finally, an evaluation of our approach with real scanned data are presented in experiment 8. Further, all simulations are performed using MATLAB 7.0.

4.1. The database

To the best of our knowledge, there are over 2000 Arabic fonts, but only a few are mostly used. The database used here includes the training and test datasets. The training dataset consists of a paragraph which contains all kinds of diacritics that may appear in an Arabic text written by 10 different fonts: Jaridah, Farsi, Kufi, Tholoth, Badr, Traditional Arabic, Andalos, Hijaz, Simplified Arabic

and Naskh font, as shown in Fig. 10. Each font is written in four sizes: 10, 11, 12 and 13 with four different font styles: normal, bold, italic and bold-italic, to end up with 160 training samples (10 fonts \times 4 sizes \times 4 styles). Also the same 160 texts were rewritten using only basic diacritics. This makes the final training samples up to 320. We extract the diacritics from the training samples and store them separately. After that, we run a test on the training diacritics which belong to the same font class and check if there exist 100% similar diacritics and exclude the duplicated diacritics from the database to reduce the training samples.

For test data, two different datasets are created: page level dataset where the image contains several text lines; and line level dataset where the image contains only one text line. Page level dataset was created by converting the pages in an Arabic eBook to digital images. The book has about 238 pages with 8 lines in each.

page. 48 pages with full diacritics and another 48 pages with basic diacritics are selected. After that the font of these pages is changed using word processor software to the selected fonts, sizes and styles to end up with 960 images, 96 images for each font.

Line level dataset is the same as the page level dataset except that each line is segmented individually. Since each page has 8 lines, consequently the line dataset will have 7680 lines image (960 pages \times 8 lines).

4.2. Page level font recognition

The aim of this experiment is to recognize the font using the diacritics extracted from the whole page with the assumption that the page is written with only one font. By using the proposed method, the average font recognition rate reaches 98.73% and the recognition rate for most font types reaches 100%. The lowest recognition rate is for the *Hijaz* font which is equal to 93.54%. The result is shown in Table 3.

4.3. Line level font recognition

We run the second experiment on the line level dataset for line font recognition with the assumption that each line was written with only one font. The average recognition rate reaches 95.42%, the highest recognition error rate occurs between the simplified Arabic font and traditional Arabic font because of the similarity between these two fonts. For OCR system, the same technique should work with both fonts. The result is shown in [Table 4](#).

While the font recognition is just a supplementary step for OCR; our approach provides a robust solution for this problem such as the book cover shown in Fig. 11(a), which contains five different fonts in just one single page. Compared to the state-of-the-art methods, the font of the book cover can be detected very easily using our approach by selecting to recognize the font in line level.

4.4. Number of diacritics influence

The difference between the above two experiments lies only in the number of the diacritics N involved in the recognition process. The recognition rate of both experiments is close, which may raise a question about the minimum number of diacritics needed to recognize a font. To answer this question, we run the third experiment on the page level dataset by selecting 1600 diacritics from each font then dividing them into 20 groups, 80 diacritics for each group. The experiment is run 80 times in each group, and each time the number of diacritics N increases by 1. We then calculate the average minimum number of diacritics needed for best font recognition rate.

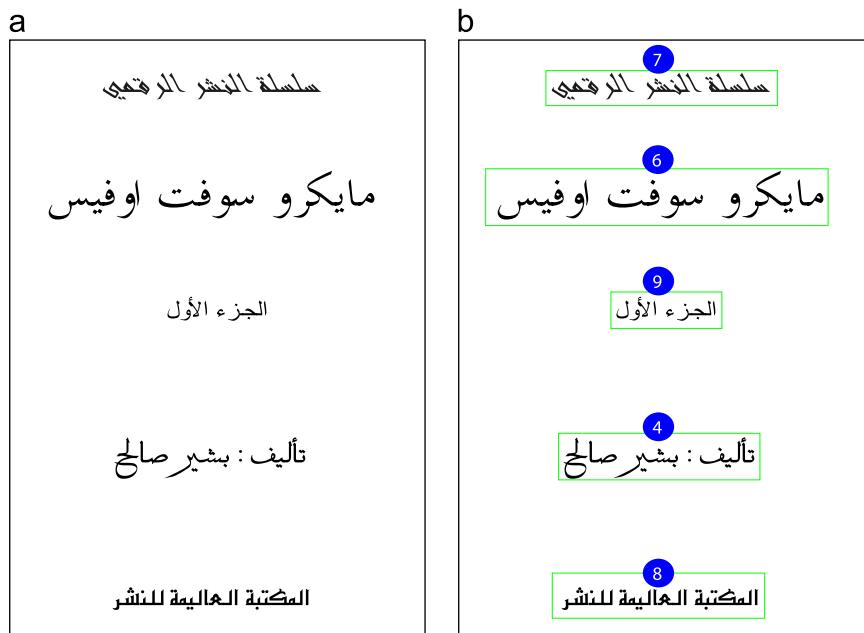
Fig. 10. Examples of the training text shown in different fonts and sizes. From top to bottom: Jaridah, Farisy, Kufi, Tholoth, Badr, Traditional Arabic, Andalos, Hijaz, Simplified Arabic and Naskh.

Table 3
Confusion matrix of 10 Arabic fonts in page level.

Table 4

Confusion matrix of 10 Arabic fonts in line level.

Font	Andalos	Badr	Farisy	Hijaz	Jaridah	Kufi	Naskh	Simplified Arabic	Tholoth	Traditional Arabic
Andalos	95.68	0.13	1.08	0.13	0.27	0.54	1.21	0	0.54	0.40
Badr	1.24	96.19	0	0.51	0	0	0	0.21	0.31	1.54
Farisy	0.14	0	88.48	0	0.41	0	1.22	0.41	9.35	0
Hijaz	0.14	0.84	0.56	93.31	0	0.28	0	0.84	4.04	0
Jaridah	0.54	0.14	0.41	0.00	97.97	0.14	0	0.14	0.68	0
Kufi	0.41	0	0	0.14	0	99.18	0	0	0.14	0.14
Naskh	0.14	0.14	0.28	0	0	0	99.15	0	0.28	0
Simplified Arabic	0	1.00	0.45	0	0.11	0	1.34	92.19	0.11	4.80
Tholoth	0	0	0	0	0	0	0.31	0	99.69	0
Traditional Arabic	0	0	0	0	0	0	0	7.62	0	92.38

**Fig. 11.** (a) A book cover where text lines are written in different fonts which can be detected using our approach by selecting to recognize the font in line level and (b) the result of applying line level font recognition, the lines fonts from top are Andalos, Traditional Arabic, Simplified Arabic, Tholoth and Hijaz.**Table 5**

The minimum number of diacritics needed for font recognition.

Font	Andalos	Badr	Farisy	Hijaz	Jaridah	Kufi	Naskh	Simplified Arabic	Tholoth	Traditional Arabic
Minimum diacritics number	5	3	7	5	4	1	1	5	1	3

Table 6

Font recognition rate given only one diacritic class.

Font\diacritic	ـ	ـ	ـ	ـ	ـ	ـ	ـ	ـ	ـ	ـ
Andalos	88.96	88.96	100	95.75	99	100	100	100	98.97	100
Badr	100	100	97.3	100	87.63	97.3	100	97.3	96.86	100
Farisy	100	100	100	100	100	100	100	100	92.43	100
Hijaz	93.33	93.33	100	100	100	100	100	100	91.26	100
Jaridah	100	100	100	98.8	98.8	100	100	100	92.86	100
Kufi	100	100	98	100	100	100	97.62	98	91.12	100
Naskh	100	96	100	95.63	100	88.67	95.42	100	89.43	100
Simplified Arabic	100	100	98	100	100	98.04	100	92.36	92.59	100
Tholoth	100	100	85.42	97.8	100	100	100	100	91.96	92.22
Traditional Arabic	100	100	100	94	100	100	100	98	94.79	89.52

As shown in Table 5, the minimum number of the diacritics needed to recognize the font of a given text is different from one font to another. Some fonts, such as *Kufi* needs only one diacritic, while *Farisy* font needs seven diacritics as the highest number.

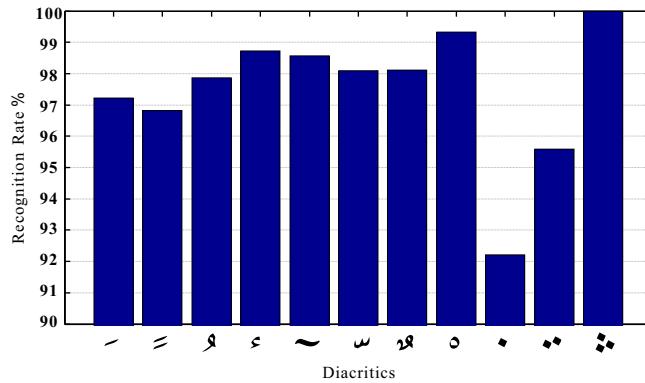


Fig. 12. Average font recognition rate given only one diacritic class.

Table 7
Confusion matrix of 10 Arabic fonts considering the style in page level dataset.

Font name	Font style	Normal	Italic	Bold	Bold-italic
Andalos	Normal	100			
	Italic		100		
	Bold	42		58	
	Bold-italic	12			88
Badr	Normal	100			
	Italic	13	87		
	Bold	17			83
	Bold-italic			4.2	95.8
Farisy	Normal	100			
	Italic		100		
	Bold	4	25	71	
	Bold-italic		92		8
Hijaz	Normal	100			
	Italic		100		
	Bold	42		58	
	Bold-italic	4			96
Jaridah	Normal	63	37	37	
	Italic		75		25
	Bold	8		92	
	Bold-italic		25		75
Kufi	Normal	67		33	
	Italic		54		46
	Bold			100	
	Bold-italic				100
Naskh	Normal	100			
	Italic		100		
	Bold			100	
	Bold-italic				100
Simplified-Arabic	Normal	96			4
	Italic		88		12
	Bold			100	
	Bold-italic			12	88
Tholoth	Normal	100			
	Italic		100		
	Bold	17	20	63	
	Bold-italic	8		4	88
Traditional-Arabic	Normal	58		42	
	Italic		100		
	Bold			100	
	Bold-italic	4	50	21	25

4.5. Diacritic class influence

The discrimination between diacritics in different fonts differs from one diacritic to another, some are well distinguishable while others are very similar. To address this issue we run an experiment for recognizing the font given only one diacritic class. From the page level dataset, we manually select 1000 diacritics from each diacritic class in each font type, then test them against the training diacritics. The result is shown in Table 6 and the average recognition rate is shown in Fig. 12. From the result we can see that the “three-dots” diacritic has the highest discrimination rate equal to 100% for all fonts, while the “one-dot” diacritic is the lowest.

4.6. Font style recognition

In this section, an experiment for recognizing the style of the font is carried out on the page level dataset. The training dataset is recreated to include the font style in the diacritic class by grouping diacritics according to the diacritic class and font style in each font. The result is 440 training class(11 diacritic class \times 10 fonts \times 4 font styles) and the average recognition rate comes to 83.02% as shown in Table 7. Naskh font has the highest discrimination rate as 100% for all four styles, while the Farsi font gets the lowest recognition rate equal 100%, 100%, 71% and 8% for normal, italic, bold and bold-italic style respectively.

4.7. Multi-language document

Another practical experiment is carried out in this section as shown in Fig. 13(a) where Arabic text is written along with another language text both in the same document. The Latin language (as in this example) does not contain any diacritics, and it does not have any effect on the extraction of the diacritics, which allows us not only to easily recognize the Arabic font but also to localize Arabic text as shown in Fig. 13(b). In this figure, the location of the Arabic text is simply the location of the diacritics.

4.8. Text orientation

Recognizing the text font type written in non-horizontal orientation is always a challenging task for many OCR systems. This is mainly due to the fact that it cannot detect the text location or it cannot recognize the font type of the text, such as maps and figures where the text may appear with skewness. But in our

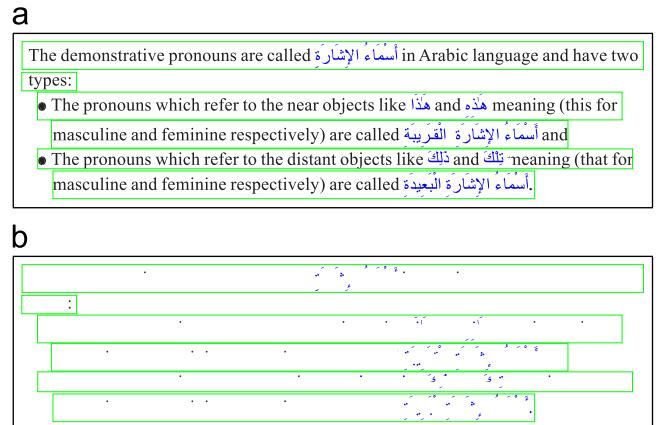


Fig. 13. Multi-language document (Arabic and English), the English text does not contain any diacritics except for letters “i” and “j” where it can be detected by calculating the ratio between the width of the diacritic and the letter (subword) associated with it, the other diacritics are belong to the Arabic text which can be used for font recognition and text localization.

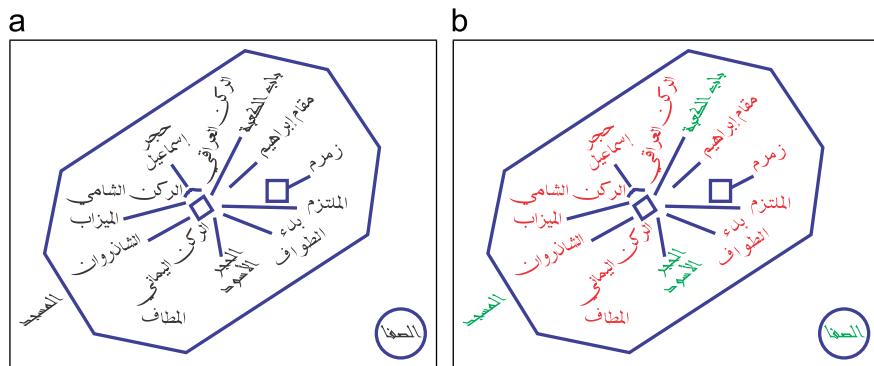


Fig. 14. Text with different orientation: (a) an illustration map of the "Holy mosque of Makkah" and (b) the result image indicates two fonts diacritics, Naskh font shown in red and Andalos font shown in green. (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this article.)

Table 8

Font recognition with 200 scanned images (20 images per font type). our system recognize the font of 191 images with top 1 result, and all the 200 images with top 3 result.

Font	Andalos	Badr	Farsi	Hijaz	Jaridah	Kufi	Naskh	Simplified Arabic	Tholoth	Traditional Arabic
Top 1 RR	20	17	19	20	20	18	20	18	20	19
Top 2 RR	20	19	20	20	20	19	20	20	20	20
Top 3 RR	20	20	20	20	20	20	20	20	20	20

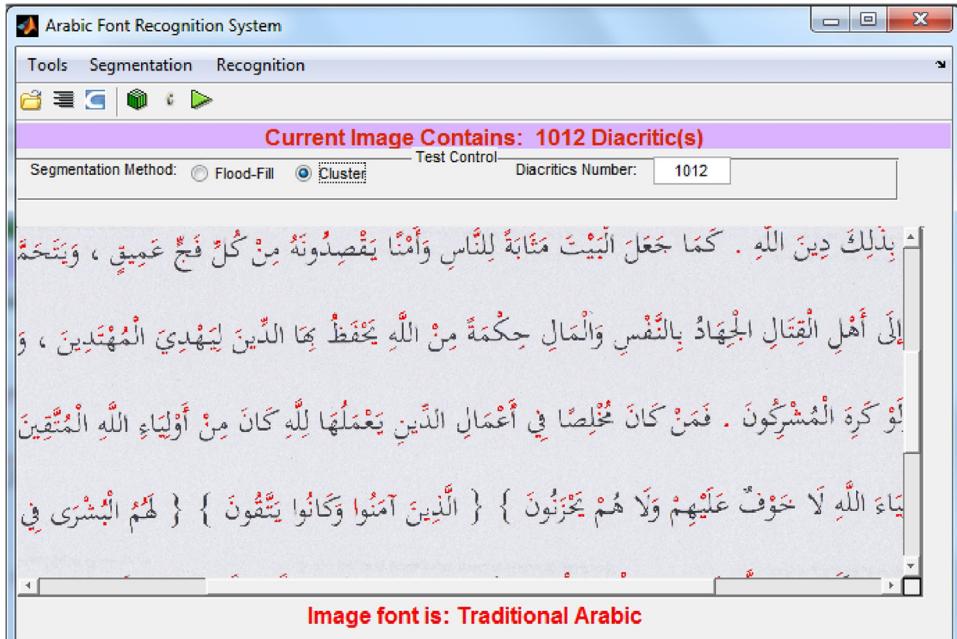


Fig. 15. Snapshot of the proposed Arabic font recognition system. The text image shown contains more than one thousand diacritics (labeled with red color). (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this article.)

approach, by segmenting the diacritics using the clustering-based algorithm, the diacritics can be located regardless of its orientation. Also, by using the rotation invariant CCRP features, we can recognize the font type of the text in any orientation. Fig. 14 shows a map used in one of our experiments, it contains about 15 labels indicating the names of places which are written in different orientations using two font families, the recognition result is shown in Fig. 14(b).

4.9. Real data evaluation

This last experiment is to test our proposed method with real data. We scanned 200 pages from our page level dataset written

by 10 different fonts, 20 pages for each font. The scan resolution was equal to 200 dpi. The test result is shown in Table 8, and a snapshot of our system with a real data image is shown in Fig. 15. From this figure we can see that almost all diacritics have been segmented correctly.

In addition, Artificial distortion has been added to some randomly selected images, and MATLAB built-in median filter and thresholding functions are used for this experiment. Our system can recognizes the font class of these images as long as the size of the noise objects is small enough to be filtered out, otherwise it will be difficult to be distinguished from the diacritics. Also, we test a few colored images with different resolutions collected from the internet as shown in Fig. 16.



Fig. 16. Examples of testing our method with real images: (a) shows the noisy image written with Jaridah font, (b) shows the result of diacritics segmentation for the noisy image, (c) shows a colored low resolution (96 dpi) image written with Tholoth font, and (d) shows the result of diacritics segmentation for image (c). (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this article.)

Table 9

Comparison study with recent proposed works.

Publication	Number of fonts	Average recognition rate	Required time per sample (ms)
Ben Moussa et al. [16]	10	96.2% using KNN 98% using RBF	N/A
Khosravi et al. [15]	10	94.16%	3.78 ms
Slimane et al. [10]	10	99%	N/A
Our	10	98.73%	1.3 ms

5. Conclusion

This paper presents a novel method for Arabic font recognition. The proposed method is based on segmenting the diacritics from the input text image, then extracting CCRP features from these diacritics. The main contribution of the proposed method is that it only uses diacritics for font recognition. We made some experiments with the available Arabic OCR products; we noticed that, the major errors come from the letters associated with diacritics. We believe that separating diacritics from letters will highly improve the accuracy of the result when processing any Arabic text image. This in turn enabled us to get a promising result in this work for font recognition and we will continue using the same approach for our future OCR. The main disadvantage of this approach is noise and broken strokes which unfortunately introduce invalid diacritics and, as a result, it might decrease the recognition rate. But, As long as the size of noise objects are smaller than the diacritics size, we can still successfully segment all diacritics in the image, but if the size of the noise objects grow to be close to diacritics size, it will be very difficult to distinguish between them, but this could be controlled during the acquisition process. A few works have been proposed for Arabic font recognition, but each of the proposed method has its own database, Table 9 is a comparison of our method with some of these methods. Accordingly, our method is the fastest one (the time indicated in the table refers to the time needed to recognize the font of a 512×512 text image), which makes it the best candidate to work with an OCR system where the time needed to identify the font type should be as short as possible.

Conflict of interest

None declared.

Acknowledgments

We want to thank the constructive comments and suggestions from the anonymous reviewers. This research was supported partially by the National Natural Science Foundation of China (Grant no. 61272203), the NSFC grant under 61272366, and the Faculty Research Grant of Hong Kong Baptist University with the Project Code: FRG2/11-12/067 and FRG2/12-13/082, the International Scientific and Technological Cooperation Project (Grant no. 2011DFA12180), National Science & Technology R&D Program (Grant no. 2012BAK31G01, Grant no. 2012BAK02B06), the Ph.D. Programs Foundation of Ministry of Education of China (Grant no. 20110142110060), and the Key Science and Technology Program of Wuhan (Grant no. 201210121021).

References

- [1] G. Nagy, Twenty years of document image analysis in pam, IEEE Transactions on Pattern Analysis and Machine Intelligence 22 (1) (2000) 38–62.
- [2] B.B. Chaudhuri, U. Garain, Automatic detection of italic, bold and all-capital words in document images, in: 14th International Conference on Pattern Recognition (ICPR), Washington, DC, USA, 1998, pp. 610–612.
- [3] A. Zramidini, R. Ingold, Optical font recognition using typographical features, IEEE Transactions on Pattern Analysis and Machine Intelligence 20 (8) (1998) 877–882.
- [4] C.B. Jeong, H.K. Kwag, S.H. Kim, J.S. Kim, S.C. Park, Identification of font styles and typefaces in printed korean documents, in: Digital Libraries—Technology and Management of Indigenous Knowledge for Global Access, vol. 2911, 2003, pp. 666–669.
- [5] Y. Zhu, T. Tan, Y. Wang, Font recognition based on global texture analysis, IEEE Transactions on Pattern Analysis and Machine Intelligence 23 (10) (2001) 1192–1200.
- [6] H. Ma, D. Doermann, Gabor filter based multi-class classifier for scanned document images, in: 7th International Conference on Document Analysis and Recognition (ICDAR), vol. 2, 2003, pp. 968–972.
- [7] C. Avilés-Cruz, R. Rangel-Kuoppa, M. Reyes-Ayala, A. Andrade-Gonzalez, R. Escarela-Perez, High-order statistical texture analysis – font recognition applied, Pattern Recognition Letters 26 (2) (2005) 135–145.

- [8] E. Rashedi, E.N. abadipour, S. Saryazdi, Farsi font recognition using correlation coefficients (in farsi), in: 4th International Conference on Machine Vision and Image Processing, Ferdosi Mashhad, 2007.
- [9] G. Joshi, S. Garg, J. Sivaswamy, A generalised framework for script identification, *International Journal of Document Analysis and Recognition* 10 (2) (2007) 55–68.
- [10] F. Slimane, S. Kanoun, A.M. Alimi, R. Ingold, J. Hennebert, Gaussian mixture models for arabic font recognition, in: 20th International Conference on Pattern Recognition (ICPR), Washington, DC, USA, 2010, pp. 2174–2177.
- [11] F. Slimane, S. Kanoun, J. Hennebert, A.M. Alimi, R. Ingold, A study on font-family and font-size recognition applied to arabic word images at ultra-low resolution, *Pattern Recognition Letters* 34 (2) (2013) 209–218.
- [12] M.B. Imani, M.R. Keyvanpour, R. Azmi, Semi-supervised persian font recognition, *Procedia Computer Science* 3 (2011) 336–342.
- [13] B. Bataineh, S. Abdulla, K. Omar, A statistical global feature extraction method for optical font recognition, in: Intelligent Information and Database Systems, vol. 6591 of Lecture Notes in Computer Science, Springer, Berlin Heidelberg, 2011, pp. 257–267.
- [14] Y. Pourasad, H. Hassibi, A. Ghorbani, Farsi font recognition using holes of letters and horizontal projection profile, in: Innovative Computing Technology, vol. 241 of Communications in Computer and Information Science, 2011, pp. 235–243.
- [15] H. Khosravi, E. Kabir, Farsi font recognition based on Sobel–Roberts features, *Pattern Recognition Letters* 31 (1) (2010) 75–82.
- [16] S. Ben Moussa, A. Zahour, A. Benabdellahid, A.M. Alimi, New features using fractal multi-dimensions for generalized arabic font recognition, *Pattern Recognition Letters* 31 (5) (2010) 361–371.
- [17] I.S.I. Abuhaiba, Arabic font recognition using decision trees built from common words, *Computing and Information Technology* 13 (3) (2005) 211–224.
- [18] I. Abuhaiba, Arabic font recognition based on templates, *The International Arab Journal of Information Technology* 1 (2003) 33–39.
- [19] M. Lutf, X. You, H. Li, Offline arabic handwriting identification using language diacritics, in: 20th International Conference on Pattern Recognition (ICPR), 2010, pp. 1912–1915.
- [20] Encyclopaedia Britannica, 2010. URL <<http://www.britannica.com/>>.
- [21] D.L. O'Leary, Arabic Thought and Its Place in History, Dover Publications, 2003.
- [22] M.C. Bateson, Arabic Language Handbook (Georgetown Classics in Arabic Language and Linguistics), Georgetown University Press, 2003.
- [23] N. Girgis, The Arabic Alphabet, International Book Centre, Inc., 1983.
- [24] R.S. Lan, J.W. Yang, Y. Tang, A composite of central and ring projection, in: International Conference on Wavelet Analysis and Pattern Recognition, 2009, pp. 200–204.
- [25] S. Theodoridis, K. Koutroumbas, Pattern Recognition, 4th edition, Academic Press, 2008.
- [26] Y.Y. Tang, H.D. Cheng, C.Y. Suen, Transformation-ring-projection (TRP) algorithm and its VLSI implementation, *International Journal of Pattern Recognition and Artificial Intelligence* 5 (1–2) (1991) 25–56.
- [27] Y. Tao, E.C.M. Lam, C.S. Huang, Y.Y. Tang, Information distribution of the central projection method for chinese character recognition, *Journal of Information Science and Engineering* 16 (1) (2000) 127–139.

Mohammed Lutf a Ph.D. student in the Department of Electronics and Information Engineering at Huazhong University of Science and Technology, Wuhan, China.. He is an active researcher in the area of Arabic character and handwriting recognition. He received his B.S. degree in Telecommunication Engineering from Dalian Maritime University, Dalian, China, in 2005 and an M.S. degree in Communication and Information Systems from Huazhong University of Science and Technology, Wuhan, China, in 2010.

Xinge You received the B.S. and M.S. degrees in Mathematics from the Hubei University, Wuhan, China, in 1990, and the Ph.D. degree in Computer Science from the Hong Kong Baptist University, Hong Kong, in 2000 and 2004, respectively. He is presently a Professor in the Department of Electronics and Information Engineering, Huazhong University of Science and Technology, Wuhan, China. His current research interests include wavelets and its application, signal and image processing, pattern recognition, machine learning, and computer vision.

Yiuming Cheung received the Ph.D. degree from the Department of Computer Science and Engineering, The Chinese University of Hong Kong in 2000. He joined the Department of Computer Science at Hong Kong Baptist University in 2001, and then became an Associate Professor in 2005. His current research interests are in the fields of machine learning and information security, particularly the topics on clustering analysis, blind source separation, neural networks, nonlinear optimization, watermarking and lip-reading. He is the founding Chairman of IEEE (Hong Kong) Computational Intelligence Chapter. Currently, he is also the Associate Editor of Knowledge and Information Systems, as well as the guest co-editor and editorial board member of the several international journals.

C.L. Philip Chen is currently a Dean and a Chair Professor of the Faculty of Science and Technology, University of Macau. He has been a Professor and the Chair of the Department of Electrical and Computer Engineering, Associate Dean for Research and Graduate Studies of the College of Engineering, University of Texas at San Antonio, Texas. His current research interests include theoretic development in computational intelligence, intelligent systems, cyber-physical systems, robotics and manufacturing automation, networking, diagnosis and prognosis, and life prediction and life-extending control. He is an elected Fellow of IEEE and AAAS. He is the President-Elect and Vice President on Conferences and Meetings of the IEEE Systems, Man and Cybernetics Society (SMCS), where he has been the Vice President of the Technical Activities on Systems Science and Engineering. Dr. Chen is a member of Tau Beta Pi and Eta Kappa Nu honor societies and has been the faculty advisor for Tau Beta Pi Engineering honor society. In addition, he is an ABET (Accreditation Board of Engineering and Technology Education) Program Evaluator for Computer Engineering, Electrical Engineering, and Software Engineering programs.