# Preservation of Historical Document Images using Phase-based Binarization

Sandhya G Kompi[1]
Department of CSE,
AMC Engineering College
Bangalore

Deepa K S[2]
Assistant Professor
Department of CSE, AMC Engineering College
Bangalore

*Abstract*— **In this paper, a binarization model based on the phase of the ancient document images is proposed. Three features are derived from the phase information of an input document image which constitute the core of the binarization model. These three features are the maximum moment of phase congruency covariance, a locally weighted mean phase angle, and a phase preserving denoised image. The proposed binarization model consists of three standard steps: 1) preprocessing; 2) main binarization; and 3) postprocessing. In the preprocessing we try to bring a rough estimate of the binarized image and in main binarization steps, the features used are mainly derived from phase. In the postprocessing step, Gaussian and median filters are considered to remove Gaussian noise and salt and pepper noise, if any.**

*Keywords*— *Historical document binarization, phase-derived features, and document enhancement.*

## I. INTRODUCTION

There are many degraded and historically-important old manuscripts and documents distributed across libraries and archives around the world. Some of these document degradation are of types, such as fading of ink, ink bleed-through, show-through and deterioration of the cellulose structure, among others. Conversion to binary form from grayscale is a common and fundamental step in almost all digitization processes. The quality of all the binarization step highly affects the performance of the subsequent document processing steps.

A key step in all the document image processing workflows is phase based binarization, but this is not a very easy process, which is unfortunate, as its performance takes a significant toll on the quality of optical character recognition results. Many research studies have been carried out to solve the problems that arise in the binarization of old document images characterized by many types of degradation [1]–[19], including faded ink, bleed-through, show-through, uneven illumination, variations in image contrast, and deterioration of the cellulose structure [1], [20]. There are also differences in patterns of hand-written and machine-printed documents, which add to the difficulties associated with the binarization of old document images. To the best of our knowledge, none of the proposed methods can deal with all types of documents and degradation. Fig. 1 shows some of the degraded document images used in this paper. In this paper, a robust and a fast

phase-based binarization method is proposed for the binarization and enhancement of ancient historical documents and manuscripts. The three main steps in the proposed method are: preprocessing, main binarization, and post-processing. The preprocessing step mainly involves image denoising with phase preservation [23], followed by some morphological operations. We incorporate the Canny edge detector [24] and a denoised image to obtain a binarized image in rough form.
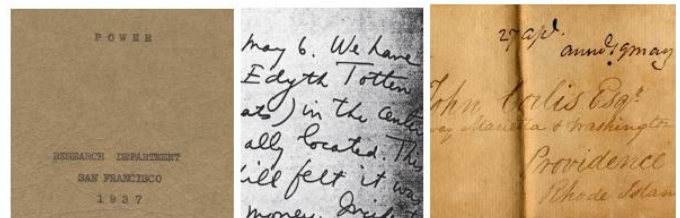


Fig. 1 Three degraded document image samples

Then, we use the phase congruency features [18], [19], [25] for the main binarization step. Phase congruency is widely used in the machine vision and image processing literature [26]–[29]; palmprint verification [26], object detection [27], finger-knuckle-print recognition [28], and biomedical applications [29] are just a few examples of the use of phase congruency as a feature detector. We show that the foreground of ancient documents can be modeled by phase congruency. The previous works [18], [19], and [30] show that phase congruency is a robust way to process historical documents, both handwritten and machine-printed manuscripts.

After completing the three binarization steps on the input images using phase congruency features and a denoised image [23], the enhancement processes are applied. A median filter and a phase congruency feature are used to construct an object exclusion map image. This map is then used to remove unwanted lines and interfering patterns. The effect of each step on the binarized output image is discussed in each associated section.

## II. LITERATURE REVIEW

In this section, we briefly describe some selected binarization methods. Gatos et al. [5] propose an adaptive binarization method based on low-pass filtering, foreground estimation, background surface computation, and a combination of these. In [6], an initial binary map is obtained using the multi-scale

Sauvola's method [1], and then statistical methods are used torestore the missed strokes and sub-strokes. In [8], Valizadeh et al. map input images into a two-dimensional feature space in which the foreground and background regions can be distinguished. Then, they partition this feature space into several small regions, which are classified into text and background based on the results of applying Niblack's method [31]. Lu et al. [9] propose a binarization method based mainly on background estimation and stroke width estimation. First, the background of the document is estimated by means of a one-dimensional iterative Gaussian smoothing procedure. Then, for accurate binarization of strokes and sub-strokes, an L1-norm gradient image is used. This method placed 1st of 43 algorithms submitted to the DIBCO'09 competition [21]. Su et al. [10] use local maximum and minimum to build a local contrast image. Then, a sliding window is applied across that image to determine local thresholds. A version of this method shared 1st place with another method, out of 17 algorithms entered in the H-DIBCO'10 contest [22]. In [2], a local contrast image is combined with a Canny edge map to produce a more robust feature map. This method performs better than those in [9] and [10]. Farrahi Moghaddam et al. [1] propose a multi-scale binarization method in which the input document is binarized several times using different scales. Then, these output images are combined to form the final output image. This method uses different parameters for Sauvola's method to produce output images of the same size, but at different scales. In contrast, Lazzara and Gerard [32] propose a multi-scale Sauvola's method which binarizes different scales of the input image with the same binarization parameters. Then, binary images with different scales are combined in some way to produce the final results. Combination methods have also attracted a great deal of interest, and provided promising results. The goal of combining existing methods is to improve the output based on assumption that different methods complement one another. In [11], several of these methods are combined based on a vote on the outputs of each. In [7], a combination of global and local adaptive binarization methods applied on an inpianted image is used to binarize handwritten document images. The results show that this method performs extremely well; however, it is limited to binarizing handwritten document images only. Learning-based methods have also been proposed in recent years. These methods are an attempt to improve the outputs of other binarization methods based on a feature map [12]–[14], or by determining the optimal parameters of binarization methods for each image [15], [16]. In [12] and [14], a self-training document binarization method is proposed. The input pixels, depending on the binarization method(s) used, are divided into three categories: foreground, background, and uncertain, based on a priori knowledge about the behavior of every method used. Then, foreground and background pixels are clustered into different classes using the k-means algorithm or the random Markov field [12], [14]. Finally, uncertain pixels are classified with the label of their nearest neighboring cluster. The features used for the final decision are pixel intensity and local image contrast. In [13], another combined method based on a modified contrast feature is proposed. Lelore and Bouchara [33] also classify pixels into three categories using a coarse thresholding method, where uncertain pixels are classified

based on super resolution of likelihood of foreground. Howe [17] proposes a method to optimize the global energy function based on a Laplacian image. In this method, a set of training images is used for optimization. In [15], Howe improved this method by tuning two key parameters for each image. In [16], a learning framework is proposed to automatically determine the optimal parameters of any binarization method for each document image. After extracting the features and determining the optimal parameters, the relation between the features and the optimal parameters is learned.

## III.   FEATURES DERIVED FROM PHASE

We use three phase-derived feature maps of the input document image in this paper: two phase congruency feature maps and a denoised image. The details are provided below.

### A. Phase Congruency Features

In [34], it is shown that the phase information of an image outweighs its magnitude information. This implicitly means that phase information is the most important feature of images. In this section, two phase congruency-based feature maps extracted from input images are discussed. These feature maps are based on the Kovesi's phase congruency model [25]. Based on the experiments, Kovesi's method worked better within our proposed binarization method. Let Me and Mo denote the even symmetric and odd symmetric log-Gabor wavelets at a scale $\rho$,

$$\left[e_\rho(x), o_\rho(x)\right] = \left[f(x) * M_\rho^e, f(x) * M_\rho^o\right].$$

where values $e\rho(x)$ and $o\rho(x)$ are real and imaginary parts of a complex-valued wavelet response. The local phase $\varphi\rho(x)$ and the local amplitude $A\rho(x)$ of the transform at a given wavelet scale $\rho$ are:

$$\phi_\rho(x) = arctan2\left(o_\rho(x), e_\rho(x)\right)$$

$$A_\rho(x) = \sqrt{e_\rho(x)^2 + o_\rho(x)^2}.$$

The phase deviation function of two dimensions is presented by considering both the scale $(\rho)$ and the orientation $(r)$ indices of the wavelet coefficients:

$$\Delta\Phi_{\rho r}(x) = \cos\left(\phi_{\rho r}(x) - \overline{\phi}_r(x)\right) - \left|\sin\left(\phi_{\rho r}(x) - \overline{\phi}_r(x)\right)\right|.$$

The two-dimensional phase congruency is calculated by:

$$PC_{2D,r}(x) = \frac{\sum_\rho W_r(x)\lfloor A_{\rho r}(x)\Delta\Phi_{\rho r}(x) - T_r\rfloor}{\sum_\rho A_{\rho r}(x)}$$

Where Tr is the threshold value.

$$I_M = \max_r PC_{2D,r}(x).$$

$$I_L(x) = arctan2\left[\sum_{\rho,r} e_{\rho r}(x), \sum_{\rho,r} o_{\rho r}(x)\right].$$

### B. Phase Preserving Denoising

An image denoising method proposed by Kovesi [24] is used in this paper, which is based on the assumption that phase information is the most important feature of images. This method also attempts to preserve the perceptually important phase information in the signal. It uses non-orthogonal, complex valued log-Gabor wavelets, which extract the local phase and amplitude information at each point in the image. The denoising process consists of determining a noise threshold at each scale and shrinking the magnitudes of the

**Special Issue - 2015**

**International Journal of Engineering Research & Technology (IJERT)**
**ISSN: 2278-0181**
**ICESMART-2015 Conference Proceedings**

filter response vector appropriately, while leaving the phase unchanged.

### C. Abbreviations and Acronyms

$A_\rho(x)$      Local amplitude;

$e_\rho(x)$      Real part of the complex-valued wavelet response;

$E(A_\rho)$      Expected value of the Rayleigh distribution at scale $\rho$;

filter      Log-Gabor filter;

$I$      Gray-level input image;

$I_{bwout}$      Final binarized output;

$I_D$      Denoised image;

$I_L$      Local weighted mean phase angle (LWMPA);

$I_{L,bw}$      Binary image corresponding to $I_L$;

$I_M$      Maximum moment of phase congruency covariance (MMPCC);

$I_{M,bw}$      Binarized $I_M$;

$k$      Number of standard deviations of noises;

$o_\rho(x)$      Imaginary part of the complex-valued wavelet response;

$PC_{2D}$      Two-dimensional phase congruency;

$\rho$      Index over filter scales;

$r$      Index over filter orientations;

$N_\rho$      Number of filter scales;

$N_r$      Number of filter orientations;

$T$      Estimated noise threshold;

$W(x)$      Weighting mean function;

$\varphi_\rho(x)$      Local phase;

## IV. PROPOSED METHODOLGY

The final binarized output image is obtained by processing the input image in three steps: preprocessing, main binarization, and postprocessing. The flowchart of the proposed method is shown in Figure 2 and the proposed binarization method algorithm is used in this methodology[35].

### A. Preprocessing

In the preprocessing step, we use a denoised image [5] instead of the original image to obtain a binarized image in rough form. A number of parameters impact the quality of the denoised output image (ID), the key ones being the noise standard deviation [9] threshold to be rejected (k), and the number of filter scales (Nρ) and the number of orientations[12] (Nr) to be used. The parameters in the experiments as follows: k = 1, Nρ = 5 and Nr = 3.

We used Otsu's method [13] on the normalized denoised image, where normalized denoised image is obtained by applying a linear image transform [17] on the denoised image. This approach can also remove noisy and degraded parts of images, because the denoising method attempts to shrink the amplitude information of the noise component. The problem with this approach is that it misses weak strokes and sub-strokes, which means that we cannot rely on its output. To solve this problem, we combine this binarized image with an edge map obtained using the Canny operator [15]. Canny operator is applied on the original document image and for combination those edges without any reference in the aforementioned binarized image are removed. We then

compute a convex hull [11] image of the combined image. At the end of this step, the structure of foreground and text is determined. However, the image is still noisy, and the strokes and sub-strokes have not been accurately binarized. Also, the binarization output is affected by some types of degradation [18]. We therefore include additional steps to deal with them.
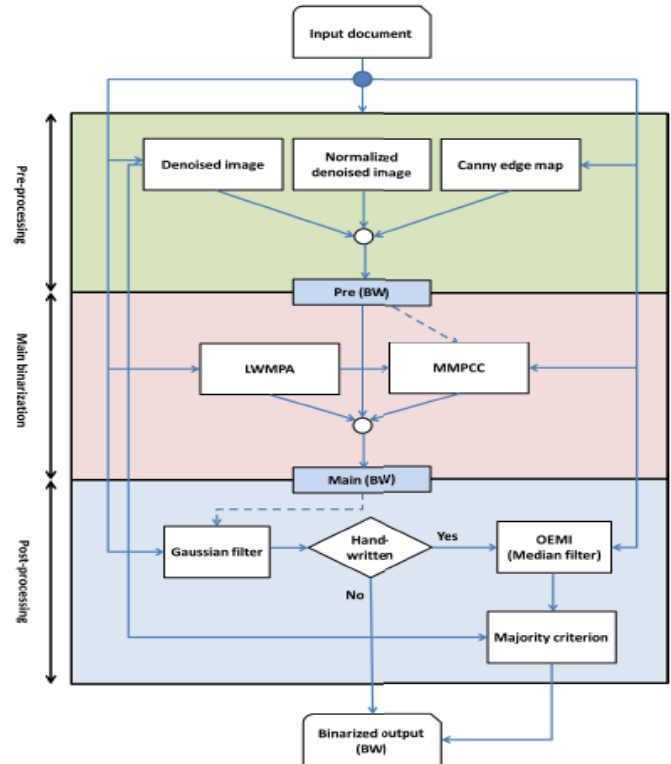


Fig 2. Flowchart of the proposed method

### B. Main Binarization

The next step is the main binarization, which is based on phase congruency features: i) the maximum moment [15] of phase congruency covariance (IM); and ii) the locally weighted [11] mean phase angle (IL).

1) IM: In this paper, IM is used to separate the background from potential foreground parts. This step performs very well, even in badly degraded documents, where it can reject a majority of badly degraded background pixels [17] by means of a noise modeling method. To achieve this, we set the number of two-dimensional log-Gabor filter [19] scales ρ to 2, and use 10 orientations of two-dimensional log-Gabor filters r.

2) IL: We consider the following assumption in classifying foreground and background pixels using IL:

$$P(x) = \begin{cases} 1, & I_L(x) \le 0 \\ 0, & I_L(x) > 0 \quad \& \quad I_{\text{Otsu},bw}(x) = 0, \end{cases}$$

where P(x) denotes one image pixel; and IOtsu,bw denotes the binarized image using Otsu's method [12].

### C. Postprocessing

In this step, we apply enhancement processes [17]. First, a bleedthrough removal [13] process is applied. Then, a Gaussian filter [17] is used to further enhance the binarization output and to separate background from potential foreground, and an exclusion process [14] is applied, based on a median filter [3] and IM maps, to remove background noise and

**Special Issue - 2015**

**International Journal of Engineering Research & Technology (IJERT)**
**ISSN: 2278-0181**
**ICESMART-2015 Conference Proceedings**

objects. Finally, a further enhancement process is applied to the denoised image. The individual steps are as follows.

*1) Global Bleed-Through Exclusion:* Bleed-through degradation [16] is a common interfering pattern and a significant problem in old and historical document images. In this paper, bleed-through is categorized in two classes: i) local bleedthrough [12]; and ii) global bleed-through [12]. Local bleed-through involves pixels located under and near foreground pixels, while global bleed-through involves pixels located far away from the foreground text. Global bleed-through is one of most challenging forms of degradation, as there is no local to enable true text to be distinguished from bleed-through.
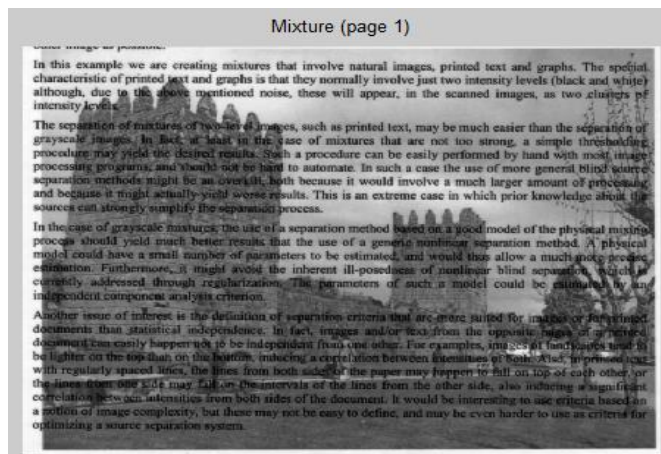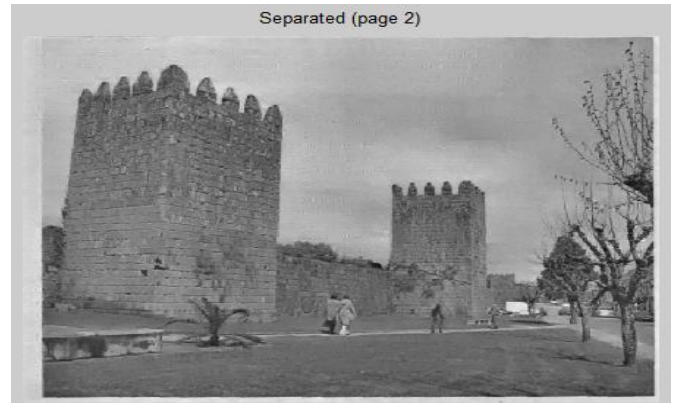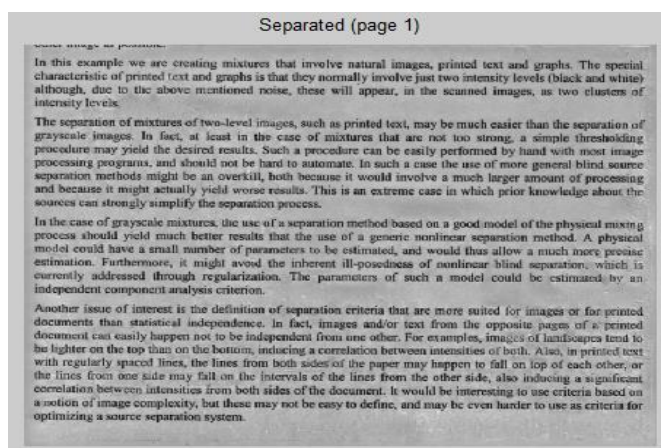


Fig.3 Bleedthrough Image





Fig. 4 Bleed through separated images 1 & 2

*2) Document Type Detection:* At this step, it is necessary need to determine the type of input document we are dealing with. We propose to apply the enhancement processes that are after this step to the handwritten documents types only, and not to machine printed documents. The method we propose for detecting the type of document is straightforward, accurate and fast. The histograms of the handwritten and machine printed documents are shown below
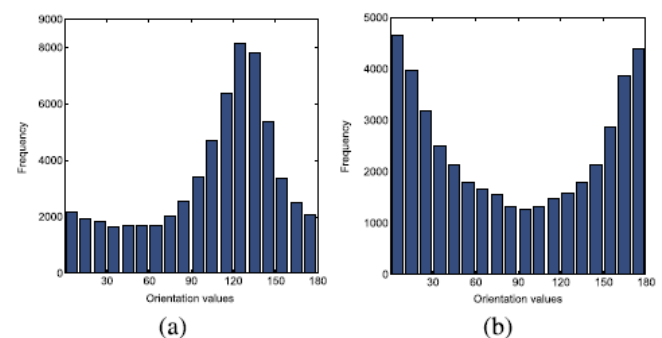


Fig 5. Histogram of handwritten and machine printed documents

*3) Object Exclusion Map (IOEM):* An object exclusion map image (IOEM) is constructed based on a combination of a median filter and a binary map of IM obtained from phase congruency. Any object without a reference in this binary map will be removed from the final post processed binarization results. This approach can remove noise, local bleed-through, and interfering patterns etc.

4) Majority Criterion: We propose a majority criterion based on the denoised image, ID. A majority criterion supposes that early binarization steps provide an optimal or at least a near optimal result. Then, based on the fact that a foreground pixel must have a lower value than its adjacent background pixels, exclusion of the foreground pixels is performed.

## V.  SIMULATION RESULTS

The figures below are obtained by simulation of the proposed methodology.

Figure 6 shows a noisy handwritten document image, which has been denoised by Kovesi's method as in Figure 7. The Figure 8 contains the preprocessed image. The IM is shown in Figure 9. Figure 10 and 11 contain main binarized image nad the post processed image.

**Special Issue - 2015**

**International Journal of Engineering Research & Technology (IJERT)**
**ISSN: 2278-0181**
**ICESMART-2015 Conference Proceedings**

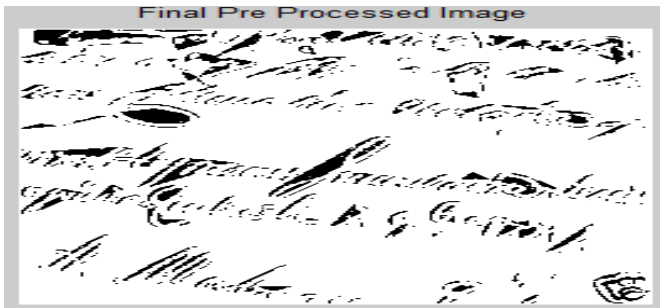Fig. 7: Phase preserving denoising by Kovesi's method
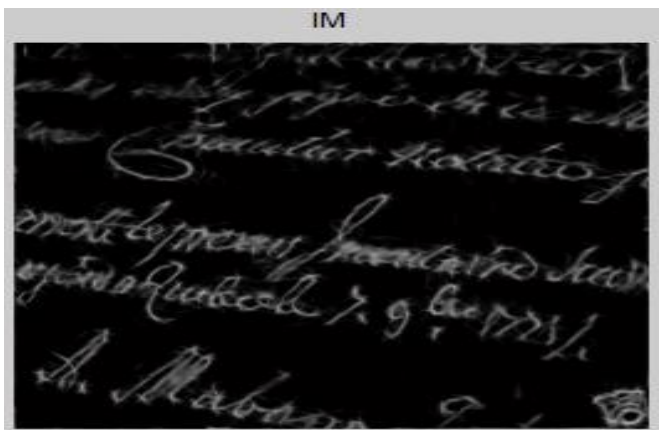

Fig.8: Preprocessed Image
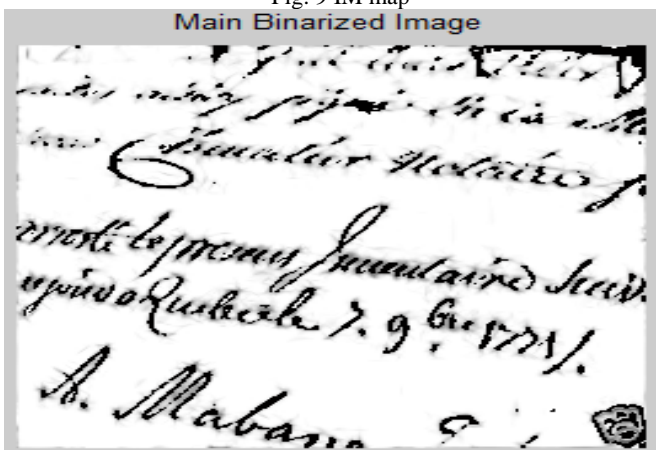

Fig. 9 IM map


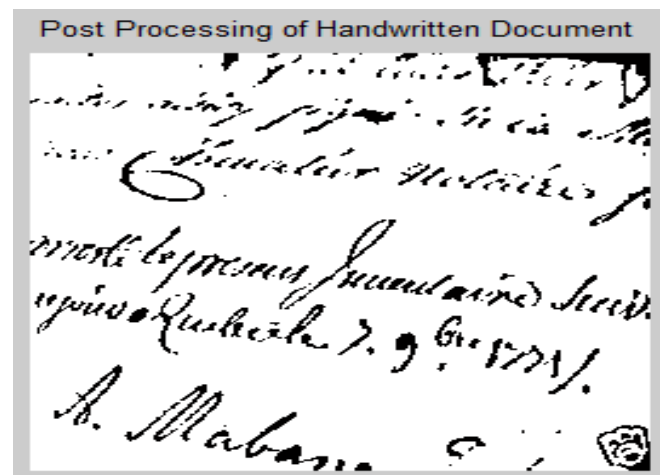Fig. 10 Main Binarized Image


Fig. 11 Image after Post Processing

## VI. CONCLUSION

In this binarization model we have introduced an image binarization method that uses the phase angle information of the input document image, and phase-based features are extracted from that image, and are used to build a model for the binarization of ancient manuscripts. Phase-preserving denoising is done by Kovesi's method and is followed by morphological operations that are used to preprocess the input image. Then, two main phase congruency features, the maximum moment of phase congruency covariance and the locally weighted mean phase angle, which are extracted to perform the main binarization. For post-processing, we have proposed to filter various types of degradation, in particular, a median filter has been used to reject salt and pepper noise, unwanted lines, and interfering patterns. Because some binarization steps work with individual objects components rather than on pixels, a Gaussian filter was used to further separate foreground from background objects by removing Gaussian noise, and to improve the final binary output. We tested the algorithm on a dataset consisting of handwritten and printed document images, and highly satisfactory results were obtained in all the cases.

## ACKNOWLEDGMENT

## REFERENCES

[1] R. F. Moghaddam and M. Cheriet, "A multi-scale framework for adaptive binarization of degraded document images," Pattern Recognit., vol. 43, no. 6, pp. 2186–2198, 2010.

[2] B. Su, S. Lu, and C. L. Tan, "Robust document image binarization technique for degraded document images," IEEE Trans. Image Process., vol. 22, no. 4, pp. 1408–1417, Apr. 2013.

[3] R. F. Moghaddam and M. Cheriet, "AdOtsu: An adaptive and parameterless generalization of Otsu's method for document image binarization," Pattern Recognit., vol. 45, no. 6, pp. 2419–2431, 2012.

[4] J. Sauvola and M. Pietikinen, "Adaptive document image binarization," Pattern Recognit., vol. 33, no. 2, pp. 225–236, 2000.

[5] B. Gatos, I. Pratikakis, and S. Perantonis, "Adaptive degraded document image binarization," Pattern Recognit., vol. 39, no. 3, pp. 317–327, 2006.

**Special Issue - 2015**

**International Journal of Engineering Research & Technology (IJERT)**
**ISSN: 2278-0181**
**ICESMART-2015 Conference Proceedings**

[6] R. Hedjam, R. F. Moghaddam, and M. Cheriet, "A spatially adaptive statistical method for the binarization of historical manuscripts and degraded document images," Pattern Recognit., vol. 44, no. 9, pp. 2184–2196, 2011.

[7] K. Ntirogiannis, B. Gatos, and I. Pratikakis, "A combined approach for the binarization of handwritten document images," Pattern Recognit. Lett., vol. 35, pp. 3–15, Jan. 2014.

[8] M. Valizadeh and E. Kabir, "Binarization of degraded document image based on feature space partitioning and classification," Int. J. Document Anal. Recognit., vol. 15, no. 1, pp. 57–69, 2010.

[9] S. Lu, B. Su, and C. Tan, "Document image binarization using background estimation and stroke edges," Int. J. Document Anal. Recognit., vol. 13, no. 4, pp. 303–314, 2010.

[10] B. Su, S. Lu, and C. Tan, "Binarization of historical document images using the local maximum and minimum," in Proc. 9th IAPR Int. Workshop DAS, 2010, pp. 159–166.

[11] Phase-Based Binarization of Ancient Document Images: Model and Applications Hossein Ziaei Nafchi, Reza Farrahi Moghaddam, Member, IEEE, and Mohamed Cheriet, Senior Member, IEEE. IEEE Transactions on Image Processing, VOL. 23, NO. 7, JULY 2014. 7, ch. 12, pp. 166–180.

[12] B. Su, S. Lu, and C. L. Tan, "A self-training learning document binarization framework," in Proc. 20th ICPR, Aug. 2010, pp. 3187–3190.

[13] B. Su, S. Lu, and C. L. Tan, "Combination of document image binarization techniques," in Proc. ICDAR, Sep. 2011, pp. 22–26.

[14] B. Su, S. Lu, and C. L. Tan, "A learning framework for degraded document image binarization using Markov random field," in Proc. 21st ICPR, Nov. 2012, pp. 3200–3203.

[15] N. Howe, "Document binarization with automatic parameter tuning," Int. J. Document Anal. Recognit., vol. 16, no. 3, pp. 247–258, 2013.

[16] M. Cheriet, R. F. Moghaddam, and R. Hedjam, "A learning framework for the optimization and automation of document binarization methods," Comput. Vis. Image Understanding, vol. 117, no. 3, pp. 269–280, 2013.

[17] N. Howe, "Document image binarization using Markov field model," in Proc. ICDAR, 2011, pp. 6–10.

[18] H. Z. Nafchi and H. R. Kanan, "A phase congruency based document binarization," in Proc. IAPR Int. Conf. Image Signal Process., 2012, pp. 113–121.

[19] H. Z. Nafchi, R. F. Moghaddam, and M. Cheriet, "Historical document binarization based on phase information of images," in Proc. ACCV, 2012, pp. 1–12.

[20] Special issue on recent advances in applications to visual cultural heritage, IEEE Signal Process. Mag., vol. 12, no. 1, pp. 234–778, Jan. 2008.

[21] B. Gatos, K. Ntirogiannis, and I. Pratikakis, "ICDAR 2009 document image binarization contest (DIBCO 2009)," in Proc. 10th ICDAR, Jul. 2009, pp. 1375–1382.

[22] I. Pratikakis, B. Gatos, and K. Ntirogiannis, "H-DIBCO 2010—Handwritten document image binarization competition," in *Proc. ICFHR*, Nov. 2010, pp. 727–732.

[23] P. Kovesi, "Phase preserving denoising of images," in Proc. Int. Conf. Digital Image Comput., Techn. Appl., 1999.

[24] J. Canny, "A computational approach to edge detection," IEEE Trans. Pattern Anal. Mach. Intell., vol. 8, no. 6, pp. 679–698, Nov. 1986.

[25] P. Kovesi, "Image features from phase congruency," Videre, J. Comput. Vis. Res., vol. 1, no. 3, pp. 1–26, 1999.

[26] V. Struc and N. Pavesic, "Phase congruency features for palm-print verification," *IET Signal Process.*, vol. 3, no. 4, pp. 258–268, Jul. 2008.

[27] A. Verikas, A. Gelzinis, M. Bacauskiene, I. Olenina, S. Olenin, and E. Vaiciukynas, "Phase congruency-based detection of circular objects applied to analysis of phytoplankton images," Pattern Recognit., vol. 45, no. 4, pp. 1659–1670, 2012.

[28] L. Zhang, L. Zhang, D. Zhang, and Z. Guo, "Phase congruency induced local features for finger-knuckle-print recognition," Pattern Recognit., vol. 45, no. 7, pp. 2522–2531, 2012.

[29] B. Obara, M. Fricker, D. Gavaghan, and V. Garu, "Contrast-independent curvilinear structure detection in biomedical images," IEEE Trans. Image Process., vol. 21, no. 5, pp. 2572–2581, May 2012.

[30] H. Z. Nafchi, R. F. Moghaddam, and M. Cheriet, "Application of phase-based features and denoising in postprocessing and binarization of historical document images," in Proc. 12th ICDAR, Aug. 2013, pp. 220–224. [31] W. Niblack, An Introduction to Digital Image Processing. Englewood Cliffs, NJ, USA: Prentice-Hall, 1986.

[32] G. Lazzara and T. Geraud, "Efficient multiscale Sauvola's binarization," in Proc. IJDAR, Jul. 2013, pp. 1–19.

[33] T. Lelore and F. Bouchara, "Super-resolved binarization of text based on the fair algorithm," in Proc. ICDAR, Sep. 2011, pp. 839–843.

[34] A. V. Oppenheim and J. S. Lim, "The importance of phase in signals," Proc. IEEE, vol. 69, no. 5, pp. 529–541, May 1981.

[35] "Phase-Based Binarization of Ancient Document Images": Model and Applications Hossein Ziaei Nafchi, Reza Farrahi Moghaddam, Member, IEEE, and Mohamed Cheriet, Senior Member, IEEE 2014.