# Page-level Handwritten Script Identification using Modified log-Gabor filter based features

Pawan Kumar Singh[1], Iman Chatterjee[2]

[1]Department of Computer Science and Engineering
Jadavpur University
Kolkata, India
Email:{pawansingh.ju, raamsarakar}@gmail.com

Ram Sarkar[1]

[2]Department of Computer Science and Engineering
Netaji Subhash Engineering College
Kolkata, India
E-mail: imanchatterjee9@gmail.com

*Abstract—* **Automatic identification of scripts, an imperative research problem during the last few decades, has posed many challenges in any multi-script environment. As India is a multilingual country, therefore, text documents containing more than one language are very familiar phenomenon here. But to digitize these multi-lingual documents using any Optical Character Recognition (OCR) engine, first it is required to recognize the scripts used to write the same. In this paper, a page-level script identification technique for eight popular handwritten scripts *namely*, *Bangla*, *Devanagari*, *Gurumukhi*, *Oriya*, *Tamil*, *Telugu*, *Urdu* along with *Roman* has been proposed. To start with, Modified log-Gabor filters based texture features are designed from each of the document pages. Then the proposed model is evaluated using multiple classifiers and based on their identification accuracies, it is found that Simple Logistic performs the best. Outcome of the present experiment reveals the usefulness of the Modified log-Gabor filters based features in recognition of handwritten *Indic* scripts. A total of 240 document pages is used to carry out the present experiment and it yields 95.57% accuracy in identifying the scripts of the documents. Even if the proposed method is assessed on limited dataset, but considering the intricacies of the scripts, the outcome can be assumed reasonably acceptable.**

**Keywords—** **Page-level script identification, Modified log-Gabor filter, Handwritten *Indic* scripts, Optical Character Recognition.**

## I. INTRODUCTION

One of the key applications in the domain of document image analysis is the OCR engine. Such engine can be defined as the software by which texts (printed or handwritten) from document images can be converted into machine editable electronic form. Till date, a large number of OCR techniques have been proposed in the literature. All these existing works on the OCR development implicitly presume that the script of the documents to be accessed by the underlying algorithm is known in advance. Moreover, OCR engines are in general script specific. As a result, such document processing methodologies would certainly require human intervention in selecting the suitable OCR package. But to develop an automated document processing scheme this is obviously incompetent and unfavorable. So, in this regard, there is a need to develop a pre-OCR module for identifying the script of the documents ready to feed to the OCR engine. Automatic identification of a script used to write a given document page paves the way for quite a few important applications, such as, automatic archiving and indexing of multilingual documents, searching necessary information from digitized archives of the multilingual document images, etc.

Script is the term for a graphic style of writing system that serves as a means to express the written languages. Languages all over the world are typeset in various types of scripts. A script may be either be used by a single language or be shared by many languages, often with minor differences from one language to other. For example, *Bangla* script is used for writing a number of Indian languages like *Bengali*, *Assamese*, *Manipuri*, *Meithei*, etc. Script identification is generally performed at three levels: (a) Page-level, (b) Text line-level and (c) Word-level. Identifying scripts at page-level is much more challenging than its text line- and word-level counterparts. This is due to the fact that script identification at text line- and word-level require fine segmentation of the document pages into their corresponding text lines and words. But each script possesses its own characteristics which are considerably different from others, and this makes a real challenge to develop a script invariant text line or word segmentation method. Therefore, in order to avoid this pitfall, script identification at page-level is recommended.

The research works on script identification, available in the literature, are mainly categorized into two group *viz*., structure-based methods [1-8] and visual appearance-based methods [9-11]. In structure-based methods, a list of connected components such as text-lines, words and characters are used to designed the useful features. This implies that the success rate of classification very much depends on the page segmentation techniques i.e. text-line, word and character segmentation steps. But, it is an agreeable fact that developing a common page segmentation process that best suits for all different script documents is almost impossible. For this reason, structure-based methods hardly meet the criterion as a generalized scheme. In contrast, visual appearance-based methods employ analysis of regions and hence fine segmentation of the underlying document is not necessary. As a result, the task of script classification task becomes significantly simple and is performed faster with the visual appearance-based methods than the structure-based.

In comparison to structure-based methods, relatively few works have been reported in the literature for visual appearance-based methods which, in general, make use of the

texture-based features. G. S. Peake *et al.* [9] reported a method for automatic script and language identification from document images using multiple channel Gabor filters and gray level co-occurrence matrices for seven languages, *viz.*, *Chinese, English, Greek, Korean, Malayalam, Persian* and *Russian*. T. N. Tan [10] developed a texture feature extraction method that was invariant to rotation for automatic script identification for six languages, *viz.*, *Chinese, Greek, English, Russian, Persian* and *Malayalam*. P. S. Hiremath *et al.* [11] proposed a novel texture feature extraction method based on the co-occurrence histograms of wavelet decomposed images, which was tested on eight printed Indian languages *namely*, *Bengali, Hindi, Kannada, Tamil, Telugu, Malayalam, Urdu,* and *English*. It is already stated that methodologies based on visual appearance have the significant relevance in script based retrieval systems because they are not only relatively faster and but also minimize the cost of document handling. So, visual appearance-based schemes can best suited for a generalized approach to the script identification problem. But unfortunately, only a few attempts [4-8] were made towards handwritten script identification of Indian documents in the literature. This has been our primary motivation behind the development of a page-level script identification technique using texture based features for eight popular handwritten scripts *namely*, *Bangla, Devanagari, Gurumukhi, Oriya, Tamil, Telugu, Urdu* and *Roman.*

## II. PROPERTIES OF SCRIPTS USED IN PRESENT WORK

India is a multilingual country with 23 constitutionally recognized languages written in 12 major scripts. The officially recognized languages [12] are *Hindi, Bengali, Punjabi, Marathi, Gujarati, Oriya, Sindhi, Assamese, Marathi, Urdu, Sanskrit, Tamil, Telugu, Kannada, Malayalam, Kashmiri, Manipuri, Konkani, Maithali, Santhali, Bodo, English* and *Dogari*. The 12 major modern scripts are currently being used: *Devanagari, Bangla, Oriya, Gujarati, Gurumukhi, Tamil, Telugu, Kannada, Malayalam, Manipuri, Roman* and *Urdu*.
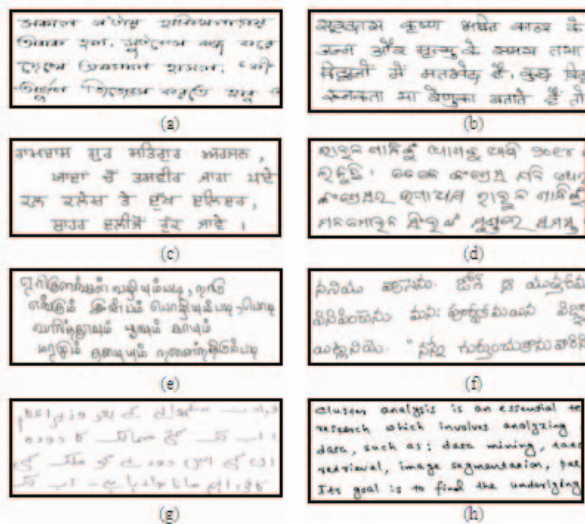


Fig. 1. Snapshots of handwritten document page written in (a) *Bangla,* (b) *Devanagari,* (c) *Gurumukhi,* (d) *Oriya,* (e) *Tamil,* (f) *Telugu,* (g) *Urdu,* and (d) *Roman* scripts.

Of these, *Urdu* is derived from the *Persian* script and *Roman* is derived from the *Latin* script. The first 10 scripts are originated from the early *Brahmi* script (300 BC) and are also referred to as *Indic* scripts [19-20]. A sample portion of handwritten document pages written in eight scripts are shown in Fig.1 and some basic information about eight scripts used for the present work is enlisted in Table I.

## III. PROPOSED WORK

Considering the fact that humans can distinguish between unfamiliar scripts just by visual inspection, script identification may be used for texture classification. Textures are complex visual patterns constituted by subpatterns [14]. Complete analysis of a texture can be done only if it has well-defined subpatterns. However subpatterns are found to lack a sound mathematical model. Thus, we propose a modified log-Gabor filter approach based on Gabor filter for handwritten script identification.

Gabor filters are local and linear band pass filters in which a Gaussian envelop is used to modulate a sinusoidal plane at a certain frequency and orientation. The impulse response of these filters is generated by multiplying a complex oscillation with Gaussian envelope function. The 2D Gabor filter function can be written as [15]:

$$\varphi(x,y) = \frac{f^2}{\pi \gamma \omega} e^{-(\frac{f^2}{\gamma^2}\acute{x}^2 + \frac{f^2}{\gamma^2}\acute{y}^2)} e^{j2\pi f \acute{x}} \qquad (1)$$

where,

$$\acute{x} = x\cos\theta + y\sin\theta$$

$$\acute{y} = -x\cos\theta + y\sin\theta$$

In the spatial domain (Eq. (1)), the product of a 2D Fourier basis function (complex plane wave) with a Gaussian centered at origin acts as the Gabor filter. Here, f is the central frequency of the filter, θ is the rotation angle, λ is sharpness (bandwidth) along the Gaussian major axis, and ω sharpness along the minor axis (perpendicular to the wave). In the given form, the aspect ratio of the Gaussian is 1/λ. This function, in the frequency domain, takes the following analytical form:

$$\varphi(u,v) = e^{\frac{-\pi^2}{f^2}(\gamma^2(u\acute{-}f)^2 + v^2\omega^2)} \qquad (2)$$

where,

$$\acute{u} = u\cos\theta + v\sin\theta$$

$$\acute{v} = -u\cos\theta + v\sin\theta$$

Gabor filters possess excellent joint localization characteristics in both the spatial and the frequency domain and its convolution kernel is obtained by multiplying a Gaussian and a cosine function. However, most applications that employ Gabor filters require a large bank of filters leading to high computational cost. Additionally, they have two main limitations:-

- The maximum bandwidth of a Gabor filter cannot exceed approximately one octave.

TABLE I.        Basic information about the scripts used in the present work.

| Scripts | Origin | Basic Character Set | | Writing Style | Number of Native Speakers in India[13] (Millions) | Used to write languages |
|---|---|---|---|---|---|---|
| | | Vowels | Consonants | | | |
| Bangla | | 11 | 39 | | 207 | Bengali, Assamese, Manipuri, etc. |
| Devanagari | | 15 | 33 | | 366 | Hindi, Nepali, Marathi, Konkani, Sindhi, etc. |
| Gurumukhi | | 12 | 30 | | 57.1 | Punjabi, Sindhi, Braj Bhasha, Khariboli, etc. |
| Oriya | | 14 | 38 | Left to right | 32.3 | Oriya, etc. |
| Tamil | Brahmi | 12 | 18 | | 66.0 | Tamil, Saurashtra, Badaga, Irula, Paniya, etc. |
| Telugu | | 16 | 37 | | 69.7 | Telugu, etc. |
| Roman/Latin | Greek | 5 | 21 | | 341 | German, English, Spanish, Indonesian, etc. |
| Urdu | Persian | 10 | 28 | Right to left | 60.3 | Urdu, Bati, Burushaski, etc. |

- Gabor filters seeking broad spectral information with maximal spatial localization do not yield optimal results.

To overcome the above limitations, an arbitrary bandwidth was used to construct the log-Gabor filter. This bandwidth helps in building a filter with minimal spatial extent. The proposed method for modified log-Gabor based feature extraction is described in the next subsection:

### A. Modified log-Gabor Filter

The log-Gabor function was proposed by Field [16] as a means to overcome the limitations of Gabor function. Field suggests that when viewed on a logarithmic frequency scale, coding of natural images by filters having Gaussian transfer functions is superior. The log-Gabor function has a transfer function of the following form:-

$$G(w) = e^{(-log(w/w_0)^2)} \Big/ (2\,(log(k/w_0)^2\,))  \qquad (3)$$

where, $w_0$ is the filter's central frequency. The term $k/w_0$ should be kept constant for varying $w_0$ in order to get constant shape ratio filters.

The log-Gabor functions possess no DC component, and their transfer function is characterized by an extended tail at the high frequency end. In the domain of statistics of ordinary images, Field's observations claim that these images have amplitude spectra falling off at approximately $1/w$. In order to encode those images that possess such spectral characteristics, filters having similar spectra must be used. Further, log-Gabor functions, with extended tails are able to encode images more efficiently than ordinary Gabor functions since these log-Gabor filters neither over-represent the lower frequency components nor under-represent the high frequency components in any encoding.

For proper texture analysis of an image using the features extracted using log-Gabor filter, the frequency components are necessary. In order to retain the spatial information eliminated by a simple Fourier transform, a Windowed Fourier Transform (WFT) is used. WFT involves multiplication of the image $X(n_0, n_1)$ by a window function $w(n_0, n_1)$. After this, the Fourier transform is applied on the result to obtain:

$$\tau\big(n_0{}^{(0)}, n_1{}^{(0)}, u_0, u_1\big) = \int_{-\infty}^{\infty} X(n_0, n_1) w\big(n_0 - n_0{}^{(0)}, n_1 - n_1{}^{(0)}\big) e^{-2\pi i(u_0 n_0 + u_1 n_1)} dn_0 dn_1$$

$$= e^{-2\pi i\big(u_0 n_0{}^{(0)} + u_1 n_1{}^{(0)}\big)} \Big[ X\big(n_0{}^{(0)}, n_1{}^{(0)}\big) * \big(w\big(n_0{}^{(0)}, n_1{}^{(0)}\big) e^{-2\pi i\big(u_0 n_0{}^{(0)} + u_1 n_1{}^{(0)}\big)}\big)\Big]$$

$$= e^{-2\pi i\big(u_0 n_0{}^{(0)} + u_1 n_1{}^{(0)}\big)} \big[X * m_{u_0, u_1}\big]\big(n_0{}^{(0)}, n_1{}^{(0)}\big) \qquad (4)$$

Here, $(n_0{}^{(0)}, n_1{}^{(0)})$ is the position in the image where the horizontal and vertical frequencies $u_0$ and $u_1$ are computed. The WFT is basically a result of the filter $m_{u_0, u1}$ convolving the image. Both frequency and spatial locations are desired for analyzing texture. Thus we try to achieve a good tradeoff between both frequency and spatial locations with respect to the Uncertainty principle of Fourier transform. A Gaussian function is employed as the optimally concentrated function in both spatial and frequency domain [16] by Gabor transforms. Here, non-isotropic Gaussian is of the form:

$$m_{f,\emptyset}(n_0, n_1) = \frac{1}{2\pi\sigma^2\lambda} e^{-\frac{1}{2\sigma^2}\left(\frac{n_0^2}{\lambda^2} + n_1'^2\right)} e^{2\pi i f n_0'} \qquad (5)$$

$$M_{f,\varphi}(u_0, u_1) = e^{-2\pi^2\sigma^2\left[(u_0'-f)^2\lambda^2 + u_1'^2\right]} \qquad (6)$$

with the center frequency $f = \sqrt{u_0{}^2 + u_1{}^2}$ and the rotated co-ordinates

$$(\acute{n}_0, \acute{n}_1) = (n_0 + \cos\emptyset + n_1\sin\emptyset - n_0\sin\emptyset + n_1\cos\emptyset)$$

Here, $1/\lambda$ is the aspect ratio. Due to the convolution theorem, the filter interpretation of the Gabor transform allows the efficient computation of the Gabor coefficients $G_{f,\varphi}(n_0,n_1)$ by multiplication of the Fourier transformed image $\tau(u_0, u_1)$ with the Fourier transform of the Gabor filter $M_{f,,\varphi}(u_0,u_1)$. The inverse Fourier transform is then applied on the resultant vector as defined below:

$$G_{f,\varphi}(n_0, n_1) = FFT^{-1}\{\tau(u_0, u_1).M_{f,\emptyset}(u_0, u_1)\} \qquad (7)$$

The images, after low pass filtering, are passed as input to a function that computes Gabor energy feature from them. This is done by multiplying in frequency domain (which is equivalent to 2D convolution in the time-space domain). The input image is then passed to a function to yield a Gabor array which is the array equivalent of the image after Gabor
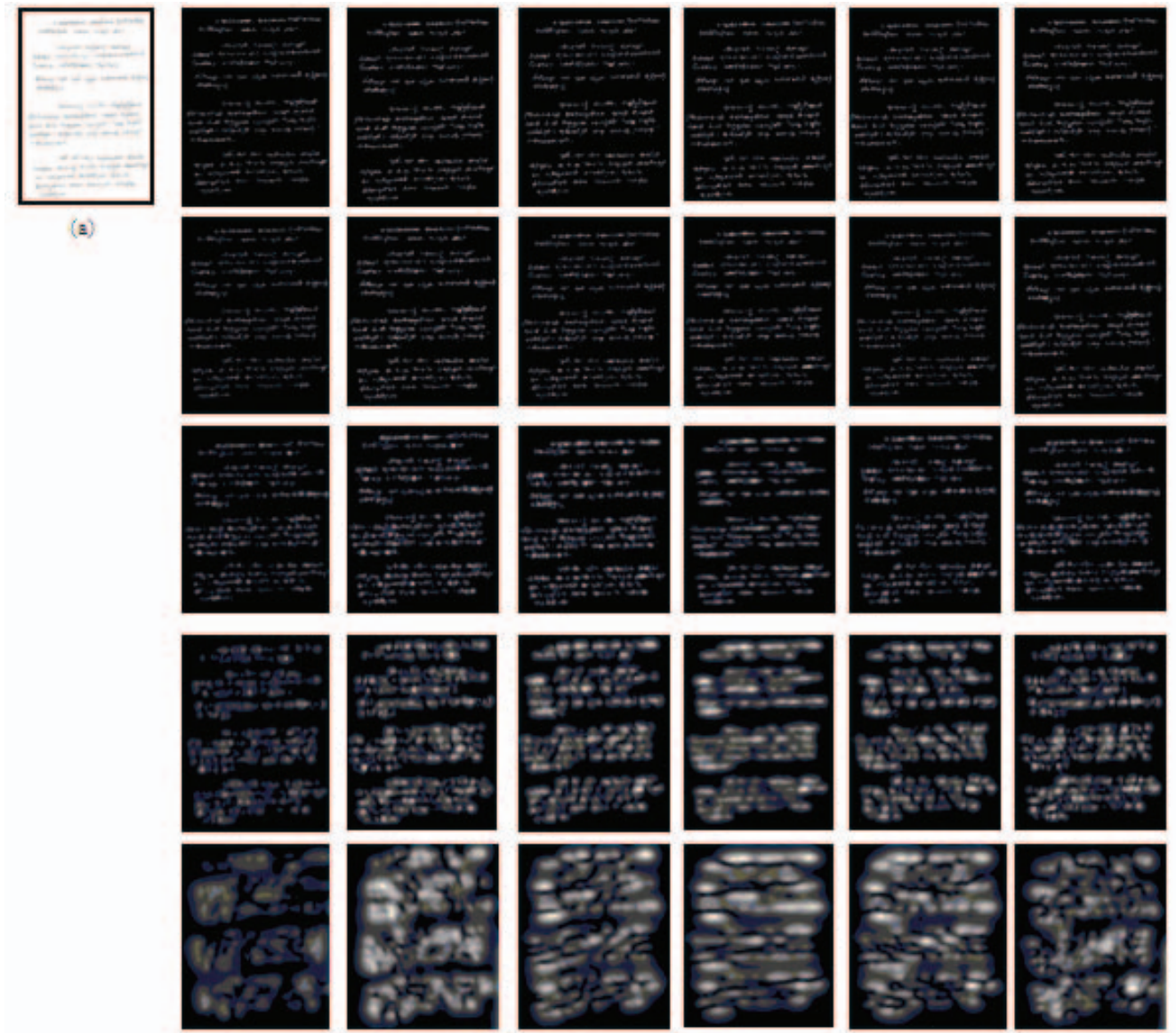
Fig. 2.    Illustration of output images after performing modified log-Gabor filter based approach on a sample *Tamil* script document (a) for 5 scales and 6 orientations (The first row shows the output for $n_s$=1 and six orientations, the second row shows the output for $n_s$=2 and six orientations, and so on).

filtering. The function displays the image equivalent of the magnitude and the real part of the Gabor array pixels.

Consider there are $n_s$ scales and $n_o$ number of orientations, resulting in $n_s \times n_o$ different filters. Let $J$ denote the input image under Fourier transform, $G_{s,o}$ the Gabor filter at scale $s$ and orientation $o$, and $V_{s,o}$, the output of the convolution of $G_{s,o}$ and $J$.

$$V_{s,o} = J * G_{s,o} \qquad (8)$$

Local responses of each of the Gabor filters can also be expressed in terms of amplitude $A_{s,o}(x, y)$ and energy $E_{s,o}(x, y)$ as defined below:

$$A_{s,o} = |V_{s,o}(x,y)| \qquad (9)$$

and

$$E_{s,o}(x,y) = |Real\{V_{s,o}(x,y)\}| - |\,Img\,\{V_{s,o}(x,y)\}| \qquad (10)$$

where, $(x, y)$ denotes the 2D position of a pixel, and *Real* and *Img* denote the real and imaginary parts of the filter responses respectively. Next, we define the median over all orientations for a fixed scale $s$ for $A_{s,o}$ and $E_{s,o}$ as follows:

$$A_s(x,y) = median\{o = 1,2...,n_o\}\,A_{s,o}(x,y) \qquad (11)$$

$$E_s(x,y) = median\{o = 1,2...,n_o\}\,E_{s,o}(x,y) \qquad (12)$$

Finally, the phase symmetry measure, denoted by $\eta(x,y)$ is defined as follows:

$$\eta(x,y) = \frac{\sum_{s=1}^{n_s} E_s(x,y)}{\sum_{s=1}^{n_s} A_s(x,y)} \qquad (13)$$

TABLE II. Success rates of the proposed script identification technique using seven classifiers and their corresponding scores at 95% confidence level (shaded cell indicates best case).

| | Classifiers | | | | | | |
|---|---|---|---|---|---|---|---|
| | Naïve Bayes | Simple Logistic | MLP | SVM | Random Forest | Bagging | MultiClass Classifier |
| Success Rate (%) | 83.09 | **95.57** | 92.30 | 93.36 | 90.34 | 88.95 | 92.25 |
| 95% confidence score (%) | 89.39 | **98.98** | 96.11 | 96.28 | 91.62 | 90.56 | 97.80 |

For the present work, features based on modified log-Gabor filter have been extracted for 5 scales ($n_s$=1, 2, 3, 4, and 5) and 6 orientations ($n_o$= $0^0$, $30^0$, $60^0$, $90^0$, $120^0$, and $150^0$), where each filter has to be convolved with the input image to obtain 30(5*6) different representations (response matrices) for the same image. These response matrices are then converted to feature vectors. Each input image provides us with one feature vector consisting of 30 elements. Application of the modified log-Gabor filter based approach on a sample *Tamil* script document for 5 scales and 6 orientations are shown in Fig. 2.

## IV. EXPERIMENTAL RESULTS AND DISCUSSION

For the experimental purpose, a total of 240 document pages written in eight scripts *namely, Bangla, Devanagari, Gurumukhi, Oriya, Tamil, Telugu, Urdu* and *Roman* have been collected. Here, each script contributes to exactly 30 documents. The document pages used in the experiment are collected from persons of different age-groups and educational backgrounds. They were told to select the content of the document as they wished to. Also, no other restriction was imposed on them except the ink of the pen. Writers contributed to the present experiment were asked to apply a black or blue ink pen while writing the text in a A-4 size page. These document pages are scanned at a resolution of 300 dpi using a flat-bed HP scanner. Among the 240 document pages, 160 (20 from each script) pages are applied to train the proposed classification algorithm and the remaining pages are applied for the testing purpose. The proposed technique is evaluated by using seven well-known classifiers such as Naïve Bayes, Simple Logistic, Multi Layer Perceptron (MLP), Support Vector Machine (SVM), Random Forest, Bagging and MultiClass Classifier. Success rates of the seven classifiers and their corresponding scores achieved at 95% confidence level are shown in Table II.

TABLE III. Confusion matrix produced by Simple Logistic (Here, A=*Bangla*, B=*Devanagari*, C=*Gurumukhi*, D= *Oriya*, E=*Tamil*, F=*Telugu*, G=*Urdu*, and H=*Roman*).

| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| A | 30 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| B | 0 | 29 | 1 | 0 | 0 | 0 | 0 | 0 |
| C | 0 | 5 | 25 | 0 | 0 | 0 | 0 | 0 |
| D | 2 | 0 | 0 | 28 | 0 | 0 | 0 | 0 |
| E | 0 | 0 | 0 | 0 | 29 | 0 | 0 | 1 |
| F | 0 | 0 | 0 | 0 | 1 | 29 | 0 | 0 |
| G | 0 | 0 | 0 | 0 | 0 | 0 | 30 | 0 |
| H | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 29 |

It is observed from Table II that the best identification accuracy is achieved as 95.57% for Simple Logistic. In the present work, detailed error analysis (showing weighted average value for all the classes for Simple Logistic) with respect to different parameters, *namely*, Kappa statistics, mean absolute error, root mean square error, True Positive Rate (TPR), False Positive Rate (FPR), precision, recall, F-measure and ROC area are computed as 0.9493, 0.0408, 0.1191, 0.956, 0.006, 0.961, 0.956, 0.956 and 0.997 respectively.

It is evident from the confusion matrix that the document pages written in *Gurumukhi* script have been generally misclassified as *Devanagari* script pages (an example is shown in Fig. 3). The main reason for this misclassification is similarity in the basic character set of both the scripts and sometimes the quality of document pages due to presence of noise, outliers, etc creates a challenge.
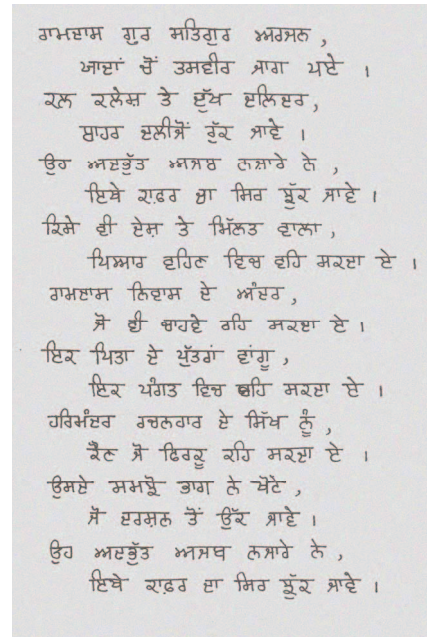


Fig. 3. A sample document page written in *Gurumukhi* script misclassified as *Devanagari* script by the present technique.

## V. CONCLUSION

Script identification, a challenging research problem in any multi-lingual environment, has got attention to the researchers few decades ago. Research in the field of script identification

aspires the visualization and foundation of an automatic system capable of differentiating all official handwritten *Indic* scripts along with *Roman* script. In this paper, we proposed a texture feature based approach to page-level script identification for some of the *Indic* script documents along with *Roman* script. At present, we applied a total of 30 features using Modified log-Gabor filter based technique and the overall accuracy of the system is found to be 95.57% which is quite impressive bearing in mind the complexities of the scripts. This technique is also useful in cases where the amount of text required for accurate recognition is quite small; with as little as few words are suffice in some cases. Though the computational complexity relating to the feature extraction is very less but, the technique has some limitations such as; it ignores the spatial relationship between the texture patterns and it is more susceptible towards noise.

In future, the technique could be extended to recognize other scripts in any multi-script environment. More data samples considering all official *Indic* scripts will be collected in future for detailed evaluation of the developed methodology. In a nutshell, the technique could be used as a general script identification module for the development of multi-script OCR system.

## REFERENCES

[1] D S. Wood, X. Yao. K. Krishnamurthi, L. Dang, "*Language identification from for printed text independent of segmentation*", In: Proc. of International Conference on Image Processing, pp. 428-431, 1995.

[2] A. L. Spitz, "*Determination of the script and language content of document images*", In: IEEE Transaction on Pattern Analysis and Machine Intelligence, vol. 19, pp.234-245, 1997.

[3] J. Hochberg, K. Bowers, M. Cannon, P. Kelly, "*Script and language identification for handwritten document images*", In: International Journal on Document Analysis and Recognition, vol. 2, pp. 45-52, Feb 1999.

[4] R. Sarkar, N. Das, S. Basu, M. Kundu, M. Nasipuri, D. K. Basu, "*Word level script Identification from Bangla and Devnagari Handwritten texts mixed with Roman scripts*", In: Journal of Computing, vol. 2, issue 2, pp. 103-108, Feb 2010.

[5] K. Roy, U. Pal, "*Word-wise Handwritten Script Separation for Indian postal automation*", In: Proc. of 10th International Workshop on Frontiers in Handwriting Recognition, La Baule, pp. 521-526, 2006.

[6] S. Chanda, S. Pal, K. Franke, U. Pal, "*Two-stage Approach for Word-wise Script Identification*", In: Proc. of 10th International Conference on Document Analysis and Recognition, pp. 926-930, 2009.

[7] P. K. Singh, R. Sarkar, N. Das, S. Basu, M. Nasipuri: "*Identification of Devnagari and Roman script from Multiscript Handwritten documents*", In: Proc. of 5th International Conference on PReMI, LNCS 8251, pp. 509-514, Dec 2013.

[8] P. K. Singh, R. Sarkar, N. Das, S. Basu, M. Nasipuri: "*Statistical Comparison of Classifiers for Script Identification from Multi-script Handwritten documents*", In: International Journal of Applied Pattern Recognition, vol. 1, No. 2, pp. 152-172, 2014.

[9] G. S. Peake, T. N. Tan, "*Script and language identification from document images*", In: Proc. of Eighth British Mach. Vision Conf., vol. 2, pp. 230-233, Sept.1997.

[10] T. N. Tan, "*Rotation invariant texture features and their use in automatic script identification*", In: IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 20, pp.751-756, 1998.

[11] P. S. Hiremath, S.Shivashankar, "*Wavelet Based Co-occurrence Histogram Features for Texture Classification with an Application to Script Identification in a Document Image*", In: Pattern Recognition Letters 29, pp 1182-1189, 2008.

[12] http://www.newworldencyclopedia.org/entry/Languages_of_India Retrieved 2015-02-15.

[13] http://timesofindia.indiatimes.com/india    Retrieved 2010-03-14.

[14] R. C. Gonzalez, R. E. Woods, "*Digital Image Processing*", vol. I. Prentice-Hall, India (1992).

[15] http://www.csse.uwa.edu.au/~pk/research/matlabfns/PhaseCongruency/Docs/convexpl.html

[16] D. J. Field, "*Relations between the statistics of natural images and the response properties of cortical cells*", In: J. Opt. Soc. Am. A, vol. 4, pp. 2379-2394, 1987.