

Camera-Based Bi-lingual Script Identification at Word Level using SFTA Features



B.V.Dhandra, Satishkumar Mallappa, Gururaj Mukarambi

Abstract: Most of the documents in various application areas like Government, Business and Research is available in the form of bi-lingual/multi-lingual text document. The multilingual documents are captured from video/camera for identification of script of the text document for automatic reading and editing. In this paper, an attempt is made to address the problem of script identification from camera captured document images with SFTA features. The input image is decomposed into a group of binary images by applying TTBD with fixing the number of the threshold as $n_i = 3$ empirically, on each decomposed binary image, Box Count, Mean Gray Level, and Pixel Count are extracted to form the feature vector. This feature vector is submitted to K-NN classifier to identify the scripts of the input document image. In all 10 scripts of the Indian languages are considered along with common English language as bi-lingual documents. The novelty of the paper is that 7 features are selected as potential features to obtain the highest accuracy. Features like Box Count (3), Mean Gray Level (2), and Pixel Count (2) have obtained the 87.02% recognition accuracy for English and Hindi Script combinations for the dataset collected and encouraging results for other combinations. These 7 potential features were selected using the technique named as feed-forward feature selection, from the set of all 18 features.

Index Terms: KNN, LBP, SFTA, SVM, TTBD.

I. INTRODUCTION

Script identification of the camera-based document images is a complex and challenging problem and it is mainly two-fold: (i) Typical challenges of the camera-captured images may have blurred, uneven illumination, complex background, etc. and (ii) challenges related to shape, size, and orientation of the text written in different scripts. In the Indian context, many documents can be found in bi-lingual/multilingual scripts e.g., signal boards on the high way, vehicle number details, route board of the buses/goods vehicles, hotels names, shops, etc. may have multi-scripts, when the scene is captured from the camera/video.

Hence, automatic processing of the camera-based multi-script documents needs to address for robust multi-script OCR system. This paper has been organized as follows: In Section-II literature survey is presented. The algorithm and block diagram of the proposed method is explained in Section-III, Section-IV contains the experimental results and discussions. The conclusion is given in Section-V.

II. LITERATURE SURVEY

The problem of script identification from the multi-script document is the thrust area of research, lot of research is carried out on script identification from the multi-script scanned document images and can be found in Dhandra et al.[4] and it is observed that very little work has been carried out on the camera captured multi-script text documents. Bhunia et al.[2], have proposed the method name attention mechanism for script identification from natural scene camera/video text images. They have converted input images into patches and feed to CNN-LSTM framework. The local features are generated from the attention-based patch weighting scheme then they have performed the dynamic weighting of local and global features by using the dynamic weighting technique. On the attention based patches, they have extracted 256 dimension latent features by using CNN network on each patch. In the classification process, the local and global features are fed to a fully connected layer. At last the attention-wise summation is performed on all the patch-wise classes. They have experimented on four publicly available benchmark datasets; SIW-13, CVSI-15, ICDAR-17, and MLe2e and achieved the recognition accuracy as 96.50%, 97.75%, 90.23%, and 96.70% respectively, but the time complexity is significantly high due to large size feature set. Recently, Madhura Jajoo et al.[11], have proposed the multi-script identification for the scripts of Bangla, Devanagari, and Roman scene text at the word level, they have extracted the shape and texture (GLCM and HoG) features, combined them to create the feature vector of size 143, and this feature vector is fed to 5 popular classifiers namely Naïve Bayes, MLP, SVM, Multi-class and simple Logistic classifiers. The MLP has exhibited a maximum accuracy of 90% among them and time complexity is significantly high. Luis Gomez et al.[10], have presented the scene text script identification by the improved patch-based method. The proposed method called patch-based classification with the ensemble of the conjoined networks. The input images are represented as the local descriptors from the patches.

Revised Manuscript Received on 30 July 2019.

* Correspondence Author

Dr. B. V. Dhandra*, Professor, Department of Computer Science, Gulbarga University, Kalaburagi, India.

Mr. Satishkumar Mallappa, Research Scholar, Department of Computer Science, Gulbarga University, Kalaburagi, Karnataka, India

Dr. Gururaj Mukarambi, Assistant Professor Department of Computer Science, Gulbarga University, Kalaburagi, India.

© The Authors. Published by Blue Eyes Intelligence Engineering and Sciences Publication (BEIESP). This is an open access article under the CC-BY-NC-ND license <http://creativecommons.org/licenses/by-nc-nd/4.0/>

Camera-Based Bi-lingual Script Identification at Word Level using SFTA Features

These patches are used to learn the discriminative stroke-parts. In training, they have used the Stochastic Gradient Descend (SGD) with the momentum and L2 regularization. They have used the 64 images to make mini batches. The discriminative stroke-parts are considered as features and fed to global classifier rule. In the experimental setup, they have considered the five benchmark datasets called SIW-13 having 13 scripts; Tibetan, Thai, Russian, Mongolian, Korean, Kannada, Japanese, Hebrew, Greek, English, Chinese, Cambodian, and Arabic, MLe2e having 4 scripts; Kannada, Chinese, Latin, and Hangul. CVSI-2015 dataset contains the 10 scripts namely; English, Hindi, Bengali, Oriya, Gujarati, Punjabi, Kannada, Tamil, Telugu, Arab, ICDAR-2013, dataset having Kannada, Tamil, Hindi, Chinese and English scripts and ALIF have the Arab script. On these datasets, they have obtained the recognition accuracy of 94.80%, 86.8%, 90.60%, 74.7%, and 100% respectively and their feature size is large. Jan Zdenek et al.[7], have introduced the new approach to identify the scripts from scene text. They have developed the triplets and extracted the local convolution features with the combination of Bag-of-Visual-Words (BoVW). They have conducted the on three public benchmark datasets namely; SIW-13, MLe2e,

and CVSI-2015. On these three datasets, they could achieve the recognition accuracy of 91.62%, 94.08%, and 96.63% respectively. Gururaj et al.[6], have proposed the LBP features based on camera captured tri-lingual scripts namely; English, Hindi, and Kannada. They have extracted 59 LBP features at the block level and reported the recognition accuracy of 96.60% for a block of 128x128, 99.90% for a block of 512x512 with K-nearest neighborhood classifier and 98.00% for a block of size 128x128, 98.07% for a block of size 256x256 and 98% for a block of size 1024x1024 with SVM classifier. Their feature vector size is large. LinLin Li et al.[9], have proposed signature and template based similarity measure for the scripts of Roman, Arabic, Chinese, Cyrillic, Greek, Hebrew, Japanese, Bengali, and Thai. They used signature generating operation, which is more able to capture the interaction without the text line direction. One template for each candidate script is generated; then, they have extracted the signature and made the clusters using a hierarchical clustering algorithm. In the experimental setup, the cosine distance measure was used with 0.02 as the maximum radius of the cluster. Hamming distance is used to measure the query document and script template for recognition and has yielded 91% accuracy.

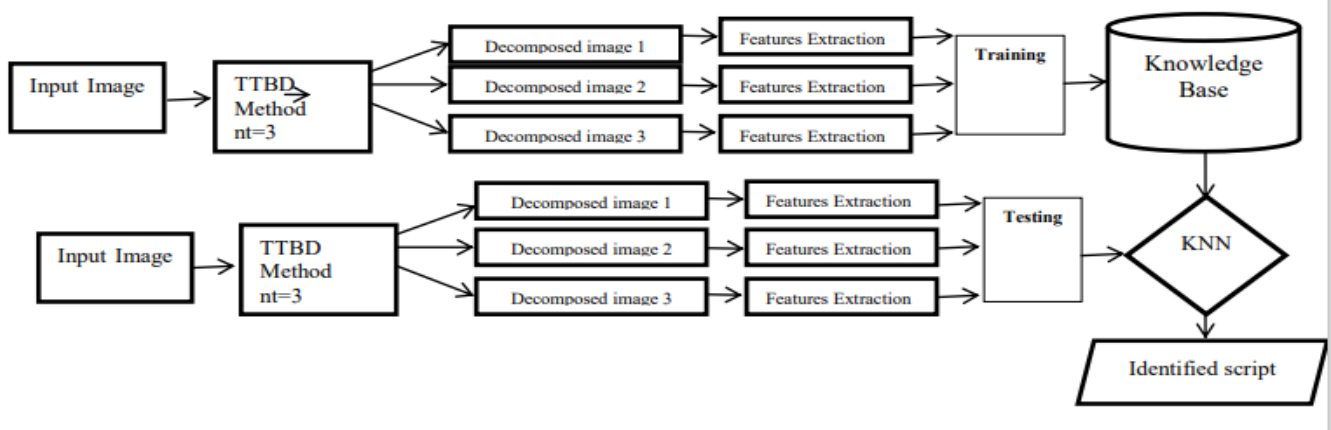


Fig 1. Block diagram of the proposed method

O.K. Fasil et al.[14], have proposed Gabor, Log-Gabor and wavelet texture features for Kannada, English, and Malayalam script identification at word level from camera captured bus signboards. They performed a series of morphological operations to localize the text from bus signboard image and extracted the 60 texture features from each word. They have applied the K-nearest neighborhood and support vector machine with Linear and Radial Gaussian Basis function for classification. Then PCA was used to reduce the features dimensional vector. They have reported maximum F-measure as 0.975 for Kannada and Malayalam with PCA, with feature size 60 and is large. Anirban Mukhopadhyay et al.[1], have attempted the problem of scene text detection in the multi-lingual scenario by using the novel one-class classifier(OCC). They have considered the English, Hindi, and Bangla scripts. They have generated the candidate text regions from scene text using MSER and SWT method. One class classifier has been trained to label the candidate regions with the handcrafted feature descriptor viz., HoG, GLCM, DCT and Gabor filter based features. The feature vector for each descriptor is for HoG 144, GLCM 22, DCT 1024, and for Gabor 640 feature adding to 1830

features. They have reported the highest recognition accuracy as 74.13% and 75.44% for OCC with GLCM and Gabor features respectively and they also reported the recognition accuracy with other four popular classifiers; One-class SVM, Decision Tree, Naïve-Bayes and AdaBoost. The highest accuracy is obtained in AdaBoost of 85.08% from GLCM features and One-class SVM 75.44% from Gabor features. However, their feature size is considerably large. From the above literature, it is observed that the script identification from the camera captured document images still suffers from the size of the feature set, recognition accuracy, and robustness. Hence, the proposed method has been carried out on word-wise script identification based on camera captured document images containing the combination of 11 scripts: English, Hindi, Kannada, Telugu, Tamil, Urdu, Malayalam, Bangla, Gujrathi, Oriya and Punjabi, taking English as the common language for all scripts.

I. PROPOSED METHOD

The problem of bi-script identification of the camera captured bi-lingual text document images is addressed by proposing the Segmentation based Fractal Texture Analysis features [5] and K-nearest neighbor classifier. In the following, the flow of the proposed system is presented with the block diagram. In this method, the input text document of the camera-based image is decomposed into a set of binary images by applying TTBD with $n_t (=3)$ thresholds that are

$t \in \{1, \dots, n_t\}$. The resultant binary image is

$$I_b(x, y) = \begin{cases} 1, & \text{if } t_l < I(x, y) \leq t_u \\ 0, & \text{otherwise} \end{cases} \quad (2.1)$$

Equally spaced threshold strategy is applied with lower and upper threshold.

$$t_i = \left\lfloor \frac{n_t}{n_t + 1} \cdot i \right\rfloor, \quad i = 1, 2, \dots, n_t \quad (2.2)$$

For each decomposed binary image, the Boxcount is obtained by using the formula

$$D_0 = \lim_{\epsilon \rightarrow 0} \frac{\log N(\epsilon)}{\log \epsilon^{-1}} \quad (2.3)$$

where $N(\epsilon)$ is the number of hyper_cubes of dimension ϵ and length ϵ that fill the object.

Then meangraylevel and pixel counting is performed. From these three features (BoxCount, Meangraylevel, Pixelcount) the feature vector is generated.

The classification process is performed by using various distance metrics like Euclidean, City block, Cosine, and correlation metric. The Correlation distance metric has outperformed as compared to other distance metric. Hence, the correlation metric is consider as the distance metric for similarity measure. The following equation (2.4) presents the correlation distance metric:

$$d_{st} = 1 - \frac{(x_s - \bar{x}_s)(y_t - \bar{y}_t)'}{\sqrt{(x_s - \bar{x}_s)(x_s - \bar{x}_s)' (y_t - \bar{y}_t)(y_t - \bar{y}_t)'}} \quad (2.4)$$

where $\bar{x}_s = \frac{1}{n} \sum_j x_{sj}$ $\bar{y}_t = \frac{1}{n} \sum_j y_{tj}$

3.1 DATA COLLECTION

The standard data sets are not available for all the scripts considered in this paper for camera captured text document images at the word level. Hence, sample data sets are collected from the various sources like News papers, Magazines, Fiction and non-fiction novels, printed documents with a camera of 4920x3264 megapixel resolutions. The

dataset consists of 11 scripts; Hindi, Kannada, English, Telugu, Tamil, Malayalam, Oriya, Bangla, Gujarathi, Punjabi and Urdu. Each with 10,000 word images contributing to 1,10,000 word images for testing the performance of the proposed script identification system. The dataset collected for the study is named as Camera Based Word Images of 11 scripts (CBWI-11).

Standard datasets are used for Identification and Recognition of scripts namely CVSI-2015.[13], SIW-13.[3], ICDAR-2017.[12], and MLe2e.[8]. The CVSI-2015 dataset contains 10 scripts namely; English, Hindi, Bengali, Oriya, Gujarathi, Punjabi, Kannada, Tamil, Telugu, Arab and the dataset size of 10,665; SIW-13 has the 13 scripts; Tibetan, Thai, Russian, Mongolian, Korean, Kannada, Japanese, Hebrew, Greek, English, Chinese, Cambodian, and Arabic. The total number of images in the dataset is 16,291; The ICDAR-2017 dataset has 9 different scripts; German, Korean, Japanese, Arabic, English, French, Chinese, Italian and Bangla. The dataset size is 68,613. The MLe2e is the scene text images of 4 scripts; Kannada, Chinese, Latin, and Hangul and the total 1,821 images. These data sets are also considered for performance evaluation of the proposed system as they have some scripts in common considered in this paper. Following Fig.1 shows the sample images collected from various sources:





(x) Telugu (xi) Urdu

Algorithm : Script Identification from multi-script document images

Input : Multi script document image.

Output: Identification of the input Script of the document.

Training phase:

Start

Fig 2. Sample Images Captured By The Camera. The Algorithm For The Proposed System Is Presented In The Following:

English				
Hindi				
Bangla				
Gujrathi				
Punjabi				
Oriya				
Urdu				
Kannada				
Telugu				
Tamil				
Malyalam				

Fig 3: (i) Input image, (ii),(iii),(iv) decomposed binary images of the input image with $n_t = 3$.

Step1: Convert the given input image into a binary image.

Step2: Apply the equally spaced threshold strategy to Decompose the binary image into a set of decomposed binary images by applying the TTBD technique.

Step 3: Compute the Boxcount for each of the decomposed binary images generated in Step-1.

Step 4: Compute the Meangraylevel

Step 5: Compute the pixelcount

Step6: Create the feature vector (Boxcount, Meangraylevel, pixelcount).

OR

Step 7: Select the three Boxcount, two Meangraylevel and two pixel count out of 18 features.

Step 8: Store the feature vector in the knowledge base by labeling the script.

Stop

Testing Phase:

Start

Step 1: Repeat the Steps 1 to 8 of the Training phase.

Step 2: Store the feature set in the test database.

Step 3: Compute the Correlation distance metric between the training and testing features and store it.

Step 4: Submit the Correlation distance metric to K-NN classifier with $k=3$ and recognize the script as the label of the trained feature corresponding to the majority voters supporting the minimum distance.

Stop

For this proposed work, the n_i value is fixed at 3 for the experiment and TTBD method is applied for input image. The TTBD method decompose into 3 binary images for each of 6(3 pairs) thresholds (Box Count, Mean Gray Level, Pixel Count) then it becomes $6 \times 3 = 18$ features. On these 18 features, 7 potential features are selected by adding the features in the feature vector in the increasing order from 1 to 18. In this

SI No.	Scripts	Own (CBWI-11) Dataset	CVSI-2015 Dataset
1	English,Hindi	86.94%	77.25%
2	English,Bangla	81.02%	81.30%
3	English,Gujrathi	86.91%	85.05%
4	English,Punjabi	82.90%	86.65%
5	English, Oriya	84.37%	77.25%
6	English,Urdu	77.62%	Urdu Script -NR-
7	English, Kannada	80.55%	89.45%
8	English, Telugu	83.69%	82.25%
9	English, Tamil	79.70%	84.65%
10	English, Malayalam	83.31%	Malayalam Script -NR-

process, the first feature is selected and observe the result and next, likewise for all features up to 18. If the result is increased by adding features, those features are saved into a feature vector and if the result is not changed and decreased then that feature is removed from the feature vector. The selected feature vector is submitted to the K-NN classifier to obtain the recognize the script.

I. EXPERIMENTAL RESULTS AND DISCUSSION

To test the performance of the proposed method, an extensive experiment is carried out on our own data set named as CBWI-11 and standard benchmark publicly available dataset named as CVSI-2015 for identification of the scripts. The CBWI-11 dataset contains the 11-scripts viz., English, Hindi, Kannada, Telugu, Tamil, Urdu, Malayalam, Bangla, Gujrathi, Oriya and Punjabi, and CVSI-2015 dataset contains the 10-scripts; English, Hindi, Arab, Bengali, Gujrathi, Punjabi, Oriya, Kannada, Telugu, Tami. The following Fig.2 presents input images taken from CBWI-11 dataset along with it's corresponding decomposed binary images obtained after applying the TTBD with $nt=3$.

The correlation distance metric is computed between the train and test images, then K-nearest neighbor classifier is applied

with $k=3$ (the optimum value of K obtained empirically). Table 1 below shows the experimental results.

Table 1: Bilingual Script Identification using 7 features.

SI No.	Scripts	Own (CBWI-11) Dataset	CVSI-2015 Dataset
1	English,Hindi	87.02%	81.35%
2	English,Bangla	64.54%	80.45%
3	English,Gujrathi	80.12%	68.45%
4	English,Punjabi	71.53%	80.05%
5	English, Oriya	74.09%	84.00%
6	English,Urdu	64.86%	Urdu Script -NR-
7	English, Kannada	63.45%	65.90%
8	English, Telugu	70.62%	78.75%
9	English, Tamil	71.08%	88.95%
10	English, Malayalam	68.25%	Malayalam Script -NR-

Table 2. Bilingual Script Identification Results in % using Set of all 18 SFTA Features

The above Table-1.shows the Bilingual Script Identification Results in % using 7 Selected Features using K-NN classifier with $k=3$, 10-fold cross validation, correlation distance metric. It is observed that the English and Hindi combination has the highest recognition accuracy of 87.02% and English with Kannada has the lowest recognition accuracy of 63.45% with 7 potential features. For the CVSI-2015 dataset English and Tamil has given the highest recognition accuracy of 88.95% and English with Kannada has the lowest accuracy of 65.90%. From the above Table.1 it is noticed that for both the datasets the lowest recognition accuracy is obtained for English and Kannada scripts. Urdu and Malayalam scripts are not reported in CVSI-2015 dataset due to this recognition is not obtained for Urdu and Malayalm. The following Fig 4 shows the graphical representation of average reconignition based on 7 features.

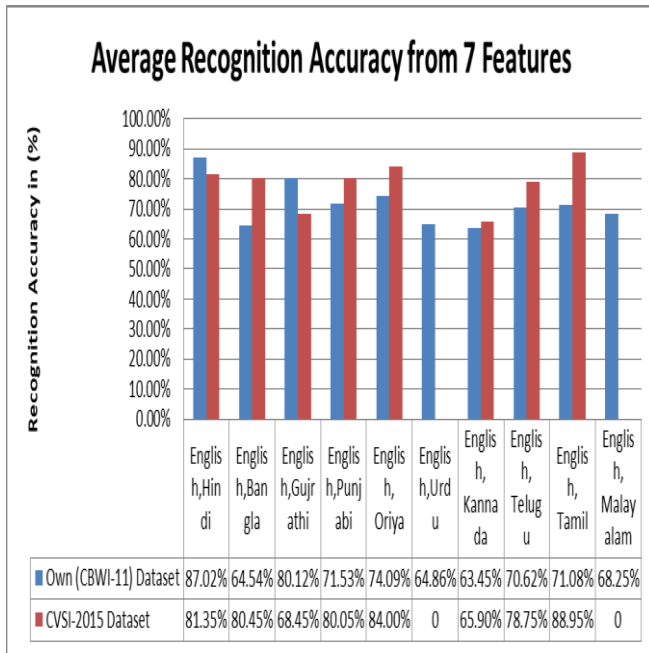


Fig 4: Average Recognitin Accuracy from using K-NN Classifier.

The following Table 2. shows the average recognition accuracy for set of all 18 SFTA features with K-NN classifier (k=3, 10-fold cross validation, correlation distance metric) for both the datasets; Own(CBWI-11) and CVSI-2015. Experimental results have shown that the script identification accuracy for English and Hindi is 86.94%, the highest accuracy among the combination of English with other Indian languages considered in the study. Further the recognition accuracies of all other bi-script combinations are also reasonably high. For the standard dataset CVSI-2015 English and Kannada is 89.45%, the highest as compare to other combination. English with Hindi and English with Oriya documents have shown 77.25% recognition accuracy. The following Fig. 4. shows the graphical representation of average recognition of bi-scripts using two datasets; own (CBWI-11) and CVSI-2015.

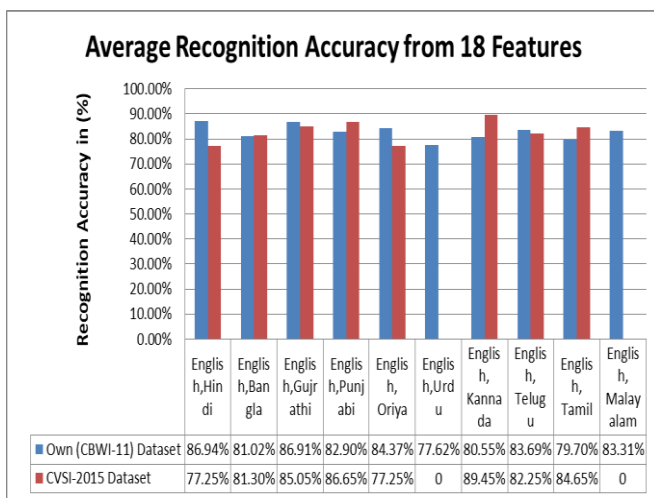


Fig 5. Average Recognitin Accuracy Using K-NN Classifier

II. CONCLUSION

This paper presents script identification from the bi-lingual document images captured from the camera for combination of 11 scripts at the word level. Tshe experiment has made by empirically setting the $n_t=3$. In K-NN classifier the different k values (3, 5, 7, and 9) and various distance metrics (Euclidean, City block, Correlation, and Cosine) are considered to obtain the highest recognition accuracy. The optimum performance is obtained for k=3 and Correlation distance metric. Hence, the proposed method has achieved the optimum recognition accuracy of 87.02%, from the English and Hindi combination with $n_t=3$ for the 7 reduced features out of 18 features and on the standard benchmark dataset; CVSI-2015 also performed well as it has 88.95% recognition accuracy which slight increased as compared to own(CBWI-11) dataset. Hence, in the future work, an attempt will be made to increase the recognition accuracy to achieve highest recognition accuracy nearing 100% for all combinations.

REFERENCES

1. Anirban Mukhopadhyay, Sourav Kumar, Souvik Roy Chowdhury, Neelotpal Chakraborty, Ayatullah Faruk Mollah, Subhadip Basu, Ram Sarkar, "Multi-Lingual Scene Text Detection Using One-Class Classifier", IJCVIP, Volume-9, pp. 48-64, 2019..
2. Ankan Kumar Bhunia, Aishik Konwer, Ayan Kumar Bhunia, Abir Bhowmick, Partha P. Roy, Umapada Pal, "Script identification in natural scene image and video frames using an attention based Convolutional-LSTM network", Pattern Recognition, Volume 85, pp. 172-184, 2019.
3. B. Shi, X. Bai, C. Yao, "Script identification in the wild via discriminative convolutional neural network", Pattern Recognition, 52, pp.448-458, 2016.
4. B.V.Dhendra, Vijayalaxmi.M.B, Gururaj Mukarambi, Mallikarjun Hangarge, "Script Identification using Discrete Curvelet Transforms", International Journal of Computer Applications, Recent Advances in Information Technology, pp.16-20, 2014
5. Costa, A. F., G. E. Humpire-Mamani, and A. J. M. Traina., "An Efficient Algorithm for Fractal Analysis of Textures", SIBGRAPI, OuroPreto, Brazil, pp.39-46, 2012
6. Gururaj Mukarambi, Satishkumar Mallappa, B.V.Dhendra, "Script Identification from Camera Based Tri-Lingual Document", IEEE, pp. 214-217, 2017..
7. Jan Zdenek, Hideki Nakayama, "Bag of Local Convolutional Triplets for Script Identification in Scene Text", ICDAR-2017, pp.369-375, 2017..
8. L. Gomez, A. Nicolaou, D. Karatzas, "Improving patch-based scene text script identification with ensembles of conjoined networks", Pattern Recognition, 67, pp.85-96, 2017.
9. Linlin Li and Chew Lim Tan, "Script Identification of Camera-based Images", IEEE, 978-1-4244-2175-6/08/, 2008
10. Lluís Gomez, Angelos Nicolaou, Dimosthenis Karatzas, "Improving patch-based scene text script identification with ensembles of conjoined networks", Pattern Recognition, pp. 85-96, 2017
11. Madhura Jajoo, Neelotpal Chakraborty, Ayatullah Faruk Mollah, Subhadi Basu, Ram Shankar, "Script Identification from Camera-Captured Multi-script Scene Text Components", AISC, pp. 159-166, 2019.
12. N. Nayef, F. Yin, I. Bizid, H. Choi, Y. Feng, D. Karatzas, Z. Luo, U. Pal, C. Rigaud, J. Chazalon, W. Khlif, "Robust Reading Challenge on Multi-Lingual Scene Text Detection and Script Identification-RRC-MLT". In Document Analysis and Recognition (ICDAR), 2017 14th IAPR International Conference IEEE, pp. 1454-1459, 2017.
13. N. Sharma, R. Mandal, R. Sharma, U. Pal, M. Blumenstein, ICDAR2015 competition on videoscript identification (CVSI 2015), In Document Analysis and Recognition (ICDAR), 2015 13th International Conference on IEEE, pp. 1196-1200, 2015.



14. O.K.Fasil,S.Manjunath,andV.N.ManjunathAradhya”Word-level script identification from scene images”, in Proceedings of the 5th International Conference on Frontiers in Intelligent Computing: Theory and Applications,pp.417-426, 2017.
15. Mallikarjun Hangarge, K.C. Santosh, Rajmohan Pardeshi, “Directional Discrete Cosine Transform for Handwritten Script Identification”, ICDAR,pp.344-348, 2013.

AUTHORS PROFILE



Dr.B.V.Dhandra is a Professor in Symbiosis Institute of Computer Studies and Research, Symbiosis International Deemed University, Pune. He obtained his PhD degree from Shivaji University, Kolhapur, India in 1993. He served as lecturer during 1979 to 1993, as a reader during 1993 to 2001 and since 2001 to 2016 he has been serving as Professor in Department of P.G. Studies and Research in Computer Science, Gulbarga University, Kalaburagi, India. He has guided 9 PhD students and 13 MPhil students. He has successfully completed 1 UGC major research project. He has published more than 150 research articles in peer reviewed national and international journals and conference proceedings. He has also authored for 2 books. His research interests are document image processing, pattern recognition and operations research.



Mr.Satishkumar Mallappa, Research Scholar in Department of Computer Science, Gulbarga University, Kalaburagi, Karnataka, India. He obtained his M.Sc degree from Gulbarga University Kalaburagi, Karnataka, India in the year 2006 and M.Phil from 2008 Gulbarga University Kalaburagi, Karnataka, India. He has published one research paper in IEEE Xplore. He has received one best research paper award. His research interests are document image processing and pattern recognition



Dr. Gururaj Mukarambi, Assistant Professor in Symbiosis Institute of Computer Studies and Research, Symbiosis International Deemed University, Pune. He is Doctorate from Gulbarga University, Gulbarga. He has 10 years of experience in Teaching and Research. He has published more than 35 research papers in peer reviewed International Journals and Conference proceedings. He has credit of published research papers in the International Journals of Elsevier, Springer and Inderscience. He has presented 15 research papers in International conferences. He has received two best research paper awards. He has been a regular reviewer and program committee member of international research conferences in India and abroad.