

# Deep learning for word-level handwritten Indic script identification

Soumya Ukil<sup>1</sup>, Swarnendu Ghosh<sup>1</sup>, Sk Md Obaidullah<sup>2</sup>, K. C. Santosh<sup>3</sup>, Kaushik Roy<sup>4</sup>,  
and Nibaran Das<sup>1</sup>

<sup>1</sup> Dept. of Computer Science & Engineering, Jadavpur University, Kolkata 700032, WB, India

<sup>2</sup> Dept. of Computer Science & Engineering, Aliah University, Kolkata 700156, WB, India

<sup>3</sup> Dept. of Computer Science, The University of South Dakota, Vermillion, SD 57069, USA

<sup>4</sup> Dept. of Computer Science & Engineering, West Bengal State University, 700126, WB, India

Corresponding authors: K. C. Santosh ([santosh.kc@usd.edu](mailto:santosh.kc@usd.edu)) & N. Das ([nibaran@gmail.com](mailto:nibaran@gmail.com))

**Abstract.** We propose a novel method that uses convolutional neural networks (CNNs) for feature extraction. Not just limited to conventional spatial domain representation, we use multilevel 2D discrete Haar wavelet transform, where image representations are scaled to a variety of different sizes. These are then used to train different CNNs to select features. To be precise, we use 10 different CNNs that select a set of 10240 features, i.e. 1024/CNN. With this, 11 different handwritten scripts are identified, where 1K words per script are used. In our test, we have achieved the maximum script identification rate of 94.73% using multi-layer perceptron (MLP). Our results outperform the state-of-the-art techniques.

**Keywords:** Convolutional neural network, deep learning, multi-layer perceptron, discrete wavelet transform, Indic script identification

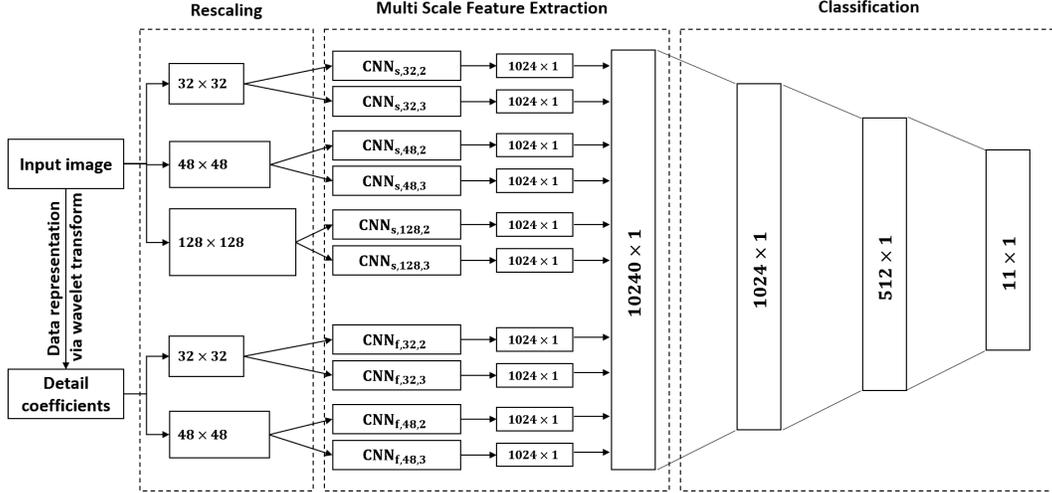
## 1 Introduction

Optical character recognition (OCR) has always been a challenging field in pattern recognition. OCR techniques are used to convert handwritten or machine printed scanned document images to machine-encoded texts. These OCR techniques are script dependent. Therefore, script identification is considered as a precursor to OCR. In particular, in case of a multi-lingual country like India script identification is the must since a single document, such as postal documents and business forms, contains several different scripts (see Fig. 1).

Indic handwritten script identification has a rich state-of-the-art literature [1–4]. More often, previous works have been focusing on word-level script identification [5]. Not stopping there, in a recent work [6], authors introduced page-level script identification performance to see whether we can expedite the processing time. In general, in their works, hand-crafted features that are based on structural and/or visual appearances (morphology-based) were used. The question is, are we just relying on what we see and use apply features accordingly or can we just let machine to select features that are required for optimal identification rate? This inspires to use deep learning, where CNNs can be used for extracting and/or selecting features for identification task(s).

Needless to say, CNNs have stood well with their immense contribution in the field of OCR. Their onset has been marked by the ground-breaking performance of CNNs on MNIST dataset [7]. Very recently, the use CNN for Indic script (Bangla character recognition) has





**Fig. 2.** Schematic block diagram of handwritten Indic script identification showing different modules: feature extraction/selection and classification.

## 2.1 CNN architecture

In general, a CNN has a layered architecture consisting of three basic types of layers namely,

1. convolutional layer (CL),
2. pooling layer (PL) and
3. fully connected layer (FCL).

CLs consist of a set of kernels that produce parameters and help in convolution operation. In CL, every kernel generates an activation map as an output. PLs do not have parameters but, their major role is to avoid possible data redundancies (still preserving their significance). In our approach, all CNNs have max-pooling operation at their corresponding PLs. In addition to these two different types of layers, FCL is used, where MLP has been in place.

In Fig. 2, we provide a complete schematic block diagram of handwritten Indic script identification showing different modules: feature extraction/selection and classification. In our study, 10 different CNNs are used to select features from a variety of representations of the studied image and we label each of them as  $CNN_{d,x,y}$ . In every  $CNN_{d,x,y}$ ,  $d$  and  $x$  respectively refer to domain representation and dimension of the studied image, and  $y$  refers to the number of convolutional and pooling layers in that particular CNN. For example, domain representation can be expressed as  $d = \{s, f\}$ , where  $s$  refers to spatial domain and  $f$ , frequency. In our case, either of these is taken into account. Note that, with the use of the *Haar wavelet transform* (HWT) (see Section 2.2), certain frequencies are removed. In case of dimension ( $x$ ), we have  $x = \{[32 \times 32], [48 \times 48], [128 \times 128]\}$ , and for simplicity,  $x = \{32, 48, 128\}$  is used. All of these dimensions signify resolution to which the input images are scaled. In case of CL and PL,  $y = \{2, 3\}$ : one of the two is taken in CNN, and  $y = 2$  means that there are two pairs of convolutional and pooling layers in the CNNs. In our model, two broad CNN architectures:  $CNN_{d,x,2}$  and  $CNN_{d,x,3}$  are used and can be summarized as in Table 1.

**Table 1.** Architecture:  $\text{CNN}_{d,x,y}$ , where  $y = \{2, 3\}$ .

Architecture	Parameter	Layer								
		CL1	PL1	CL2	PL2	CL3	PL3	FCL1	FCL2	Softmax
$\text{CNN}_{d,x,2}$	Channel	32	32	64	64	—	—	1024	512	11
	Filter size	5×5	2×2	5×5	2×2	—	—	—	—	—
	Pad size	2	—	2	—	—	—	—	—	—
$\text{CNN}_{d,x,3}$	Channel	32	32	64	64	128	128	1024	512	11
	Filter size	7×7	2×2	5×5	2×2	3×3	2×2	—	—	—
	Pad size	3	—	2	—	1	—	—	—	—

Index

CL = convolutional layer

PL = pooling layer

FCL = fully connected layer

1.  $\text{CNN}_{d,x,2}$ : We have six different layers, i.e. two CLs, two PLs and two FCLs. The first two CLs are followed by PLs. In FCLs, the first layer takes image representation that has been generated by CLs and PLs and reshapes them in the form of a vector, and the second FCL produces a set of 1024 features.
2.  $\text{CNN}_{d,x,3}$ : Every CNN has eight different layers: three CLs, three PLs and two FCLs. In general, such an architecture is very much similar to previously mentioned  $\text{CNN}_{d,x,2}$ . The difference lies in additional pair of CL and PL that follows second pair of the same. Like the  $\text{CNN}_{d,x,2}$ , these CNNs produce a set of 1024 features for any studied image.

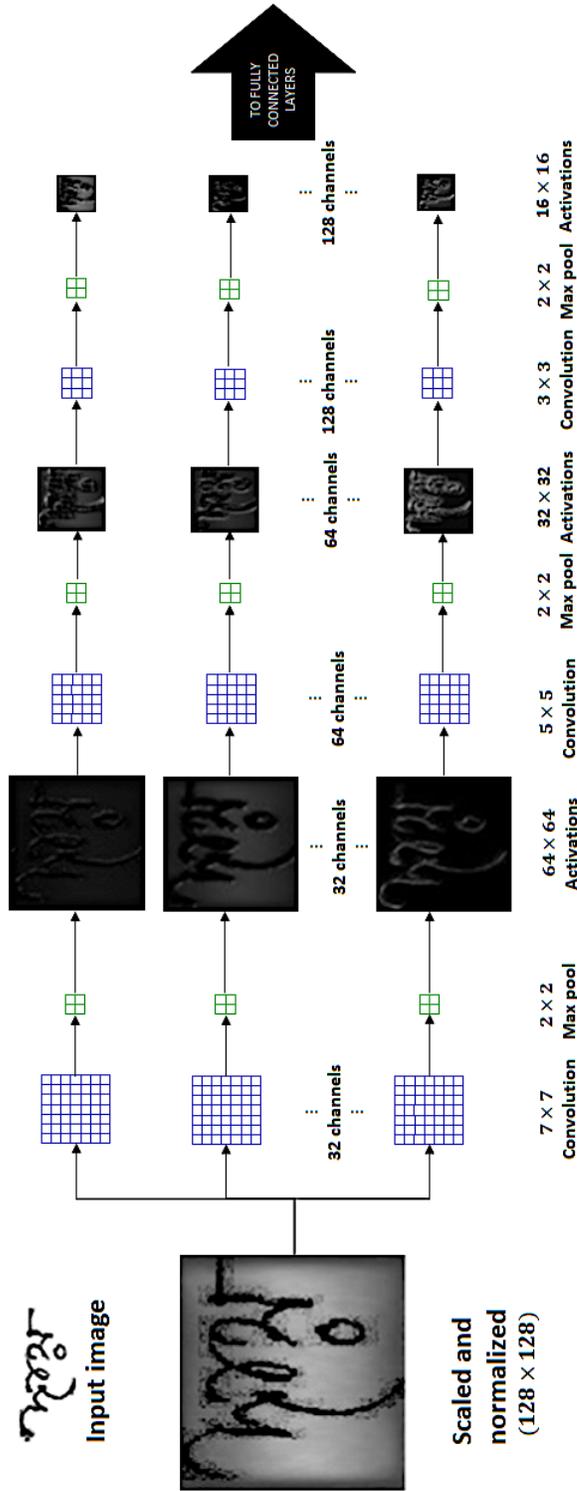
Once again, the architectural details of aforementioned CNNs are summarized in Table 1 that follows schematic block diagram of the system (see Fig. 2). For more understanding, in Fig. 3, we provide the activation maps for  $\text{CNN}_{s,128,3}$ . This means that it uses spatial domain image representation with the dimensionality of 128 and three pairs of convolutional and pooling layers in the CNNs.

## 2.2 Data representation

In general, since Fourier transform [11] may not work to successfully provide information about what frequencies are present at what time, wavelet transform (WT) [12] and short time Fourier transform [13] are typically used. Both of these help identify the frequency components present in any signal at any given time. WT, on the other hand, can provide dynamic resolution.

We consider an image a 2D time signal that can be resized. We then use multilevel 2D discrete WT on a scaled/resized image ( $128 \times 128$ ) to generate frequency domain representation. To be precise, we use the Haar wavelet [14] with seven different level decomposition that can generate approximated and detailed coefficients. Since the approximated coefficients are equivalent to zero, we use the detailed coefficients in addition to modified approximated coefficients to reconstruct the image. In the modified approximated coefficients, we consider only high frequency components.

In our method, using a variety of different WTs, such as the Daubechies [15] and several decomposition levels, the best results were observed with the Haar wavelet and a decomposition level of 7. This reconstructed image is further resized to  $32 \times 32$  ( $x = 32$ ) and  $48 \times 48$



**Fig. 3.** Illustrating the activation maps for  $CNN_{s,128,3}$ : spatial domain image representation with the dimensionality of 128 and three-layered convolutional and pooling layers.

( $x = 48$ ), and are fed into multiple CNNs as mentioned in Section 2.1. Like in the frequency domain representation, spatial domain representations are resized/scaled to  $32 \times 32$  ( $x = 32$ ),  $48 \times 48$  ( $x=48$ ) and  $128 \times 128$  ( $x = 128$ ) are fed into multiple CNNs.

### 3 Experiments

#### 3.1 Dataset, and evaluation metrics and protocol.

To evaluate our proposed approach for word-level handwritten script identification, we have considered a dataset named PHD\_Indic.11 [6]. It is composed of 11K scanned word images (grayscale) from 11 different Indic script, i.e. 1K per script. A few samples are shown in Fig. 4. For more information about dataset, we refer to recently reported work [6]. The primary reason behind considering PHD\_Indic.11 dataset in our test is, no big size data has been reported in the literature for research purpose, till this date.

Using the exact same CNN representation as before,  $C_{d,x,y}[i][j]$  represents the count where an instance with label  $i$  is classified as  $j$ . Accuracy (acc) of the particular CNN can then be computed as

$$\text{acc}_{d,x,y} = \frac{\sum_{i=1}^{11} C_{d,x,y}[i][i]}{\sum_{i=1}^{11} \sum_{j=1}^{11} C_{d,x,y}[j][i]}.$$

Precision (prec) can be computed as

$$\text{prec}_{d,x,y} = \frac{\sum_{i=1}^{11} \text{prec}_{d,x,y}^i}{11} \quad \text{and} \quad \text{prec}_{d,x,y}^i = \frac{C_{d,x,y}[i][i]}{\sum_{j=1}^{11} C_{d,x,y}[i][j]},$$

where  $\text{prec}_{d,x,y}^i$  refers to precision for any  $i$ -th label. In a similar fashion, recall (rec) can be computed as

$$\text{rec}_{d,x,y} = \frac{\sum_{i=1}^{11} \text{rec}_{d,x,y}^i}{11} \quad \text{and} \quad \text{rec}_{d,x,y}^i = \frac{C_{d,x,y}[i][i]}{\sum_{j=1}^{11} C_{d,x,y}[i][j]},$$

where  $\text{rec}_{d,x,y}^i$  refers to recall for any  $i$ -th label. Having both precision and recall, f-score can be computed as

$$\text{f-score}_{d,x,y} = \frac{\sum_{i=1}^{11} \text{f-score}_{d,x,y}^i}{11} \quad \text{and} \quad \text{f-score}_{d,x,y}^i = 2 \times \frac{\text{prec}_{d,x,y}^i \times \text{rec}_{d,x,y}^i}{\text{prec}_{d,x,y}^i + \text{rec}_{d,x,y}^i}.$$

Following conventional 4:1, i.e. train:test evaluation protocol, we have separated 8.8K images for training and the remaining 2.2K images for testing. We ran our experiments by using a machine: GTX 730 with 384 CUDA cores and 4GB GPU RAM. Besides, it has Intel Pentium Core2Quad Q6600 and 4GB RAM.

#### 3.2 Experimental set up

As mentioned earlier, are are required to train the CNNs first before testing. In other words, it is important to see how training and testing have been performed.

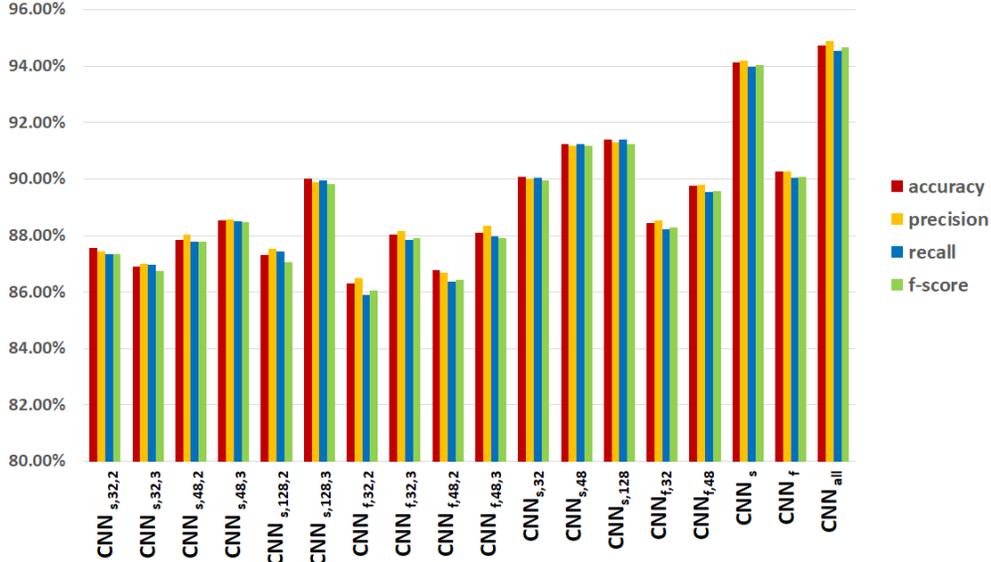
Bangla	কাজ	রাপগাছ	এক	প্রোগ্রাম	প্রকাশ
Devanagari	कर	मुश्किलों	दूरने	लडने	कामो
Gujarati	મનથી	અભે	થી ૨૨૬	ભાભાની	મીપુ
Gurumukhi	ਯੋਗੀ	ਰਿਤ	ਪੜੀ	ਮਾਰਦਾਣੀ	ਦੁਖਾਣੀ
Kannada	ಚಿಂತೆ	ನಿರೀಕ್ಷಕರ	ನೀಡ	ಗುಂಪಾನೆ	ನಿರ್ವಹಣೆ
Malayalam	അറി	ചെയ്	രരന്ന	തന്ന	അവര
Oriya	କୋଳା	ଭାବି	ସାଧାରଣ	ଉତ୍ତର	ଦୁଇ
Roman	Labour's	charger	African	went	seeking
Tamil	உயி	அவை	கிணிய	கிணிய	காணீ
Telegu	అది	అనంది	మమ్మి	మీ అన	ఉచిత
Urdu	نان لکیرا	راشکانه ها	ایران مهند	بسی اقامت	نه ترانهای

Fig. 4. Illustrating few samples from dataset named PHDIndic\_11 used in our experiment.

Our CNNs in this study represented by  $CNN_{d,x,y}$  are trained independently using the training dataset (as mentioned earlier). To clarify once again, these CNNs has either two or three pairs of consecutive convolutional and pooling layers. Besides, each of them has three fully connected layers. The first of these three layers function as an input layer with number of neurons that depends on the size of the input image specified by  $x$ . The second layer in CNNs has 1024 neurons, and during training, we apply a dropout probability of 0.5. The final layer has 11 neurons, whose outputs are used as input to an 11-way soft-max classifier that provides us with classification/identification probabilities for each of 11 possible classes. Our training set of 8.8K word images are split into multiple batches of 50 word images and the CNNs are trained accordingly. For optimization, Adam optimizer [16] was used with learning rate of  $1 \times 10^{-3}$  having default parameters:  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . This helps apply gradients for loss to the weight parameters during back propagation. We computed the accuracy of the CNN as training proceeds by taking ratio of the of images successfully classified in the batch to the total number of images in the studied batch.

After training CNNs with 8.8K word images, we evaluated/tested each of them independently with the test set that is composed of 2.2K word images.

More specifically, for each input size specification, we have two CNNs, i.e. for domain (d) and input size ( $x \times x$ ):  $CNN_{d,x,2}$  and  $CNN_{d,x,3}$ . Altogether, we have 10 different CNNs since we have three different input sizes ( $x = \{32, 48, 128\}$ ) for raw image and two different input sizes ( $x = \{32, 48\}$ ) for wavelet transformed image/data. For better understanding, we refer readers to Fig. 2. Note that we have trained the CNNs to extract 1024 features from each one, i.e.  $120 \times 1$ . These are then concatenated to form a single  $10240 \times 1$  vector, where ten different CNNs are employed. Like in the conventional machine learning classification, these features are used for training and testing purpose using MLP classifier.



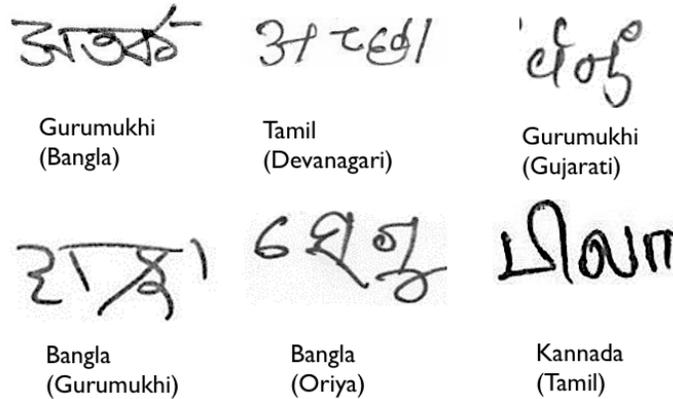
**Fig. 5.** Our results (in terms of accuracy, precision, recall and f-score) for all networks: CNNs and their possible combinations.

### 3.3 Our results and comparative study

In this section, using dataset, and evaluation metrics (see Section 3.1), and experimental setup (see Section 3.2), we summarize our results and comparative study as follows:

1. We provide results that have been produced from different architectures (CNNs), and select the highest script identification rate from them; and
2. We then take highest script identification for a comparative study, where previous relevant works are considered.

**Our results:** Fig. 5 shows the comparison of the individual CNNs along with the effect of combining them. The individual CNN<sub>d,x,y</sub> produced the maximum script identification rate of 90% for CNN<sub>s,128,3</sub>. As we ensemble two- and three-layered networks of corresponding domain (d) and input size (x), we observe a positive correlation between input size and accuracy. CNN<sub>s,128</sub> provides the maximum script identification rate of 91.41% in this category. The primary reason behind the increase in accuracy is that the ensemble of two- and three-layered networks suggests that these networks complement each other. Further, to study an effect of the spatial (s) and frequency (f) domain representation in CNN<sub>d,x,y</sub>, we ensemble networks across all input sizes and depth of network. The spatial representation, CNN<sub>s</sub>, have produced script identification rate of 94.14% and the frequency domain representation, CNN<sub>f</sub>, have escalated up to 90.27%. However, the frequency domain representations learning can be complimented and it has been clearly seen when we combined them. In their combination, we have achieved the highest script identification rate of 94.73%. Since we



**Fig. 6.** Misclassified samples, where script names in the bracket are the actual scripts but, our system identified them incorrectly. For example, in the first case, word image has been identified as Gurumukhi and it actually is Bangla.

have not received 100% script identification rate, it is wise to provide a few samples where our system failed to identify correctly (see Fig. 6).

Like we have mentioned in Section 3.1, we also provide precision, recall and f-score for all architectures in Fig. 5. In what follows, the highest script identification rate, i.e. 94.73% will be taken for a comparison.

**Comparative study:** For a fair comparison, widely used deep learning methods, such as LeNet [7] and AlexNet [9] were taken. In addition, recently reported work on 11 handwritten Indic script dataset [6] (including their baseline results) was considered. In Table 2, we summarize the results. Our comparative study is focused on accuracy (not precision, recall and f-score), since other methods reported accuracy, i.e. identification rate. Of course, in Fig. 5, we are not just limited to accuracy.

In Table 2, our method outperforms all other methods. Precisely, it outperforms Obaidullah et al. [6] by 4.7%, LeNet [7] by 2.73% and AlexNet [9] by 2.61%.

## 4 Conclusion

In this paper, we have proposed a novel framework that uses convolutional neural networks (CNNs) for feature extraction. In our method, in addition, to conventional spatial domain representation, we have used multilevel 2D discrete Haar wavelet transform, where image representations have been scaled to a variety of different sizes. Having these, several different CNNs have been used to select features. With this, 11 different handwritten scripts: Bangla, Devnagari, Gujarati, Gurumukhi, Kannada, Malayalam, Oriya, Roman, Tamil, Telugu and Urdu, have been identified, where 1000 words per script are used. In our test, we have achieved the maximum script identification rate of 94.73% using multi-layer perceptron (MLP). To the best of our knowledge, this is the biggest data for Indic script identification work. Considering the complexity and the size of the dataset, our method outperforms the previously reported techniques.

**Table 2.** Comparative study.

Method	Accuracy
Obaidullah et al. [6] (Hand-crafted features)	91.00%
LeNet [7] (CNN)	82.00%
AlexNet [9] (CNN)	92.14%
Our method (Multiscale CNN + WT)	94.73%

## References

1. Ghosh, D., Dube, T., Shivaprasad, A.: Script recognitiona review. *IEEE Transactions on pattern analysis and machine intelligence* **32**(12) (2010) 2142–2161
2. Pal, U., Jayadevan, R., Sharma, N.: Handwriting recognition in indian regional scripts: a survey of offline techniques. *ACM Transactions on Asian Language Information Processing (TALIP)* **11**(1) (2012) 1
3. Singh, P.K., Sarkar, R., Nasipuri, M., Doermann, D.: Word-level script identification for hand-written indic scripts. In: *Document Analysis and Recognition (ICDAR), 2015 13th International Conference on*, IEEE (2015) 1106–1110
4. Hangarge, M., Santosh, K., Pardeshi, R.: Directional discrete cosine transform for handwritten script identification. In: *Document Analysis and Recognition (ICDAR), 2013 12th International Conference on*, IEEE (2013) 344–348
5. Pati, P.B., Ramakrishnan, A.: Word level multi-script identification. *Pattern Recognition Letters* **29**(9) (2008) 1218 – 1229
6. Obaidullah, S.M., Halder, C., Santosh, K., Das, N., Roy, K.: Phdindic.11: page-level hand-written document image dataset of 11 official indic scripts for script identification. *Multimedia Tools and Applications* (2017) 1–36
7. LeCun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. *Proceedings of the IEEE* **86**(11) (1998) 2278–2324
8. Roy, S., Das, N., Kundu, M., Nasipuri, M.: Handwritten isolated bangla compound character recognition: A new benchmark using a novel deep learning approach. *Pattern Recognition Letters* **90** (2017) 15–21
9. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: *Advances in neural information processing systems*. (2012) 1097–1105
10. Sarkhel, R., Das, N., Das, A., Kundu, M., Nasipuri, M.: A multi-scale deep quad tree based feature extraction method for the recognition of isolated handwritten characters of popular indic scripts. *Pattern Recognition* (2017)
11. Smith, S.W., et al.: *The scientist and engineer’s guide to digital signal processing*. (1997)
12. Daubechies, I.: The wavelet transform, time-frequency localization and signal analysis. *IEEE transactions on information theory* **36**(5) (1990) 961–1005
13. Portnoff, M.: Time-frequency representation of digital signals and systems based on short-time fourier analysis. *IEEE Transactions on Acoustics, Speech, and Signal Processing* **28**(1) (1980) 55–69
14. Sundararajan, D.: *Fundamentals of the discrete haar wavelet transform*. (2011)
15. Vonesch, C., Blu, T., Unser, M.: Generalized daubechies wavelet families. *IEEE Transactions on Signal Processing* **55**(9) (2007) 4415–4429

16. Kingma, D., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)