

# Creation of data resources and evaluation tool for multi-lingual OCR

S. Setlur      S. Kompalli      R. Vemulapati      V. Govindaraju

Center of Excellence for Document Analysis and Recognition,  
University at Buffalo, Buffalo, NY - 14260  
{kompalli, setlur, govind, raman-v}@cedar.buffalo.edu

## Abstract

*Indic scripts contribute to a large number of languages and dialects used by people all over the world with Devanagari being most widely used. To date, OCR techniques of different Indic scripts have not yet been widely disseminated or evaluated independently. Automated evaluation tools are currently not available for lack of a standard representation of ground-truth and result data.*

*A key reason for the absence of sustained research efforts in off-line Devanagari OCR appears to be the paucity of data resources. Ground truthed data for words and characters, on-line dictionaries, corpora of text documents and reliable, standardized statistical analysis and evaluation tools are currently lacking. So, the creation of such data resources will undoubtedly provide a much needed fillip to researchers working on Devanagari OCR.*

*This paper describes a National Science Foundation sponsored project under the International Digital Libraries program to create data resources that will facilitate development of Devanagari OCR technology and provide a standardized test bed and evaluation tools for Devanagari script recognition.<sup>1</sup>*

## 1 Introduction

Automatic text recognition using OCR is the process of converting the image of a textual document into its digital textual equivalent, the advantage being that the textual material can be edited. The document image itself can be either machine-printed or handwritten, or a combination of the two.

The recognition of any unit of a document is in the context of the rest of the document and can benefit from the context. For all practical purposes, the units of recognition range from a single character

---

<sup>1</sup>This material is based upon work supported by the National Science Foundation under Grant No. 0112059. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

to a complete sentence. While character images are intuitively the default domain in which OCR's operate, grammatical and linguistic constraints from the sentence level can drive the recognition of the characters and words [1, 12].

The development of OCRs requires vast amounts of training data with ground-truth to achieve acceptable levels of performance. Manual truthing of data is expensive and labor intensive. Moreover, a standardized, public-domain truth set is essential to to justly compare the relative performance of different OCRs. Such data is sorely lacking for Indic Script OCR researchers at present. In this paper, we discuss the architecture of a system to collect data for *Devanagari*, the most widely used Indic Script. Section 2 provides a brief overview of the *Devanagari* script. Section 3 reviews some *Devanagari* OCR techniques that have been published. Section 4 describes the data collection and ground-truthing methodology and the tool designed for this purpose. Section 6 describes the future work planned under this initiative.

## 2 Overview of Devanagari script

The basic character set of *Devanagari* script is made up of 16 vowels (Table 1) and 36 basic consonants (Table 2). A vowel following a consonant takes a modified shape, which is placed to the left, right, top or bottom of the consonant. Each of the consonants combines with vowel modifiers for 15 of the 16 vowels to create new character shapes (Table 1). More than one consonant can combine with each other to form conjunct consonants (Some examples can be seen in Table 3). In addition to these, there are also a few special characters, which are shown in Table 2 and 10 numerals shown in Table 2. Thus, the total number of character shapes encountered in *Devanagari* script is well over 300 and it is evident from looking at the characters that there would be a significant confusion between similar characters. A word in *Devanagari* script consists of composite charac-

Table 1: *Devanagari* vowels and their corresponding consonant modifiers

Vowel	Modifier	Modified Consonant	Vowel	Modifier	Modified Consonant
अ		क	ऌ	३	
आ	ा	का	ॡ	३	
इ	ि	कि	ए	॑	के
ई	ी	की	ऐ	॑	कै
उ	ु	कु	ओ	ो	को
ऊ	ू	कू	औ	ौ	कौ
ऋ	ॠ	कृ	अं	ं	कं
ॠ	ॡ	कृ	अः	ः	कः

Table 2: Simple Consonants (with transliteration), Digits and Special Characters in *Devanagari*

Simple Consonants										
क	ख	ग	घ	ङ	च	छ	ज	झ	ञ	
ka	kha	ga	gha	ṅa	cha	Cha	ja	jha	Jṅa	
ट	ठ	ड	ढ	ण	त	थ	द	ध	न	
Ta	Tha	Da	Dha	Ṇa	ta	tha	da	dha	na	
प	फ	ब	भ	म						
pa	pha	ba	bha	ma						
य	र	ल	व	श	ष	स	ह	ळ	क्ष	ज्ञ
ya	ra	la	va	sha	Sha	sa	ha	La	xa	Gya
Digits										
०	१	२	३	४	५	६	७	८	९	

Special Characters	
◌ं	Chandrabindu
◌ः	Avagraha
	Danda
॥	Paragraph end
ॐ	AUM

Table 3: Examples of conjunct consonant shapes in *Devanagari*

Consonants	Conjunct form	Consonants	Conjunct form	Consonants	Conjunct form
क त	क्त	द द	द्द	ल ल	ल्ल
क य	क्य	द ध	द्व	श च	श्च
क व	क्व	द न	द्वन	श न	श्न
घ न	घ्न	द म	द्वम	श र	श्चर
ङ क	ङ्क	द य	द्वय	श ल	श्चल
ङ ग	ङ्ग	द व	द्वव	श व	श्चव
ञ च	ञ्च	ध न	ध्न	ष ट	ष्ट
ट ट	ट्ट	न न	न्न	ह न	ह्न
त त	त्त	प त	पत्त	ह य	ह्य
त न	त्न	प न	पन्न	ह र	ह्र
त र	त्र	प ल	प्ल	ह ल	ह्ल
द ग	द्ग	म ल	म्ल	ह व	ह्व

Table 4: Indic scripts

Script	Official	Transliteration	Unicode	Vowels	Consonants	Vowel Signs	Spl Characters
Benagali	Yes	Yes	Yes	11	32	11	25
Devanagari	Yes	Yes	Yes	16	37	14	27
Gujarati	Yes	Yes	Yes	13	34	13	8
Kannada	Yes	Yes	Yes	14	35	13	8
Tamil	Yes	Yes	Yes	12	22	11	7
Telugu	Yes	Yes	Yes	14	35	13	8
Gurumukhi	Yes	Yes	Yes	10	32	10	11

ters joined by a horizontal line at the top which is referred to as the header line or *shirorekha*. A single or double vertical line called *danda* is traditionally used to indicate end of phrase or end of sentence. In modern practice, inter-word spacing and use of Latin punctuation prevail. Likewise, although *Devanagari* has a native set of symbols for numerals, Arabic numerals are sometimes used in modern documents.

## 2.1 Similarity of other Indic Scripts to Devanagari Script

Bengali, Gujarati, and Gurumukhi (Punjabi) along with Devanagari (Hindi) are classified into a group called *Aryan* scripts and hold similar characteristics with respect to their rendering, formation of compound/conjunct characters and modifiers. Tamil, Telugu and Kannada are classified as *Dravidian* scripts and their character shapes are distinct from those of *Aryan* scripts. Though there exist some differences between the two groups of scripts, for example - fewer consonants and lack of consonant forms in Tamil, due to the structural similarity of the languages, similar technologies support several Indic Scripts (Table 4). The Multi-Lingual truthing tool has been designed taking advantage of this similarity, leading to a single technique for handling the truthing of documents in these scripts.

## 3 Devanagari OCR - Background

We surveyed the available *Devanagari* OCR literature to study the techniques used by researchers in the area to determine the nature of data we would be required to collect for the test bed. *Devanagari* text can have compound characters, characters with vowel modifiers and half characters. Due to the complexity of the character set, Researchers pre-process *Devanagari* characters in different ways. While some attempt to separate the vowel modifiers and the core characters [18], recognizing them separately, others attempt to recognize each compound consonant or consonant with vowel modifier as a separate class [19, 14].

Research [24] into compound and fused character segmentation of Devanagari script indicates an accuracy level of 85%, using a hybrid, multi-pass approach that relies on structural features. Many other techniques [11, 26] also exist that give similar results, but using different algorithms. There is a strong indication that the character-segmentation level has significant issues, and algorithms handling those issues need to be analysed and tested independently.

Researchers [18] using structural features and implementing a blackboard architecture have shown solutions that provide about 70% accuracy on *Devanagari* text containing unknown font and about 80% accuracy for text with a known font. Statistical and Structural features, starting from low-level parameters like average word and character dimensions to higher level information like stroke structures were extracted from the document image. These were taken together and further analyzed to obtain improved results; identifying first the base characters, then the compound characters etc. Thus, it was a kind of re-entrant algorithm. Other researchers [14], used a top-down architecture, to achieve about 96% accuracy at the character recognition level. Various structural features of *Devanagari* script were used to create a heuristic and obtain character stroke features. These feature sets were then processed using a Decision-tree classifier. Other techniques include a Binary tree classifier based on basic primitive constructs for the script [9], and syntactic pattern analysis systems, which used a picture language for *Devanagari* OCR [10].

Current research in *Devanagari* OCR does not indicate a single view of what the granularity of a character class should be, with some adopting the notion of treating the base character distinct from the modifier [8, 18] and others treating the base character with modifiers and the conjunct characters as distinct classes [11, 14].

Given these issues, it was decided to collect word level data and ground-truth as a start and words would be later segmented into their component

characters and ground-truthed as necessary. This is a significant departure from previous databases [3, 2, 7, 6] which mostly provided character level data only. However, most of the scripts listed in section 3.1 by nature, do not have compound characters and construction of their words usually involves simple, sequential ordering of individual characters.

### 3.1 Current OCR databases

Current OCR training data sets provide samples of all characters in the script under consideration. A few of the databases listed below are indicative that such an exhaustive coverage of all characters is not limited to Latin-character based scripts [3, 2], but is retained across other, complicated scripts. [7, 6]

- UW English Document Image Database I  
1147 document page images from English Scientific and Technical Journals, with text zone bounds, ground truth data for each zone, finer attributes and qualitative information. [3]
- Handwritten Korean character image database  
Handwritten Cursive Hangul character images - 2,350 classes \* 100 samples, 100dpi grayscale/binary images. [7]
- UW-II English/Japanese Document Image Database  
624 English pages from journals, 477 Japanese pages from journals, with text zone bounds, ground truth data for each zone, finer attributes and qualitative information. [4]
- ERIM Arabic Document Database  
Over 750 pages of Arabic text at 300dpi, all character information stored as Unicode, Truth for characters as well as ligatures, Over 200 distinct Arabic ligatures Approximately 1,000,000 characters. [5]
- Japanese Character Image Database  
Approximately 180,000 Kanji, Hiragana, Katakana, alphanumeric, and symbolic characters extracted from a variety of machine printed documents, with quality varying from clean to degraded digitized as binary tiff images at 400 dpi on a flatbed scanner. More than 3,300 character categories in JIS level-0 and level-1. Every image supplied with a manually determined truth-value (4-byte JIS code). [6]

### 3.2 Other data Visualization tools

This section summarizes the tools that were studied in order to ensure platform independence and ease of user interaction. The survey included Multi-Lingual

text editors, OCR's and other data visualization toolkits.

- iLEAP, a Indic language word processor  
Developed by CDAC, India, iLEAP is an Intelligent, Internet ready, Indian language Word Processor on Windows. The desktop view of iLEAP is self explanatory, with language-sensitive tool layout and operation. When the choice of language is Hindi, the entire context of the system, including keyboard layouts, available fonts, and the word reference is changed to Hindi. It is possible to switch between the languages, and hence the context of operation at the click of a mouse.
- Chitrakan, a *Devanagari* OCR  
A *Devanagari* OCR package developed by CDAC, India, Chitrakan provides Text and Picture segmentation for grayscale/color, .bmp/.tiff images scanned at 300 dpi and gives results in ISCII or .RTF format. In addition, rotation, color inversion and file printing, image editing tools are provided.
- Transcriber, for speech signal segmentation  
Transcriber is a tool for segmenting, labeling and transcribing speech and provides management of several layers of segmentation, for different applications. It synchronizes the cursors in the text editor and the signal window, ensuring always that control remains on the signal that is being processed. Transcriptions are stored in XML format for easier automatic processing and exchange; syntactic validation of files upon their DTD. Transcriptions are encoded using the most current encodings including Unicode.
- ViPER Video Processing Evaluation Resource  
ViPER is a toolkit to automate video processing evaluation tasks, aimed at developing a ground truth representation which can be maintained and manipulated by a number of different tools and applications. The system is modularized into three parts - ViPER-GT, for creating frame-level ground truth; ViPER-PE, providing performance evaluation of any automated video analysis system; and ViPER-Viz for visualization that contains programs which enable a user to compactly visualize ground truth results, and the performance evaluation results.
- GATE - General Architecture for text encoding.  
GATE comprises an architecture, development environment, created by the Sheffield NLP group, and has been used in many language

processing projects, particularly for Information Extraction in many languages. GATE includes resources for common Language Engineering data structures and algorithms, including documents, corpora and various annotation types, a set of language analysis components for Information Extraction and a range of data visualisation and editing components. Input can be documents in a variety of formats including XML, RTF, email, HTML, SGML and plain text, and are converted into a single unified model of annotation.

## 4 Collection of Devanagari documents

Character and word training data with ground truth from documents that are representative of those encountered in the real world are basic to the development of recognition systems. However, collection and truing of data is a laborious and expensive process. Designing an evaluation test bed for recognition systems also requires creation of test sets that account for most, if not all of the variability seen in the real world [17].

Simulated data has often been proposed as an alternative to real data for OCR evaluation. An analysis of Asian and Latin script separation using simulated data [22] indicated that at the level of separating languages, degradation models would work well. Also, research has been done in the direction of multilingual OCR evaluation using multiple, multilingual versions of a single document set (the bible) as groundtruth [23]. Multiple copies were obtained from a computer generated common dataset by using degradation models.

Earlier research [21] describes a Morphological degradation model that takes into consideration pixel inversion and blurring as modeling parameters. Another degradation model [20] relies on image and script features such as Resolution, Size, Scaling, Blur, Offset, Skew, Threshold, Jitter and Sensitivity as salient features suitable for generating simulated documents.

A survey of feature extraction techniques for OCR [25] indicates that OCR researchers make use of a very rich set of feature vectors, encompassing both variant features as well as invariant features like pixel density, rotation invariant Zernike Moments, Geometric Moment Invariants, Normalized projection histograms etc. Since different models operate on different image properties, and because different feature vectors are used by classifiers for different scripts, it is not likely that all simulation models would work equally well for OCRs of all scripts. To the best knowledge of the authors, a detailed analysis of simulation models in the lines of work done

for Latin or Japanese [22, 20] has not yet been attempted for *Devanagari*. It would be inaccurate to assume that the same set of degradation parameters would work equally well for testing Devanagari OCR's. The design of the truing tool has been such that various degradation parameters for the documents can be measured and recorded along with the truth files. This data can be analyzed at a later stage to determine any new degradation parameters for *Devanagari*. We hope that an analysis of such real-life data for *Devanagari* would go a long way in determining a more standard technique of identifying script-specific and script-independent simulation parameters.

By virtue of the rich character set of *Devanagari*, it would be infeasible to collect document images containing samples of every character. This is also because some of the character forms in *Devanagari* are not commonly used. To cover most of the character set, we are collecting a wide range of documents, ranging from Sanskrit scriptures to contemporary novels, magazines and newspapers, with varying print quality and font. To this end, 400 printed document pages in Hindi and Sanskrit are being scanned. Estimating an average of 300 words per document, a data set of 120,000 words would be created.

## 5 Truing tool Design

We have developed a framework (see fig: 1) for truing based upon our previous experience in designing ergonomic truing tools, to handle the new challenges posed by *Devanagari* document truing,

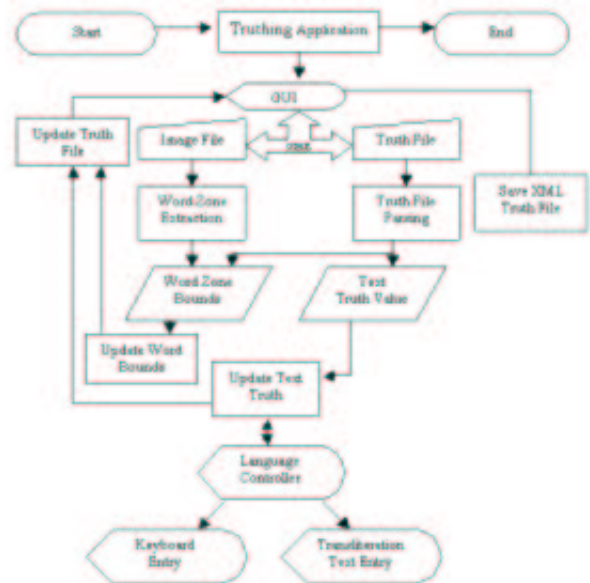


Figure 1: Design of the system

An interface for ground-truthing documents typically should be able to display the document image and segments such as words and characters and allow the truther to record the ground-truth of the document segments. The truth should be recorded using a standard representation. It is our intent that *Devanagari* OCR developers will provide results of their system on the standard test set being created, using the same representation as the truth and use the evaluation tools developed to evaluate the performance of their system. This will allow multiple OCRs to be bench-marked and assessed.

## 5.1 Devanagari Truth Representation

To satisfy the requirements outlined in section 3, it was decided to represent the truth in a way that would allow the user to query both word level and glyph level information. Unicode allows us to represent most *Devanagari* characters - right from the basic glyphs, to more complicated composite/compound forms.

Unicode has emerged as the de-facto standard for encoding characters of all languages and hence was selected as the ideal standard to represent both the ground-truth as well as the results of recognition systems. Unicode uses a hexadecimal code for representation of the basic character, and the modifiers.

$$\begin{array}{l}
 JA_d + NYA_l \rightarrow J.NYA_n \\
 ज् + ज \rightarrow ज्ञ \\
 RA_d + GA_d + GHA_l \rightarrow GA_n + GHA_l + RA_{sup} \quad \text{Displayed Output} \\
 र् + ग् + घ \rightarrow र + घ + ँ \rightarrow र्घ
 \end{array}$$

Figure 2: Compound character representation in Unicode

Detailed rules (fig:2 shows two rules) exist to generate characters with variations like the viraama, compound characters, characters with modifiers etc [15]. These rules specify the character that needs to be rendered when a given string of unicode values is provided.

Though there still exist rare characters that cannot be represented using Unicode, in most cases it provides a bijective relation between the character code and the compound character. Representations like transliteration are very elegant, phonetic and more intuitive, but they give an injective relation from the compound character code set to the characters. Representing the truth in a standard format such as Unicode allows results and truth to be easily compared.

However, while the representation of the truth

would be in Unicode, the process of generating the truth would require the use of a more intuitive interface.

## 5.2 Transliteration Scheme - An efficient technique for Devanagari input

To access a character from a font, one needs a way to specify that character in the input source. This is trivial in the case of the Latin alphabet where you can just type **A** on the keyboard when you want to get a typeset **A** in the output and the data can be saved as an ASCII string. In the case of non-Latin character sets, however, one would need either a keyboard designed for the specific script (hard to obtain and difficult to train users) or a transliteration scheme (usually phonetic) that would be able to map ASCII strings to corresponding non-Latin characters. So, to get the Devanagari character क, one could type in the phonetically equivalent (in English) ASCII string **ka**. Since the sounds associated with many Devanagari characters cannot be represented in phonetic English, approximations such as **RRi** for ऋ are used. Also, many sounds can be phonetically represented in English in more than one way. This has led to people adopting a variety of different transliteration schemes and no single standard has found universal acceptance. However, to the best knowledge of the authors, among the available transliteration schemes, the ITRANS technique supports the most number of languages and hence, was deemed more popular than other schemes. It is also a fairly intuitive scheme that can be learnt quickly.

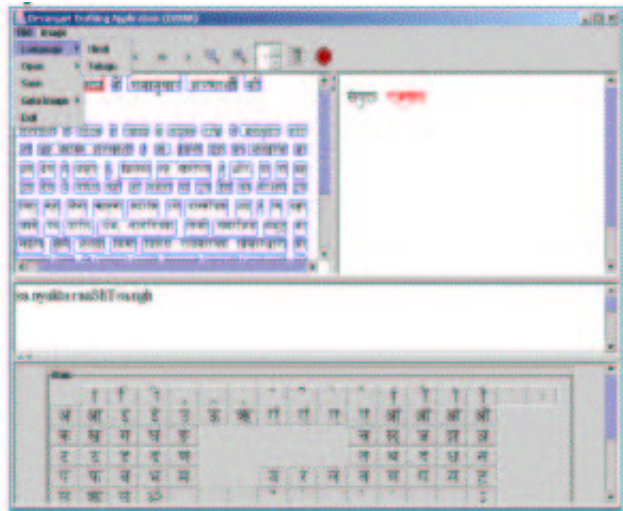


Figure 3: GUI that is used in order to truth the documents

Using transliteration as a truthing aid enables

users to record the truth using a phonetic representation of the *Devanagari* text and the transliterated text can be parsed and converted into the corresponding Unicode for storage. A tool for truthing *Devanagari* documents was designed. Figure 3 shows a screen dump of the truthing tool interface.

### 5.3 Image Processing

The tool accepts images in a wide variety of popular image formats including GIF, TIFF and JPEG. Some basic image processing tools such as line and word separation are integrated with the tool. Controls are provided in the tool bar to change thresholds for the word separation. Word boundaries can also be manually selected to correct errors if any, in the automated process. Users can also demarcate non-*Devanagari* or spurious components on the image. Once the words are satisfactorily separated, the image is ready for truthing.

Line separation is based on average pixel intensity (as measured by the grayscale values) in each row. First, sum of the grayscale values along each row of the image is calculated. Next, the average grey scale value is obtained by dividing this sum with the number of pixels (in that row) whose grey scale value is greater than 0. Regions between 2 successive minima of these intensity values are considered to be the line boundaries. After obtaining line information from the image, each line is considered separately for word separation. The procedure for word separation is very similar to the one outlined for line separation, except that the sum is taken along each column of the line. Spurious lines/words are eliminated by using a predetermined threshold for the height/width. The advantage of this method (over simple summing of the grayscale values along each row/column) is that it produces much sharper histogram profiles in the presence of uniform background noise.

The keyers view the segmented document image in the box at the top left of the screen. The truth can either be input using the simulated keyboard in the bottom panel or by typing in transliterated text using the ITRANS convention. The truth is stored along with the word co-ordinates in Unicode format as described earlier. The truth will be rendered from the Unicode values using a *Devanagari* font in the top right panel to give the keyer visual feedback so that the keyer can compare the truth and the *Devanagari* text side-by-side. The current word being truthed is highlighted both in the image as well as in the rendered truth.

### 5.4 Truth File Format

The truth will be stored in XML format, which is a very popular markup language for representing data in an organized fashion. Commercial as well as pub-

lic domain parsers are available to query collections of XML documents. The use of XML would provide researchers with a means to design parsers and query engines tailored to their specific needs.

```

<truth documentImage=".\\gaban001.tif">
- <physicalInfo>
  <skew angle="0" />
</physicalInfo>
- <linguistics>
  <script value="Devanagari" />
</linguistics>
- <word number="0">
  <bounds x1="713" y1="250" x2="950" y2="320" LineNumber="0" />
  <char>2327</char>
  <char>2348</char>
  <char>2344</char>
  <char>10</char>
  <transliteration>gaban</transliteration>
</word>
</truth>

```

गबन

Figure 4: Truth File Format

Currently, the truth file for each document contains the image name, script of the document, the skew angle, word boundaries and two different formats of truth for each word - Unicode and ITRANS Transliteration. Future ideas are to identify and include other relevant physical and linguistic features, like document degradation index, document source information etc.

The tool exploits the Platform Independence of Java and has a modular design; allowing scaling up to other scripts. The tool has been Currently implemented for Hindi, Telugu and English. The data set contains Document Data, for OCR purposes, with truth files in the popular XML format. We also propose to provide query tools to retrieve word images based on a search string in transliterated text as well as Unicode values.

## 6 Conclusions and future work

The initial datasets, consisting of documents with their word-level truth have been made freely available to the public in early 2003. Distribution via the world wide web is being planned. The truthing tool and query engines will also be made available to the public. We would like to acknowledge the assistance provided by Prof. B.B.Chaudhuri and U. Pal of the Indian Statistical Institute, Calcutta in the design of the composition of the data sets.

The following data resources will also be created during the course of this effort.

- Create standard training and testing sets of machine-printed and handwritten . *Devanagari* characters for training and evaluation of character recognizers with ground-truth.
- Create and/or consolidate any existing machine readable dictionaries for Sanskrit and Hindi words that can be used for various Natural Language Processing (NLP) applications.

- Create and/or consolidate any machine readable corpora of texts for Sanskrit and Hindi that can be used for various NLP applications.
- Develop software for grammatical tagging of corpora for use in various NLP applications.

## References

- [1] Hanhong Xue, Venu Govindaraju *On the Dependence of Handwritten Word Recognizers on Lexicons*. IEEE Pattern Analysis and Machine Intelligence. Vol 24, No 12. pp. 1553-1564.
- [2] J. J. Hull. *A database for handwritten text recognition research*. IEEE Transaction on Pattern Analysis and Machine Intelligence, 16(5):550-554, May 1994.
- [3] Isabelle Guyon, Robert M. Haralick, et al. *Data Sets For OCR And Document Image Understanding Research* Proceedings of the 2nd International Conference on Document Analysis and Recognition.
- [4] Haralick Robert M *UW-II English/Japanese Document Image Database*. Intelligent Systems Laboratory, University of Washington.
- [5] Schlosser Steven G. *ERIM Arabic Database* Document Processing Research Program. Information and Materials Applications Laboratory, Environmental Research Institute of Michigan.
- [6] *Japanese Character Image Database* Center for Excellence in Document Analysis and Research,
- [7] D.H. Kim, Y.S. Hwang, S.T. Park, et al. *Handwritten korean character image database* Proceedings of the 2nd International Conference on Document Analysis and Recognition.
- [8] Veena Bansal, R.M.K. Sinha. *On How to Describe Shapes of Devanagari Characters and Use Them for Recognition*. Fifth International Conference on Document Analysis and Recognition.
- [9] Sethi, I.K., Chatterjee, B., *Machine Recognition of Constrained Hand Printed Devanagari*. Pattern Recognition(9), 1977, pp. 69-75.
- [10] R.M.K. Sinha, Mahabala H. *Recognition of Devanagari Script*. IEEE Trans on Systems, Man, and Cybernetics, Vol. 9, pp 435-441, 1979.
- [11] A. Bishnu, B.B. Chaudhuri. *Segmentation of Bangla Handwritten Text into Characters by Recursive Contour Following* Fifth International Conference on Document Analysis and Recognition.
- [12] J.C. Perez-Cortes, J.C. Amengual, J. Arlandis, and R. Llobet. *Stochastic error correcting parsing for ocr post-processing*. In *International Conference on Pattern Recognition ICPR-2000*, Barcelona, Spain.
- [13] *Online Database of languages: www.ethnologue.com*
- [14] Chaudhuri, B.B.; Pal, U. *An OCR system to read two Indian language scripts: Bangla and Devanagari*. Proceedings of the 4th International Conference on Document Analysis and Recognition (ICDAR '97) Volume: 2 , 1997 Page(s): 1011 -1015 vol.2
- [15] *The Unicode Standard Version 3.0*. Addison-Wesley, 2000. ISBN 0-201-61633-5.
- [16] S. Setlur and V. Govindaraju, *Translingual OCR by Template Correlations* Proceedings of SPIE, Vol. 3964, 1999, pp. 305-312.
- [17] S. Setlur, A. Lawson, V. Govindaraju, and S. Srihari, *Large scale address recognition systems - Truthing, testing, tools and other evaluation issues*, International Journal of Document Analysis and Recognition, Vol. 4, No. 3, March 2002, pp. 154-169.
- [18] V. Bansal, R.M.K. Sinha, *Integrating Knowledge Sources in Devanagari Text Recognition*, IEEE Transactions on Systems, Man and Cybernetics, Vol. 30, No. 4, pp. 500-505, 2000.
- [19] B. Chaudhuri and U. Pal, *A complete printed Bangala OCR system*, Journal of Pattern Recognition, Vol. 31, No. 5, pp. 531-549, 1998.
- [20] T. K. Ho and Henry S. Baird. *Large-Scale Simulation Studies in Image Pattern Recognition* IEEE Trans. PAMI, Vol. 19, No. 10, pp. 1067-1079, October 1997.
- [21] Henry S. Baird, T. Kanungo, R. M. Haralick. *Validation and Estimation of Document Degradation Models* Proc., 4th UNLV Symp. on Document Analysis and Information Retrieval, Las Vegas, Nevada, April 24-26, 1995.
- [22] D.-S. Lee, C. Nohl, Henry S. Baird, *Language Identification in Complex, Unoriented, and Degraded Document Images*. Proc., IAPR 1996 Workshop on Document Analysis Systems, Malvern, PA, October 14-16, 1996.



- [23] I. Kanungo and P. Resnik, *The Bible, Truth, and Multilingual OCR Evaluation*. Proc. of SPIE Conference on Document Recognition and Retrieval (VI), vol. 3651, San Jose, CA, January 27-28, 1999
- [24] Bansal, Veena and R.M.K. Sinha *Segmentation of touching and fused Devanagari characters* Pattern Recognition, volume 35 (2002), number 4 pp. 875-893.
- [25] Oivind Due Trier, Anil K Jain, Torfinn Taxt, *Feature Extraction Methods For Character Recognition A Survey (1995)* Pattern Recognition Vol 29, No. 4 pp. 641-662, 1996.
- [26] U. Garain, B.B. Chaudhuri, *Segmentation of Touching Characters in Printed Devnagari and Bangla Scripts Using Fuzzy Multifactorial Analysis*. Proceedings of the 6th International Conference on Document Analysis and Recognition. (ICDAR '01).