

Creation of data resources and design of an evaluation test bed for Devanagari script recognition

Srirangaraj Setlur, Suryaprakash Kompalli, Ramanaprasad Vemulapati, and Venugopal Govindaraju

Center of Excellence for Document Analysis and Recognition,

University at Buffalo, Buffalo NY 14260

{setlur, kompalli, raman-v, govind}@cedar.buffalo.edu

ABSTRACT

The Indian subcontinent has a large number of languages, dialects, and scripts and the *Devanagari* script is the primary and most widely used of all the scripts. Much of the *Devanagari* OCR research has been restricted to a handful of groups and has not 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 analyses and evaluation tools are currently lacking and creation of such data resources will undoubtedly provide a much needed fillip to researchers working on *Devanagari* OCR.

This paper describes the initial steps of a National Science Foundation sponsored effort under the International Digital Libraries program to create data resources to facilitate development of *Devanagari* OCR technology and to provide a standardized test bed and evaluation tools for *Devanagari* script recognition.

Keywords: *Devanagari* OCR Data Resources, Multi-lingual OCR, OCR Evaluation, Ground-truthing

1. INTRODUCTION

Devanagari is the script used by a number of Indian languages, prominent among which are Sanskrit, Hindi and Marathi. Many other Indian languages use close variants of this script which is an alphabetic script. Over half a billion people around the world read and write in Hindi. Sanskrit is one of the ancient languages of the world and even though it is no longer a spoken language, masterpieces of Sanskrit literature such as the *Mahabharata*, the *Ramayana* and the *Bhagavadgita* are very well known internationally and there is a great interest in digitizing and archiving different versions of these and other lesser known works of Sanskrit literature. As such, research on *Devanagari* OCR is attracting a lot of interest.

Automatic text recognition using OCR is the process of converting the image of a textual document into its textual equivalent, the advantage being that the textual material is editable. The document image itself can be either machine-printed or handwritten, or perhaps a mixture of the two. Generally, speaking the recognition of any unit of a document is in the context of the rest of the document and can benefit from the context. However, for all practical purposes, the units of recognition range from a character to a sentence. Grammatical and linguistic constraints from the sentence level can drive the recognition of the characters and words.

The development of OCRs requires vast amounts of training data with ground-truth to achieve acceptable levels of performance. Collection and truthing of data is expensive as well as labor intensive. Standardized test data is vital to compare the relative performance of different OCRs. Such data is sorely lacking for *Devanagari* researchers.

Section ?? provides a brief overview of the *Devanagari* script. Section ?? reviews some *Devanagari* OCR techniques that have been published. Section ?? describes the data collection and ground-truthing methodology and the tool designed for this purpose. Section ?? describes the future work planned under this initiative.

2. OVERVIEW OF DEVANAGARI SCRIPT

The basic character set of *Devanagari* script is very large comprising of 16 vowels (Table ??) and 36 basic consonants (Table ??). A vowel following a consonant takes a modified shape, which is

Vowel	Modifier	Modified Consonant
अ		क
आ	ा	का
इ	ि	कि
ई	ी	की
उ	ु	कु
ऊ	ू	कू
ऋ	ृ	कृ
ॠ	ॠ	कॠ
ऌ	ॡ	कॡ
ॢ	ॣ	कॣ
अं	ं	कं
अः	ः	कः

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

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 ??). More than one consonant can combine with each other to form conjunct consonants (Some examples can be seen in Table ??). In addition to these, there are also a few special characters used in *Devanagari* text which are shown in Table ?? and the 10 numerals shown in Table ??). 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 characters 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.

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. Due to the complexity of the *Devanagari* character set, the text is pre-processed in different ways by different researchers. *Devanagari* text can have compound characters, characters with vowel modifiers and half characters. Researchers deal with such characters in different ways. While some attempt

Simple Consonants										
क	ख	ग	घ	ङ	च	छ	ज	झ	ञ	
ka	kha	ga	gha	ṅa	cha	Cha	ja	jha	JNa	
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 2. Simple Consonants (with transliteration), Digits and Special Characters in *Devanagari*

Consonants	Conjunct form
क त	क्त्
क य	क्य
क व	क्व
घ न	घ्न
ङ क	ङ्क
ङ ग	ङ्ग
ञ च	ञ्च
ट ट	ट्ट
त त	त्त
त न	त्न
त र	त्र
द ग	द्ग

Consonants	Conjunct form
द द	द्द
द ध	द्ध
द न	द्न
द म	द्म
द य	द्य
द व	द्व
ध न	ध्न
न न	न्न
प त	प्त
प न	प्न
प ल	प्ल
म ल	म्ल

Consonants	Conjunct form
ल ल	ल्ल
श च	श्च
श न	श्न
श र	श्च
श ल	श्ल
श व	श्व
ष ट	ष्ट
ह न	ह्न
ह य	ह्य
ह र	ह्र
ह ल	ह्ल
ह व	ह्व

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

to separate the vowel modifiers and the core characters, recognizing them separately, others attempt to recognize each compound consonant or consonant with vowel modifier as a separate class. Thus, after extracting words from a test document, different segmentation methods are used to divide the word into its constituent characters or shapes. It is essential that the data collection be driven by the needs of a majority of the developers.

A study of the many available standard OCR databases for other scripts reveal that for Latin-character based scripts, most databases provide samples of all characters in the script under consideration. However, due to the presence of compound characters and modifiers in Devanagari, providing sufficient samples of every possible character shape in the data set would be infeasible. Since researchers in *Devanagari* OCR appear to not have 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 and others treating the base character with the modifiers or conjunct characters as distinct classes, 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 later, as necessary.

4. DATA COLLECTION AND GROUND TRUTHING 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 truthing 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.

Simulated data is a poor alternative to real data since it is very hard to model all the variations in print, handwriting and quality of documents that are encountered. To this end, 400 printed document pages in Hindi and Sanskrit are being scanned. These documents include a variety of printed matter such as pages from old and new books, newspapers, magazines of varying print quality, font types and sizes. Estimating an average of 300 words per document, a data set of 120,000 words would be created.

We have developed a framework for truthing based upon our previous experience in designing ergonomic truthing tools, to handle the new challenges posed by *Devanagari* document truthing. 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 benchmarked and assessed.

Unicode has emerged as the de-facto standard for encoding characters of all languages and hence would be the ideal standard to represent both the ground-truth as well as the results of recognition systems. 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 and natural interface.

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. The ITRANS transliteration scheme is a popular and fairly intuitive scheme that can be learnt quickly.

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.

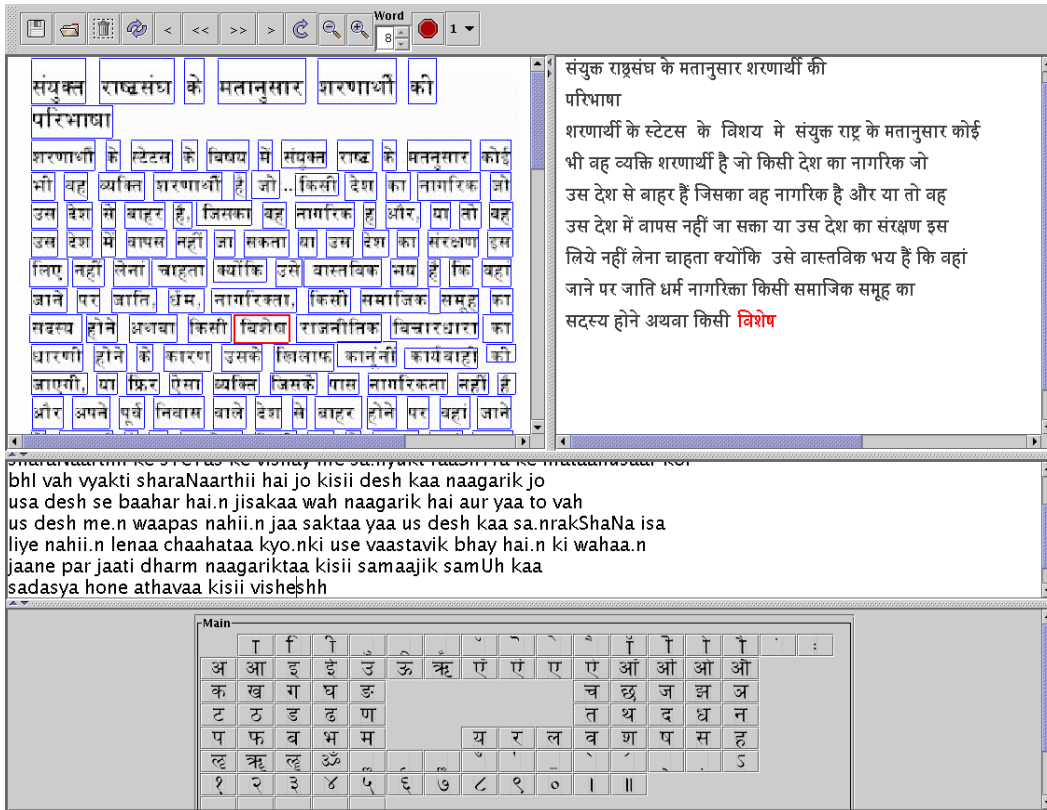


Figure 1. Truthing tool for *Devanagari* documents

Figure ?? shows a screen dump of the truthing tool interface.

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.

Representing the truth in a standard format such as Unicode allows results and truth to be easily compared. We also propose to provide query tools to retrieve word images based on a search string in transliterated text as well as Unicode values. The database would initially contain word images tagged with the corresponding word truth in Unicode format and also the transliterated text for that word in ITRANS format. We also propose to make available the truth files in XML format to enable users to design customized parsers and other tools for analyzing the data. Unicode uses a hexadecimal code for representation of the basic character, and the modifiers. Detailed rules exist to generate the full character set [14]. Thus, it can (theoretically) represent the entire set of *Devanagari* characters.

5. CONCLUSIONS AND FUTURE WORK

The initial datasets that will be made freely available to the public by early 2003 and distributed via the world wide web will consist of word images with their ground-truth. 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. Haralick Robert M *UW-I English Document Image Database*. Intelligent Systems Laboratory, University of Washington.
2. Haralick Robert M *UW-II English/Japanese Document Image Database*. Intelligent Systems Laboratory, University of Washington.
3. Phillips Ihsin *UW-III English/Technical Document Image Database*. Department of Computer Science, Queens College, The City University of New York.
4. Schlosser Steven G. *ERIM Arabic Database* Document Processing Research Program. Information and Materials Applications Laboratory, Environmental Research Institute of Michigan.
5. Senior Andrew *Ground Truthed Handwritten Word Images* Handwriting Recognition Group, IBM T.J.Watson Research Center
6. Alpaydin E., Alimoglu F. *Pen-Based Recognition of Handwritten Digits*. Bogazici University, Department of Computer Engineering
7. 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

8. Malaviya, Ashutosh. Leja, Christoph. Peters, Liliane. *Multi-script handwriting recognition with FOHDEL*. New Frontiers in Fuzzy Logic and Soft Computing Biennial Conference of the North American Fuzzy Information Processing Society - NAFIPS 1996. IEEE, Piscataway, NJ, USA.. p 147-151
9. Sethi, I.K., Chatterjee, B., *Machine Recognition of Constrained Hand Printed Devanagari*. Pattern Recognition(9), 1977, pp. 69-75
10. Connell, S.D., Sinha, R.M.K., Jain, A.K. "Recognition of Unconstrained On-line *Devanagari* Characters" ICPR00(Vol II: 368-371).
11. *The Unicode Standard Version 3.0*
12. S. Setlur and V. Govindaraju, "Translingual OCR by Template Correlations", *Proceedings of SPIE*, Vol. 3964, 1999, pp. 305-312.
13. 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.
14. 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.
15. V. Bansal, "Integrating Knowledge Sources in Devanagari Text Recognition", *Ph.D. Thesis*, 1999.
16. B. Chaudhuri and U. Pal, "Skew angle detection of digitized Indian script documents", *IEEE transactions on Pattern Analysis and Machine Intelligence*, Vol. 19, No. 2, 1997.
17. B. Chaudhuri and U. Pal, "A complete printed Bangala OCR system", *Journal of Pattern Recognition*, Vol. 31, No. 5, pp. 531-549, 1998.