Enabling digital access to multi-lingual Indian documents

Venu Govindaraju, Srirangaraj Setlur, Swapnil Khedekar, Suryaprakash Kompalli,
Faisal Farooq, Ramanaprasad Vemulapati
govind@cedar.buffalo.edu
Center of Excellence for Document Analysis and Recognition,
University at Buffalo, NY - 14260

Abstract

In this paper we present methodologies for three important tasks that will eventually enable digital access of multilingual Indian document images. First, we describe several document image analysis techniques necessary to prepare Devanagari document images for OCR. The second task is OCR for machine printed Devanagari words without the help of a lexicon. We describe the OCR methodology and show how it is being extended to other Indian languages. Finally, we describe a versatile platform that facilitates automatic segmentation of document images in multiple Indian languages and an interface to capture the ground truth corresponding to the text. We use transliterated English text and virtual keyboards in a range of Indian languages for this purpose. The multi-lingual data entry capabilities of the tool and its underlying UNICODE data representation within a structured XML document also allow users to annotate passages of text in one language in other languages using a markup scheme to switch between scripts. Text and annotations are rendered in the appropriate scripts as the text is being annotated, thus providing users prompt and natural feedback. The XML back-end allows meta-data to be recorded describing the annotated document.

1. Introduction

Online collections of digitized texts have facilitated world-wide access to materials traditionally inaccessible for a variety of reasons: the expense of print publication, the fragility of the documents, their remoteness, or the reticence of their proprietors. Efforts to promote access to two millenia of Indian written culture have been frustrated by the limitations of the technology. While manuscript transcription and indexing can pro-

vide ready access to these documents, the work is arduous and labor-intensive. To date, an effective tool for recognition of Indic scripts, such as Devanagari, has not been developed. Our goal is to make significant advances in automatic recognition of Devangari and other Indic scripts making it possible to access Indian documents with ease. Figure 1 illustrates the features of the envisoned system. Query terms are presented by the user in either the original script or transliterated English text. The response to the search is enabled by document image analysis and OCR techniques designed for Indian languages.

The Indian sub-continent lays claim to literature spanning 11 official scripts (henceforth refered to as Indic scripts) that support 18 official languages. These can be broadly grouped into (i) Aryan scripts - Devanagari (supporting Sanskrit, Hindi, Nepali, Marathi, and Konkani), Bengali (Bengali, Manipuri, and Assamese), Gujarati (Gujarati), Punjabi (Punjabi) and Oriya (Oriya), (ii) Dravidian scripts - Tamil (Tamil), Telugu (Telugu), Kannada (Kannada), and Malayalam (Malayalam), and (iii) Arabic-based scripts - Urdu (Urdu, Kashmiri, and Sindhi). Indian languages share a phonetically common character set and other similarities. Figures 2 and 3 show the characters and their symbolic representation in various Indic scripts. All Indic scripts except Urdu are written from left to right.

The basic Devanagari language character set consists of about 16 vowels and about 36 basic consonants. Unlike English, text in Indic languages is not simply a juxtaposition of character symbols. Vowel modifiers combine with one or more consonants to form conjuncts that are formed by modifying the constituent consonant shapes. Vowel modifiers usually attach to the consonants they modify on the top, bottom or right. Conjunct consonants are formed by one of the consonants assuming a half-consonant or a special form placed either to the left, top or bottom of the other

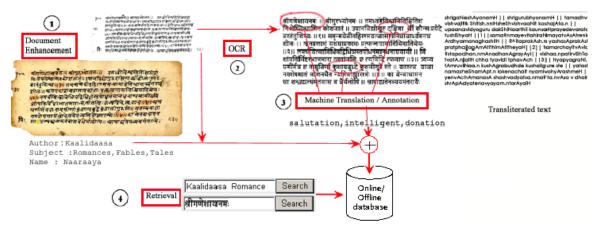


Figure 1. Multi-lingual Digital Library: (1),(2): Script/Language Identification (3): Meta-data - Document info, Translation, Annotation (4): Search and retrieval

consonant. The conjuncts can be either touching or separated from each other and more than two consonants can also combine, albeit rarely. Hence, despite a small basic character set and the absence of upper and lower cases as in English, the number of classes that an Indic OCR has to consider is very large. The units of recognition tend to be characters or words. Grammatical and linguistic constraints from the sentence level can help post-process the recognition of characters and words to eliminate ambiguities and correct errors. Character and word training data with ground truth from documents that are representative of those encountered in the real world is essential for the development of the OCR systems. The need for tools that can facilitate the collection of ground-truthed data for Indic OCR is apparent.

In this paper we present methodologies for three important tasks that will eventually enable digital access of multilingual Indic document images.

- 1. Document Image Analysis. It includes binarization, skew correction, text block extraction, line separation, and word segmentation. We show methods developed for Devanagari documents (Section 2).
- 2. OCR for machine-printed Devanagari text. Automatic text recognition using OCR is the process of converting a textual document into its digital equivalent to enable editing, efficient storage, and querying (Section 3).
- 3. Design and development of a versatile tool for ground-truthing and annotating multi-lingual texts. It includes image processing modules that can generate word and character data as well as a semi-automated interface (using the OCR) to collect character data for training. The modular design of the tool naturally lends itself to other Indic scripts as well (Section 4).

2. Devanagari Documents' Image Analysis

The process of converting a physical document into a corresponding entity in a digital library involves document layout analysis, separation of text and graphics, OCR of text blocks, annotations and other enhancements.

Devanagari documents are scanned at 300dpi resolution with 256 grayscale levels from different sources (books, magazines, newspapers, etc.). Since most segmentation methods use projection profiles which work best on binary images, binarization is the first step of processing. It is done using a neighborhood adaptive thresholding technique. We use morphological opening to remove spurious black pixels which may interfere with the segmentation techniques and join any broken strokes. Skew correction is performed if needed [2].

2.1. Separation of Text and Graphics

This is accomplished by smearing combined with horizontal projections. We take advantage of the headline or *Shirorekha* which is a distinctive feature of Devanagari Script. In the horizontal histogram, rows where the *Shirorekha* is present corresponds to the bin with the maximum number of black pixels. This shows up as a sharp peak in the histogram plot and gives the profile generated from text a distinctive pattern (Figure 2). Another important feature of Devanagari text is its consistent frequency and orientation. Adjacent lines and words in a line are of similar sizes respectively.

Thus algorithmically, Devanagari text pattern is distinguished from graphics by (i) the shape of the horizontal profile, and (ii) the regularity of intervals between the peaks in the profile and their widths. Figure 2 illustrates this idea. The method is recursively applied [8] to obtain the final output (Figure 5). Figure ?? gives the psuedocode of the implemented algorithm.

	का 8	ৰ ক ু	ಕ್ ಕ ಈ ಹಿ	<u>5</u>	
	ৰা খে	য খ ও	ఖ្ಖ_ ഖ	كمو	
	ग्र		តែក_ ហ	گ	
	घध				
			*	گو	
	ङ ১	হ ওড়ি	ឧស្ធន	_	
चिथ्च চ 5	ವಚ್ ಈ	<u>2</u> æ			
छঔਛছ⊋	ಭ ಛ _	چمو 20	प्प्यिপ्ट	ಏಪ⊔	پ ل
जिश्न छ ्	ස ස වූ	ன் <u>உ</u>	फाध्डक 🖰	ಘ ಫ ⊥	پھ ∆ہ
झअङबा छ	ဏ္ ညာ	جو ران	ৰাওঘিব 🛭	ಬ ಬ _	ബ 💛
যাগছ আই	మి. ఉు. @⊇ యా యి =	ଗ୍ରେ	দি এভড় ও	ಭ ಟ _	بمو 6
21216 00	× 1 9		मिभागघरा	*	מ 🗅
टिट्टिउंड	ಬಟ∟	క త			_
ರದರಾರ	ಠ ಠ _	ٹو ٥	मियजग्री	ထေတ္မ	ی رو
<u>इ</u> ১३ডৢঞ	ದ ವ _	ഡ 🚡	र्थिवत्र	ದರ ೯	თ კ
೭೭೯ ೯೪	್ ಫ _	ڏنو سک	लिपहन्	ဗဗ၈	ല ၂
णिश्टिष	ස ස හො		ंविव्ह	ವವ ಮ	വ
4 1 01					
तित्डछुछ	ಕತ ಕ	ග ප	হাখন শ গ	ช ଅ _ '	ش (۵)
थथ्यथ् 2		10	ष्प _ य ६		ഷി —
			संसग्य ६		സ സ
दिधिच्य	ದವ_	B 3	l .	(0 -	
ម ឧប ម ប	女 돠 _	دم و (10	<u>ইএটিক চ</u>	హ్ జ్ ஹూ	പറ ക
ਜਿਮਨਕਨ	ನನ ೨	თ ა			

Figure 2. Consonant set- Indic scripts share a common phonetic character set with minor differences.

The title is identified by its large font size and location in the upper half of the document. The horizontal and vertical gaps between profiles are used to separate paragraphs in a column. The reading order of text is obtained by sorting text blocks by their left boundaries as major key and top boundary as minor key (Figure 5). This simple approach works most of the time except when graphics blocks are larger than the column size. In such cases, additional page layout knowledge is required.

2.2. Line separation

Researchers [1, 4] have used horizontal projection profiles for line separation as well. However, these techniques assume that there is a consistent white gap between consecutive lines. However, presence of ascenders and descenders violate this assumption frequently (Figure 6).

We reuse the horizontal profile from the text block segmentation stage and take advantage of the constraint that each prominent sharp peak corresponds to a single text line. First we find the connected components (words in Devanagari text) in a region and sort by their left boundaries. Each component in the list is linked to components with horizontal overlap (threshold is 60% of component height). Figure 8(a) illustrates the technique. There are a few drawbacks to this approach. Small components such as dashes, Anuswar, and noise can create problems (Figure 8(b)). Furthermore, if words from two different lines are physically touching each other creating a single connected component the results will be incorrect. To address this issue, we consider the median of components which are highly probable words. Components with height less than 0.5 x median or greater than 1.5 x median are not included in the grouping. The deatils of the algorithm are shown in psuedocode



Figure 3. Set of vowels and vowel modifier illustration for Indic scripts from Devanagari (Dev) to Urdu (Arabic).



Figure 2. Sample Image and histogram showing the widths.

form in Figure 7.

2.3. Word separation

Devanagari words are single connected components with a horizontal line called headline or *Shirorekha* touching all characters. Each word is separated from adjacent word by a vertical gap. Thus simple vertical projection profile gives the separated words (Figure 9). To avoid intereference from ascenders and descenders, components touching either upper or lower boundary of the line separating adjacent lines (derived by method described above) and with less than 15% word height are removed before computing the profile. The deatils of the algorithm are shown in psuedocode

form in Figure 10.

2.4. Character separation

Character shapes are modified by vowels (Figure 3) now referred to as ascenders and descenders. Since all characters in a word are connected by the *Shirorekha* the obvious solution is to remove the *Shirorekha*.

2.4.1. Removal of Shirorekha Shirorekha exhibits the following properties in the horizontal projection profile: (i) it corresponds to the strongest peak except in cases of single characters such as \(\frac{1}{2}\), (ii) it corresponds to the widest spread owing to its thickness, and (ii) fewest number of

Take horizontal projection: profile Find regions between two valleys: width for each width

Determine if it has distinctive shape Calculate width of base and peak

end

Cluster adjacent widths showing regular pattern if regular pattern not found

Take vertical histogram and find gaps if no gap found

Label current block as IMAGE

end

end

for each regular pattern

Take vertical histogram and find gaps Segment blocks and label as TEXT

end

for each irregular pattern

Take vertical histogram and find gaps Segment blocks and call routine for each

end

Sort the TEXT blocks

Figure 4. Pseudocode for Separating and Text and Graphics

black-white transitions. Let, profile(i) denote the number of black pixels in the row i; run.no(i) denote the number of black to white transitions in row i; and extent(i) denote the distance between the leftmost and the rightmost black pixel in row i.

- **2.4.2.** Ascender Separation Ascenders are present above the *Shirorekha* by definition. Thus finding the connected components in rows 0 (first row of word) to i_1 isolates the ascenders (Figure 11). This approach is superior to using vertical profiles because the latter cannot separate combined ascenders such as in \mathbf{H} .
- 2.4.3. Character Body Separation Gaps in the the vertical profile of a word (with *Shirorekha* removed) correspond to inter character spacing (Figure 12). Characters such as **A** and **A** have significant features in the *Shirorekha* region that distinguishes it from **A** and **A**. Thus the *Shirorekha* is retained in the final character images after the character boundaries are identified. Characters such as **A**, **A** and **A** are split into two components and must be recognized separately during OCR.
- **2.4.4. Descender Separation** Descender separation is non-trivial because the character heights vary considerably (Figure 13). For example characters \overline{c} , \overline{n} and

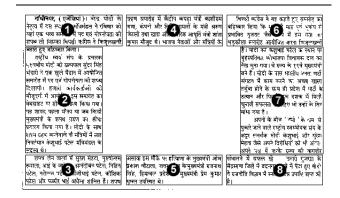


Figure 5. Text blocks in reading order.

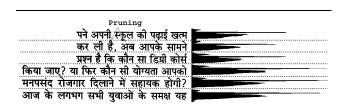


Figure 6. Closely spaced lines with ascenders and descenders (vowel modifiers) result in line segmentation errors.

\(\varepsilon\) all have different heights. Accordingly, the descender connects at different heights for different characters. Many researchers [1, 4] have ignored this factor altogether.

The proportion of characters with descenders is small. So, only a few characters would have a height greater than normal for a considered font. Heights of all characters are measured and quantized into 10 bins. The most populated bin height is taken as the average character height, say h. A character connected to a descender must have a significant portion of strokes below this baseline h. Empirically, we have found this ratio to be 30% of the character height. Figure 13(a) illustrates the portion of \mathbf{E} and below \mathbf{E} baseline is very small and thus not segmented but \mathbf{E} and \mathbf{E} are segmented.

Thus the presence of a descender is hypothesized if $h_c - h \ge 0.3 \times h$. h_l and h_h are the width of the bins adjacent to h. The region between h_l and h_h is searched for the descender (Figure 13b). Pseudocode of the algorithm is given in Figure ??.

2.4.5. Shadow Characters Individual characters which overlap (vertically) are said to be in the shadow of one another. These cannot be separated by the method described.

```
Take horizontal projection: profile
Find rows between two zeros in profile: width
for each width
      Find number of peaks in width
      if single
            Label width as LINE.
      else
            region_list=null
            Find Connected components: CC
            for each CC
                 Find region with hori overlap of >60%
                 if found
                       Grow region to fully include CC
                 else
                       Create new region=CC, add to list
                 end
            end
            Label each region as LINE
      end
end
```

Figure 7. Pseudocode for Line Separation.

Moreover, descenders might create the shadow effect as well (see Figure 14). We deal with this case as in [5]. Descenders are not separate but rather recognized as whole along with the base character.

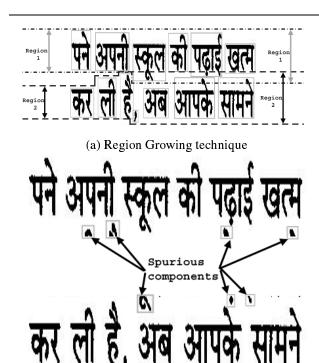
3. Devanagari OCR

The image analysis routines described in the previous section return 4 types of components: ascenders, descenders, base characters, and conjuncts. These will now be the input to the OCR system (Figure 15).

The reassembly procedure is applied at the word level. Once all the characters, ascenders and descenders are recognized and their relative positions are known, the relationship between the modifiers and base characters must be established. For example, the descender can attach to single character and thus the character most probable must be determined by choosing the candidate with maximum vertical overlap. In case of c, if the bottom touches a character it remains but if it connects a Γ , it becomes Γ . Also characters like Γ must be completed as Π or form a conjunct.

3.1. Feature Selection

Since the conjuncts are not split, the number of shape classes is large. Also, the character shapes are complex as in $\overline{4}$ and $\overline{4}$, and the features must capture nuances for accurate classification.



(b) Resulting segmentation and spurious components

Figure 8. Region Growing technique and resulting lines



Figure 9. Word segmentation using vertical profile.

We use gradient features as described in [9, 10] extracted from the original grayscale images. Given the wide variations in sizes of characters even within a font family, as in sizes of modifiers compared to base characters, all images are normalized to 32×32. The gradient computed using a Sobel operator measures the magnitude and direction of intensity changes in a small neighborhood of each pixel. Example of the gradient magnitude and direction of a character image is illustrated in Figure 16. Higher gradient magnitudes are indicated by darker shades of gray. The gradient map is computed by retaining only those gradient magnitudes which are above a threshold. The gradient map is put into a 3 X 3 grid and the gradient directions are quantized into 8 bins: horizontal, vertical, and the two diagonals in each direction. Thus each cell contributes 8 features into

```
for each row
      Initialize profile=0, extent=1, runno=1
end
for each row i
      for each column j
           if image(i,j)=black
                  Increament profile(i)
           end
           if image(i,j)=white and image(i,j-1)=black
                  Increament runno(i)
           end
      end
      Find leftmost black pixel start
           and rightmost black pixel stop
      extent(i)=stop-start+1
      f(i) = profile(i) * extent(i) / runno(i)
end
Find i such that max(f)=f(i)
for x=i to 1
      if f(x) < 0.7 * f(i)
           i_1 = x + 1
           break
      end
end
for x=i to word_height
      if f(x) < 0.7 * f(i)
           i_2=x-1
           break
      end
end
```

Figure 10. Pseudocode for Word Separation

the feature vector. Concatenation of the gradient vectors for all 9 cells forms the 72 features in the vector.

In Figure 17(a), we show the feature vector generated for $\overline{\mathbf{q}}$ and $\overline{\mathbf{q}}$. Since these two characters shapes are the same except for cell 5, it is reflected in the plots of the figure. The low standard deviation values over the entire set (approx. 0.012/sample) of same character with different fonts and sizes, indicate the robustness of these features (Figure 17(b)).

3.2. Training

The training dataset is formed from 10 different books, magazines and newspapers. Documents have different languages, font styles and sizes. Also they differ in page quality and layout. We extract character images from these documents using the tool [11] described in section 4. 40 images

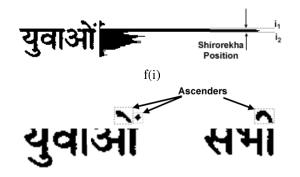


Figure 11. Shirorekha detection and Ascender separation

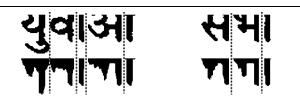


Figure 12. Character separation using vertical profile

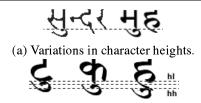
for each character are taken from each source giving us a total of 400 images per class.

3.3. Classifier Design

We use same features and extraction method for the classes of ascenders, descenders and characters. However, the classifier for each set is different.

3.3.2. Characters We consider 38 characters and 83 frequently occurring conjunct character classes in a multistage classification approach. First we preclassify these classes into 4 categories depending on their structural properties such as position of the vertical line *danda* which is straightforward.

• Characters with No danda: ত ত হ হ হছ etc.



(b) Different height bins and probable segmentation region.



(c) Probable region and f(i) for ξ , \mathfrak{F} and ξ . Dotted line denotes $\min(f(i))$, segmentation boundary.

Figure 13. Shirorekha detection and Ascender separation

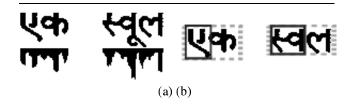


Figure 14. Shadow characters can not be separated by vertical profiles.

- Characters with Middle danda : ক फ ক্ट etc.
- Characters with Right danda : त अ व्य स्त etc.
- Characters with Multiple dandas : क्य फ्त etc.

Each category is then classified using a separate classifier. We use a three level artificial neural network for each of the category. Since we have a 72 features floating point vector we have 72 perceptrons in input layer. Depending on the number of output classes, number of nodes in the hidden and output layer varies. The network is trained by using standard back-propagation algorithm. We use upto top 3 choices which can be used by lexicon-based driven post-processing methods.

3.4. Preliminary results

The classification accuracy on characters (including conjuncts), ascenders, and descenders is about 98.5% on the training set and 95% on a blind test set. Extensive test results will be reported in the final paper. Failure analysis

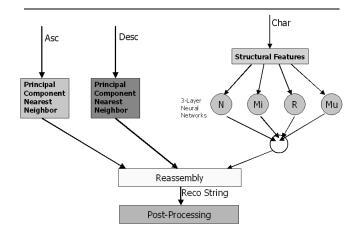


Figure 15. Classifier Design: Deals with classes of ascenders, descenders, base characters and conjuncts separately.

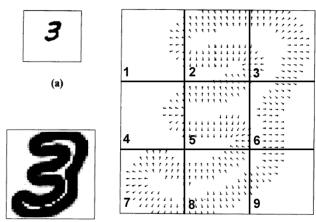


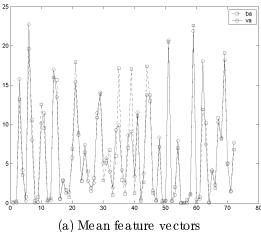
Figure taken from [9].
(a) Original Image. (b) Gradient magnitude.
(c) Gradient directions. 3 X 3 Grid and ordering.

Figure 16. Gradient magnitude and direction.

indicates that confusions are among ব, ৰ and इ, ড, ह (Figure 18).

4. Multi-lingual truthing, annotation and OCR analysis tool

We have developed a versatile multi-function tool for multi-lingual OCR researchers that provides an interface for generating ground-truth for a document image at multiple levels of granularity, annotate or record translations for documents in multiple languages as well as provides a graphical analysis tool for trouble-shooting and evaluating OCR



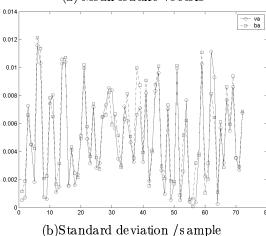


Figure 17. Graphs showing values for ব and ব



Figure 18. Character confusions

problems.

4.1. Ground truthing

Ground-truthing document images in a non-Roman language necessitates the adoption of a non-ASCII representation for the truth. 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. While the representation of the truth would be in Unicode the pro-

cess 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 one 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 and other Indic script characters cannot be represented in phonetic English, approximations such as RRi for $\frac{1}{2}$ are used. Also, many sounds can be phonetically represented in English in more than one way. This has led to the adoption of many different transliteration schemes and no single standard has found universal acceptance. After analyzing these transliteration schemes, ITRANS, a purely ASCII and intuitive transliteration scheme for Indic languages was chosen as the transliteration standard for the tool.

Using transliteration as a truthing aid enables users to record the truth using a phonetic representation of the Indic text and the transliterated text can be parsed and converted into the corresponding Unicode for storage.

4.2. Annotation

The multi-lingual data entry capabilities of the tool and its underlying UNICODE data representation within a structured XML document also allow users to annotate passages of text in one language in other languages using a markup scheme to switch between scripts. Text and annotations are rendered in the appropriate scripts as the text is being annotated, thus providing users prompt and natural feedback. The annotation features of the tool will be used by the LitGloss project [?] at the University at Buffalo, which is a collection of texts of literary or cultural interest, written in languages other than English, and expertly annotated so as to facilitate comprehension by English-speaking readers.

4.3. OCR analysis

5. Summary

References

[1] B.B. Chaudhari, U. Pal, An OCR system to read two Indian language scripts: Bangla and Devanagari (Hindi), Proc. 4th

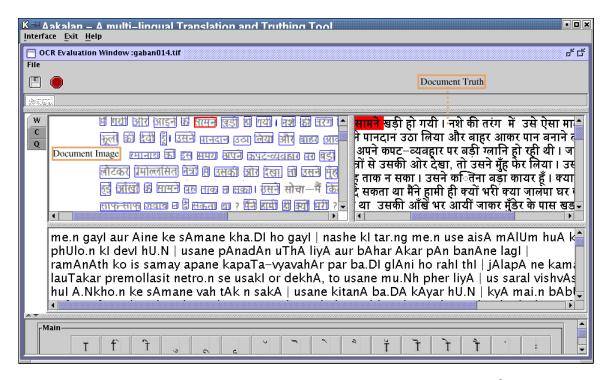


Figure 19. Truthing of document images: Page layout analysis provides word boundaries.

Transliteration and GUI keyboard used for text input.

- International Conference on Document Analysis and Recognition, Vol 2, pp 1011 -1015
- [2] B.B. Chaudhari, U. Pal, Skew Angle Detection of Digitized Indian Script Documents, IEEE Transactions on Pattern Analysis and Machine Intelligence, vol 19, No 2, Feb 1997 pp. 182-186
- [3] B.B. Chaudhari, U. Pal, A Complete Printed Bangla OCR System, IEEE Transactions on Pattern Analysis and Machine Intelligence, vol 31, No 5, 1998, pp. 531-549
- [4] V. Bansal, R.M.K. Sinha, Integrating knowledge sources in Devanagari text recognition system, IEEE Transactions on Systems, Man and Cybernetics, Vol. 30, No. 4, pp. 500-505, 2000.
- [5] V. Bansal, R.M.K. Sinha. On How to Describe Shapes of Devanagari Characters and Use Them for Recognition, Proc. Fifth International Conference on Document Analysis and Recognition.
- [6] V. Bansal, R.M.K. Sinha. Segmentation of Touching and Fused Devanagari Characters, Pattern Recognition, vol 35, 2002, pp. 875-893.
- [7] A. Bishnu, B.B. Chaudhuri. Segmentation of Bangla Hand-written Text into Characters by Recursive Contour Following, IEEE Transactions on Systems, Man, & Cybernetics Part A: Systems & Humans. Vol 30, No 4, Jul 2000, pp. 500-505.
- [8] S. Khedekar, V. Ramanaprasad, S. Setlur, V. Govindaraju*Text Image Separation in Devanagari Documents*, Proc. International Conference on Document Analysis and Recognition, 2003, pp.1265-1269.

- [9] G. Srikantan, S. Lam, S. Srihari Gradient based Contour Encoding for Character Recognition, Pattern Recognition, Vol. 29, No. 7, 1996, pp. 1147-1160.
- [10] S. SinghShape Detection using Gradient Features for Handwritten Character Recognition, Proc. 13th International Conference on Pattern Recognition, (ICPR'96), Austria, Vol. III, 1996, pp. 145-149.
- [11] S. Kompalli, S. Khedekar, F. Farooq, V. Ramanaprasad, S. Setlur, V. GovindarajuMultilingual Script Truthing, Annotation and Recognition Platform for Digital Libraries, International Workshop on Document Image Analysis for Libraries, DIAL'04, (submitted).
- [12] Sethi, I.K., Chatterjee, B., Machine Recognition of Constrained Hand Printed Devanagari, Pattern Recognition(9), 1977, pp. 69-75.
- [13] R.M.K. Sinha, Mahabala H. Recognition of Devanagari Script, IEEE Trans on Systems, Man, and Cybernetics, Vol. 9, pp 435-441, 1979.
- [14] P.O. Duda, P. E. Hart and D. G. Stork Pattern Classification Wiley, New York, 2001.
- [15] The Unicode Standard Version 4.0. http://www.unicode.org/, http://www.unicode.org/onlinedat/languages-scripts.html Addison-Wesley, 2000. ISBN 0-201-61633-5.

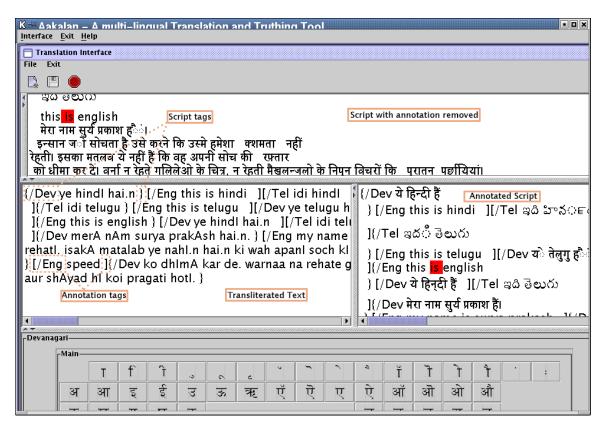


Figure 20. Multi-lingual annotation: Devanagari, Arabic, English and Telugu data entry using language codes and annotation tags.

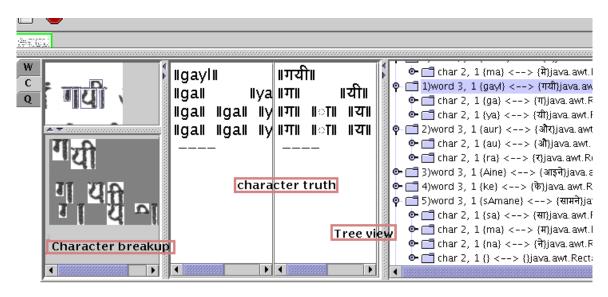


Figure 21. Interface to study the breakup of a word: correct breakup